

1 **Short title:** Arabidopsis Fe deficiency response pCREs

2

3 **Corresponding authors:**

4 Petra Bauer

5 Heinrich Heine University

6 Institute of Botany

7 Universitätsstraße 1, Building 26.13, Floor/Room 02.36

8 40225 Düsseldorf

9 Tel.: +49-211-81-13479

10 E-mail: petra.bauer@hhu.de

11

12 Shin-Han Shiu

13 Michigan State University

14 Plant Biology Laboratories

15 612 Wilson Road, Room 166

16 East Lansing, MI 48824-1312

17 Tel.: +1-517-353-7196

18 E-mail: shius@msu.edu

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38 **Putative *cis*-regulatory elements predict iron deficiency responses in Arabidopsis**
39 **roots**

40
41 Birte Schwarz¹, Christina B. Azodi², Shin-Han Shiu^{2,3}, Petra Bauer^{1,4}

42
43 ¹Institute of Botany, Heinrich Heine University, Universitätsstr. 1, Düsseldorf, Germany

44 ²Department of Plant Biology, Michigan State University, East Lansing, MI, USA

45 ³Department of Computational, Mathematics, Science, and Engineering, Michigan State
46 University, East Lansing, MI, USA

47 ⁴Cluster of Excellence on Plant Science (CEPLAS), Heinrich Heine University, Düsseldorf,
48 Germany

49
50 **One sentence summary**

51 >100 putative *cis*-regulatory elements robustly predict Arabidopsis root Fe deficiency-
52 responses in computational models, and shed light on the mechanisms of transcriptional
53 regulation.

54
55 **Funding:** This work was supported by the German Research Foundation grant through the
56 DFG International Research Training group 1525 to P.B., a NSF Graduate Research
57 Fellowship to C.B.A, and by grants to S.-H.S. from the US National Science Foundation IOS-
58 1546617, DEB-1655386, and DGE-1828149, and the US Department of Energy (DOE) Great
59 Lakes Bioenergy Research Center (DOE Office of Science BER DE-SC0018409). This work
60 received funding from Germany's Excellence Strategy, EXC 2048/1, Project ID: 390686111.

61
62 **Author contributions:** B.S., P.B., and S.-H.S. conceived the project; B.S., and S.-H.S.
63 designed the research plan; B.S., and C.B.A. analyzed the data; B.S. wrote the original draft;
64 B.S., C.B.A., S.-H.S., and P.B. reviewed and edited the article; P.B., C.B.A, and S.-H.S.
65 acquired funding. P.B. and S.-H.S. agree to serve as the authors responsible for contact and
66 ensure communication.

67
68 **Corresponding authors:** Petra Bauer, petra.bauer@hhu.de; Shin-Han Shiu, shius@msu.edu

69
70 **Keywords:** Arabidopsis, Fe deficiency, machine learning, Random Forest, co-expression
71 clustering, *cis*-regulatory element, transcription factor binding motif, FIT, IDE1, coumarin

72
73
74
75

76 **Abbreviations**

77	CIS-BP	Catalog of inferred sequence binding preferences
78	CNS	Conserved non-coding sequence
79	DAP-seq	DNA affinity purification sequencing
80	Fe / -Fe	Iron / Iron deficiency
81	FeS	Iron-Sulfur
82	FET	Fisher's exact test
83	FIT	FE DEFICIENCY-INDUCED TRANSCRIPTION FACTOR
84	freq-pCRE	Frequent pCRE
85	GO	(Biological process) Gene ontology
86	GS	Gold standard
87	IDE	Iron Deficiency-responsive Element
88	log ₂ FC	log ₂ fold-change
89	MA	Mugineic acid
90	min-pCRE	Minimum set pCRE
91	PCC	Pearson's correlation coefficient
92	pCRE	Putative <i>cis</i> -regulatory element
93	PWM	Position weight matrix
94	RF	Random Forest
95	TF	Transcription factor
96	TFBM	Transcription factor binding motif
97	TSS	Transcription start site
98	TTS	Transcription termination site
99	Zn	Zinc

100

101

102

103

104

105

106

107

108

109

110

111 **Abstract**

112 Iron (Fe) is a key cofactor in many cellular redox processes, including respiration and
113 photosynthesis. Plant Fe deficiency (-Fe) activates a complex regulatory network which
114 coordinates root Fe uptake and distribution to sink tissues, while avoiding over-accumulation
115 of Fe and other metals to toxic levels. In *Arabidopsis* (*Arabidopsis thaliana*), FIT (FER-LIKE
116 FE DEFICIENCY-INDUCED TRANSCRIPTION FACTOR), a bHLH transcription factor (TF), is
117 required for up-regulation of root Fe acquisition genes. However, other root and shoot -Fe-
118 induced genes involved in Fe allocation and signaling are FIT-independent. The *cis*-regulatory
119 code, i.e. the *cis*-regulatory elements (CREs) and their combinations that regulate plant -Fe-
120 responses, remains largely elusive. Using *Arabidopsis* genome and transcriptome data, we
121 identified over 100 putative CREs (pCREs) that were predictive of -Fe-induced up-regulation
122 of genes in root tissue. We used large-scale *in vitro* TF binding data, association with FIT-
123 dependent or FIT-independent co-expression clusters, positional bias, and evolutionary
124 conservation to assess pCRE properties and possible functions. In addition to bHLH and MYB
125 TFs, also B3, NAC, bZIP, and TCP TFs might be important regulators for -Fe responses. Our
126 approach uncovered IDE1 (Iron Deficiency-responsive Element 1), a -Fe response CRE in
127 grass species, to be conserved in regulating genes for biosynthesis of Fe-chelating
128 compounds also in *Arabidopsis*. Our findings provide a comprehensive source of *cis*-regulatory
129 information for -Fe-responsive genes, that advances our mechanistic understanding and
130 informs future efforts in engineering plants with more efficient Fe uptake or transport systems.
131

132 **Introduction**

133 The micronutrient iron (Fe) is crucial for survival of all organisms. Plants encounter Fe
134 deficiency (-Fe) on calcareous and alkaline soils or during developmental phases with
135 increased sink demands. As a central component of heme and Fe-sulfur (FeS) clusters, Fe
136 acts in redox processes in plants in basically all important metabolic processes, such as the
137 respiratory and photosynthetic electron transport chains, chlorophyll biosynthesis, DNA
138 replication and repair, and nitrogen and sulfur assimilation. Consequently, plants react to -Fe
139 with a range of molecular, physiological and morphological adjustments, which is reflected in
140 transcriptional alterations of more than 1000 genes in *Arabidopsis* (*Arabidopsis thaliana*)
141 (Dinneny et al., 2008; Rodríguez-Celma et al., 2013; Mai et al., 2016). In the shoots, the
142 photosynthetic machinery is remodeled, leading to visible leaf chlorosis symptoms, and
143 essential Fe-requiring processes are prioritized, which can be achieved through break-down
144 of dispensable Fe-bound proteins and Fe redistribution between organelles (Blaby-Haas and
145 Merchant, 2013; Balk and Schaedler, 2014; Hantzis et al., 2018). In the roots, genes controlling
146 soil Fe uptake and detoxification of other transition metal ions acquired along with Fe are up-

147 regulated. Additionally, Fe is mobilized from internal storages and distributed to Fe sinks. -Fe
148 also leads to changes in root architecture and root hair morphology (Brumbarova et al., 2015;
149 Curie and Mari, 2017; Jeong et al., 2017; Li and Lan, 2017).

150 To acquire soil Fe, grasses secrete mugineic acid (MA) family phytosiderophores and
151 import Fe³⁺-MA complexes into the root ("Strategy II"). In contrast, non-grass monocots and
152 dicots, such as Arabidopsis, acquire Fe via a reduction-based mechanism, in which soil Fe³⁺
153 is solubilized by lowering the local pH through proton extrusion, followed by reduction to plant-
154 accessible Fe²⁺ at the root epidermis and Fe²⁺ uptake ("Strategy I") (Marschner and Römheld,
155 1994). In Strategy I, secreted chelators (mainly phenylpropanoid-derived coumarins or
156 riboflavin derivatives) aid efficient Fe³⁺ solubilization and reduction (Fourcroy et al., 2014;
157 Schmid et al., 2014). Thus, Fe chelation is important during acquisition in both strategies.

158 Transcriptional control plays an important role in -Fe responses. A regulatory cascade
159 ultimately controls a set of -Fe response genes. In both, Strategy I and II, the current cascade
160 model involves related subgroups of basic helix-loop-helix (bHLH) transcription factors (TFs).
161 When rice and Arabidopsis plants experience -Fe, subgroup IVc bHLH proteins activate
162 subgroup Ib and IVb *BHLH* genes (Zhang et al., 2015; Liang et al., 2017). Downstream from
163 IVc bHLH TFs (ILR3/bHLH34/bHLH104/bHLH115 in Arabidopsis, PRI1 in rice), subgroup Ib
164 bHLH TFs (bHLH38/39/100/101 in Arabidopsis, IRO2 in rice) and subgroup IVb bHLH TFs
165 (PYE in Arabidopsis, IRO3 in rice) regulate responses further downstream (Ogo et al., 2007;
166 Yuan et al., 2008; Long et al., 2010; Zheng et al., 2010). In addition, IVc bHLH protein levels
167 are controlled by Fe-regulated E3 ligases (Selote et al., 2015; Zhang et al., 2017).

168 Despite these conserved regulatory and functional interactions of subgroup IVc, Ib, and
169 IVb bHLH TFs between grass and non-grass species, it remains unclear if other regulatory
170 components between Strategy I and II are conserved. For example, in grasses, IDEF1 (IRON
171 DEFICIENCY-RESPONSIVE ELEMENT BINDING FACTOR1, ABI3VP1 subfamily of B3 TF)
172 and IDEF2 (NAC TF) coordinate -Fe responses through binding to IDE1 (Iron Deficiency-
173 responsive Element 1) and IDE2 (Kobayashi et al., 2007; Ogo et al., 2008). IDE1 has been
174 connected to induction of genes involved in Strategy II MA biosynthesis and Fe-MA uptake
175 (Kobayashi et al., 2005; Ogo et al., 2007). However, while barley IDE1 can drive reporter gene
176 expression in tobacco in a -Fe-dependent manner and IDE1-like motifs are present in several
177 Arabidopsis -Fe response genes, a function for IDE1 has not been shown in Strategy I plants
178 (Kobayashi et al., 2003; Kobayashi et al., 2005; Kobayashi et al., 2007; Murgia et al., 2011).
179 Strategy I Fe acquisition requires the bHLH TF FIT (FER-LIKE IRON DEFICIENCY-INDUCED
180 TRANSCRIPTION FACTOR) that is absent in rice (Colangelo and Gueriot, 2004; Jakoby et
181 al., 2004), and is activated upon -Fe mainly through interaction with subgroup Ib TFs (Yuan et
182 al., 2008; Sivitz et al., 2012; Wang et al., 2013). FIT is essential for up-regulation of Fe³⁺

183 reduction, Fe²⁺ uptake, and chelator biosynthesis and export (Colangelo and Guerinot, 2004;
184 Jakoby et al., 2004; Sivitz et al., 2012; Schmid et al., 2014; Mai et al., 2016).

185 A co-expression network built with -Fe-responsive genes gives insight into the complex
186 -Fe regulatory system in Arabidopsis (Ivanov et al., 2012). Among co-expression clusters, one
187 contains root-specific and FIT-dependent genes involved in Fe acquisition, while another one
188 is composed of root- and shoot-expressed FIT-independent genes. In this work, we refer to
189 robust (i.e. consistently identified in different studies) FIT-dependent and FIT-independent
190 genes as the “gold standard” (GS) -Fe-induced genes. The concept to discriminate FIT-
191 dependent and FIT-independent co-expression clusters has proven very informative for
192 interpreting mutant phenotypes and to place novel regulators into the -Fe response cascade
193 (e.g. Zhang et al., 2015; Liang et al., 2017; Gratz et al., 2019). FIT-independent network genes
194 mostly act in sub-cellular and long-distance transport and distribution of Fe and in Fe signaling
195 and they include subgroup Ib *BHLH* genes and *PYE* (Ivanov et al., 2012). Only few upstream
196 regulators for FIT-independent gene expression have been identified yet, namely bHLH IVc
197 TFs, controlling Ib *BHLH* and *PYE*, and *PYE* controlling *NAS4*, *ZIF1* and *FRO3* of the same
198 co-expression regulon (Long et al., 2010).

199 For most -Fe-responsive genes, including reliable marker genes (Ivanov et al., 2012;
200 Mai et al., 2016), the specific *cis*-regulatory elements (CREs) which coordinate their expression
201 are unknown. Computational approaches uncover regulatory connections on a genome-wide
202 scale, such as through elucidating the *cis*-regulatory code, i.e. the collection of CREs and the
203 genes they regulate in a given regulatory context (Yáñez-Cuna et al., 2013). Putative CREs
204 (pCREs) could be identified computationally by the over-representation of sequences in the
205 promoter regions of co-regulated genes. Combining with data for TF binding motifs (TFBMs)
206 in Arabidopsis (Weirauch et al., 2014; O'Malley et al., 2016), regulatory connections can be
207 made between TFs, binding sequences, and target genes. To further improve the confidence
208 of computationally derived *cis*-regulatory code, machine learning algorithms (reviewed in Ma
209 et al., 2014) can be applied to build models with pCREs to predict gene expression or
210 transcriptional responses. These models can be evaluated by making predictions on
211 expression of genes that are not part of the model training. Most importantly, a good model
212 indicates that the pCREs used are most likely important for regulating the expression/response
213 of interest. Previous work has demonstrated the suitability of machine learning for elucidating
214 the *cis*-regulatory code of environmental stress responses in Arabidopsis (Zou et al., 2011;
215 Uygun et al., 2017).

216 To get a deeper understanding of -Fe response regulation, we elucidate the underlying
217 *cis*-regulatory code. Because some TFs have well established roles in -Fe response, we can
218 use these to validate our findings. We combined genome, transcriptome, and *in vitro* protein-
219 DNA interaction data to uncover links between pCREs controlling -Fe responses and their

220 upstream TFs. With pCREs over-represented in promoters of co-expressed genes we modeled
221 -Fe-induced up-regulation and identified over 100 informative pCREs of -Fe-responsive
222 processes.

223

224 **Results and Discussion**

225 **Overview of approach and functions of -Fe-responsive genes**

226 To identify root -Fe-associated CREs at a genome-wide scale, we defined root -Fe
227 response co-expression clusters, then we identified *k*-mers enriched in the promoter regions
228 of those genes, and finally we modeled -Fe response on the basis of the enriched promoter *k*-
229 mers. An overview of our complete workflow including functional analysis of the identified
230 pCREs is shown in **Figure 1A**. Because many factors, such as the choice of data set or the
231 measure used to define expression similarity, impact the discovery of functional connections
232 between genes (Uygun et al., 2016), we used multiple expression data combinations and
233 algorithms with varying parameters (see **Methods**).

234 -Fe-responsive genes (\log_2 fold-change (\log_2FC) >1 or <-1 , $q<0.05$) were identified
235 using transcriptomic data available for six time points after an -Fe treatment in Arabidopsis
236 seedling roots (Dinneny et al., 2008). Enrichment analysis (Fisher's exact test (FET), $q<0.05$)
237 of biological process gene ontologies (GOs) showed that, next to Fe-related GOs (e.g. Fe
238 transport, homeostasis and FeS cluster assembly), responses to several hormones, including
239 auxin, ethylene, abscisic acid and jasmonic acid, were over-represented (**Figure 1B**). This is
240 consistent with the roles of hormones in -Fe response (Brumbarova et al., 2015) and in root
241 and root hair morphology and development (Schmidt et al., 2000), which were also enriched
242 GOs. -Fe affects the photosynthetic machinery and often correlates with oxidative stress
243 responses (Rodríguez-Celma et al., 2013), which is reflected in enrichment of GOs regarding
244 oxidative stress, photosynthesis, and primary metabolism even in roots (**Figure 1B**;
245 **Supplemental Figure S1**).

246 **Using multiple expression data sets to define -Fe co-expression clusters**

247 We next grouped differentially regulated -Fe response genes into co-expression
248 clusters using two approaches: *k*-means clustering and correlation to gold standard. *K*-means
249 clustering was based on the transcriptional responses to -Fe alone and combined with different
250 responses to other stress and developmental conditions (**Figure 1C**) (Schmid et al., 2005;
251 Kilian et al., 2007; Dinneny et al., 2008; Goda et al., 2008). Correlation-based clusters were
252 generated for each gene in our curated list of gold standard (GS) -Fe response genes (see
253 **Methods; Supplemental Table S1**), by selecting the differentially regulated -Fe response
254 genes with a significantly similar (Pearson's Correlation Coefficient (PCC); see **Methods**)
255 expression pattern to the GS gene, also using the different combinations of transcriptional data

256 **(Figure 1C)**. To identify co-expression clusters with similar biological functions, we grouped
257 them according to their enriched GOs (FET, $q < 0.05$) into “superclusters” **(Figure 2A;**
258 **Supplemental Figure S2A;** see **Methods**), which were defined as groups of at least 20
259 clusters with significantly higher similarity to each other than the average similarity of all
260 clusters (all Mann-Whitney U, $p < 2.2e-16$; **Supplemental Figure S2B, C**). While *k*-means
261 supercluster C was enriched in an Fe-related GO (cellular response to Fe, GO shared by $\geq 75\%$
262 of co-expression clusters within each supercluster), *k*-means superclusters A and B shared
263 GOs related to different stress responses.

264 Because the current TAIR GO annotation for -Fe response-related processes does not
265 contain all -Fe-responsive genes of interest (e.g. *MYB10*, *UGT72E1*, *AT3G07720*, *FEP3* or
266 *NAS4*, (Ivanov et al., 2012)), we also determined if *k*-means-generated co-expression clusters
267 were enriched for GS genes (“GS-enriched”; FET, $q < 0.05$; **right, Figure 2A**). While many of
268 these clusters were part of the Fe-related GO supercluster A, the GS approach allowed us to
269 identify an additional 23 Fe-related co-expression clusters that would have been overlooked
270 by conventional GO enrichment analysis. In total, 7% of the *k*-means-generated clusters were
271 GS-enriched **(Figure 2B)**. Applying this same analysis to the correlation-based clusters
272 **(Figure 2A; Supplemental Figure S2A)**, we found higher levels of similarity between
273 correlation-based clusters compared to *k*-means clusters (Mann-Whitney U, $p < 2.2e-16$;
274 **Supplemental Figure S2D**), because we pre-condition their identification on GS genes, some
275 of which are tightly co-regulated (Ivanov et al., 2012). Accordingly, we found that 93% were
276 GS-enriched **(right, Figure 2B)**.

277 GS genes are either FIT target (“FIT-dependent”) or FIT-independent Fe homeostasis
278 (“FIT-independent”) genes, which we found reflected in our GS-enriched clusters: 71% of the
279 GS-enriched clusters were more specifically enriched for FIT-dependent (39%) and/or FIT-
280 independent genes (32%) (FET, $q < 0.05$; **bottom, Figure 2B**). The remaining GS-enriched
281 clusters were enriched for both (“mixed”). Interestingly, clusters based on combined
282 expression data (i.e. data combinations (dc) 2, 3, 5a/b, 6) were more often enriched for FIT-
283 dependent or FIT-independent genes, while -Fe time course data alone (dc1) produced mainly
284 mixed category clusters **(Figure 2C)**. The utility of including spatial or developmental data (dc2,
285 6) to define co-expression clusters reflects that -Fe-responsive genes act at different time
286 points and in different tissues and organs (Dinneny et al., 2008; Ivanov et al., 2012; Jeong et
287 al., 2017). Finally, genes in clusters not enriched for GS-genes (“non-enriched”) tended to
288 respond to particular abiotic stresses, for example cold **(Supplemental Figure S8;** e.g.
289 clusters 937, 973) or salt (e.g. clusters 900, 915, 936, 987), whereas gene expression for GS-
290 enriched cluster genes tended to randomly oscillate under different abiotic stresses (e.g.
291 clusters 818, 835, 858, 889), which might indicate different regulatory networks and highlight

292 the usefulness of incorporating additional abiotic stress data (as in dc3, 5, 6) when defining co-
293 expression clusters that are likely co-regulated.

294 In summary, by using different expression data sets and clustering methods we defined
295 1,959 -Fe co-expression clusters, many of which were enriched for FIT-dependent and/or FIT-
296 independent GS genes. These represent possibly co-regulated functional units in Fe
297 acquisition and Fe homeostasis processes, well-suited to identify pCREs which can
298 explain -Fe-induced up-regulation. Genes in co-expression clusters that were enriched in -Fe-
299 responsive genes but not GS genes (non-enriched clusters) are presumably regulated by
300 mechanisms different from GS-enriched clusters.

301 **A machine learning approach to model regulation of -Fe responsive co-expression** 302 **clusters**

303 The machine learning algorithm Random Forest (RF) has been successfully used to
304 model stress transcriptional response using *cis*-regulatory sequences in plants (Zou et al.,
305 2011; Deng et al., 2017; Uygun et al., 2017). Here, for each co-expression cluster, we used
306 pCREs (enriched *k*-mers in putative promoter sequences; see **Methods**) to build a RF model
307 that classifies genes as belonging to the cluster in question or as a non-responsive gene (see
308 **Methods**). The pCREs from models performing above a defined threshold ($F1 \geq 0.7$; see
309 **Methods**) were then considered further. Out of 1,959 co-expression clusters, 28% of the
310 models passed the performance threshold, 60% performed poorly, and for 12% no model could
311 be built due to small size (median size=2 genes; **Supplemental Figure S3A, B**). Poor
312 performing models (median $F1=0.62$) were mostly for small clusters (median size=12)
313 (**Supplemental Figure S3B; Supplemental Table S2**) likely due to the lack of training data.
314 Nonetheless, 66 large clusters (>100 genes, median size=135) also performed poorly (median
315 $F1=0.64$) – this is likely because these large clusters are too heterogeneous containing genes
316 with multiple regulatory codes (Uygun et al., 2016; Uygun et al., 2017), and/or are co-regulated
317 but not at the transcriptional level (e.g. post-translationally controlled). Interestingly, of the 28%
318 of clusters with models above the threshold, only 36% were GS-enriched clusters
319 (**Supplemental Figure S3A**). Nonetheless, models built for GS-enriched clusters (median
320 $F1=0.68$) tended to perform better than models built for non-enriched clusters (median
321 $F1=0.65$; Mann-Whitney U, $p < 2.358e-09$; **Figure 3A-C**).

322 Good model performance indicates that genes in a cluster are more likely co-regulated,
323 and, because pCREs were used to build the model, these pCREs are likely the regulatory
324 sequences contributing to the co-regulation. Taken together, we identified 5,639 pCREs
325 enriched in promoters of -Fe-responsive genes that may be predictive of -Fe-induced up- or
326 down-regulation. To further evaluate the biological relevance of pCREs, in the following
327 sections, we assess pCREs based on their association with GS-enriched or non-enriched
328 clusters, importance for model performance, and similarity to known TF binding sites. Known

329 -Fe CREs from Arabidopsis and also from grasses, for example E-/G-boxes (bHLH TF binding
330 sites) and IDE1, will serve as positive controls.

331 **Identifying common pCREs across co-expression clusters**

332 We expect true -Fe response CREs to be: (i) important for building models with good
333 performance in predicting -Fe response, and (ii) reliably identified in co-expression clusters
334 with similar gene content. Therefore, for each pCRE we calculated the proportion of clusters
335 enriched for the pCRE and its average importance rank across those clusters (**Supplemental**
336 **Table S3**). The importance rank of a pCRE was derived from an importance score for the
337 pCRE in question from the trained RF models that reflects how useful a pCRE was for
338 predicting -Fe response genes in a cluster. This allowed us to get a snapshot of pattern of
339 presence and absence of important pCREs for genes correctly predicted (true positives (TP))
340 in co-expression clusters with good (**Figure 3D, E**) and poor (**Figure 3F**) performance. The
341 pCREs were enriched in between 1 (0.6%) and 56 (35%) GS-enriched clusters and in between
342 1 (0.3%) and 54 (15%) non-enriched clusters, with 173 pCREs considered frequent pCREs
343 (freq-pCREs, enriched in >5% of GS-enriched or non-enriched clusters) (**Supplemental Table**
344 **S3**). Across GS-enriched clusters, pCREs tended to have higher proportions with higher
345 importance ranks than across clusters that were not GS-enriched (Mann-Whitney U, $p < 2.2e-$
346 16 ; **Figure 4A, B**; **Supplemental Figure S4A, B**). The higher proportion of GS-enriched
347 cluster pCREs can be explained partly by the fact that those clusters are more homogenous
348 in terms of gene contents than non-enriched clusters (Mann-Whitney U, $p < 2.2e-16$;
349 **Supplemental Figure S3C**).

350 We next determined whether GS-enriched clusters are regulated by a different set of
351 pCREs than non-enriched clusters. Out of the 5,639 pCREs, 15% (n=860) were unique to GS-
352 enriched clusters, while 73% (n=4109) were unique to non-enriched clusters and 12% (n=670)
353 were found in both GS-enriched and non-enriched clusters (**inset, Figure 4A**). This indicates
354 that GS-enriched clusters and non-enriched clusters are regulated partly by different pCREs,
355 but also by a fraction of shared pCREs. However, 43% (n=286) of the 670 shared pCREs were
356 predominant to GS-enriched clusters (i.e. having only low proportion and low importance rank
357 in non-enriched clusters; **top, Supplemental Figure S4C**). This indicates that pCREs that
358 were categorized as shared might not be equally important for regulating both GS-enriched
359 and non-enriched clusters. Interestingly, unique GS-enriched freq-pCREs represented 59%
360 (n=102) of the 173 freq-pCREs, while 35.5% (n=58) were unique non-enriched, and only 7.5%
361 (n=13) freq-pCREs were shared between GS-enriched and non-enriched clusters, indicating
362 that pCREs with high proportion are also the ones which seem to regulate almost exclusively
363 either GS-enriched or non-enriched cluster functions, but not both (**inset; Figure 4A; bottom,**
364 **Supplemental Figure S4C**). Furthermore, freq-pCREs tended to have higher importance
365 ranks than non-frequent pCREs (Mann-Whitney U, $p < 1.924e-14$; **Supplemental Figure S5D**).

366 Together, this suggests that freq-pCREs could be particularly relevant for regulation of -Fe
367 response mechanisms.

368 To characterize the freq-pCREs, we grouped them according to sequence similarity
369 using pair-wise PCC distances of pCRE position weight matrices (PWM; see **Methods**). 62%
370 (n=107) of all freq-pCREs could be placed into one of eight pCRE groups (**Figure 4C, D**;
371 **Supplemental Figure S4E**). Freq-pCREs of the same group tended to be predictors of the
372 same cluster category (GS-enriched/non-enriched).

373 In summary, we identified more than 100 -Fe pCREs that were reliably associated
374 either exclusively to GS-enriched or non-enriched co-expression clusters or with high
375 preference for one of the categories. Those pCREs were also ranked as important for machine
376 learning models and might therefore be candidates for functionally relevant motifs to different
377 responses to -Fe.

378 **Similarity of -Fe pCREs to known TFBMs**

379 CREs are recognized by TFs to modulate gene expression. To identify what types of
380 TFs may bind to the identified pCREs, we examined the similarities between the -Fe pCREs
381 and known TF binding motifs (TFBMs) from two sources (see **Methods**). Based on threshold
382 similarities, we were able to match a specific TF and/or a specific TF family to each of the 173
383 freq-pCREs (see **Methods; Figure 5A; Supplemental Figure S5**). To gain an overview which
384 TF families might be associated with GS-enriched clusters and how specific these TF families
385 are, we asked which families contained over-represented numbers of TFs that likely bound
386 pCREs from GS-enriched and non-enriched cluster categories. We found that most TF families
387 were found with higher proportion in either GS-enriched clusters (14 TF families) or non-
388 enriched clusters (12), while only four were similarly distributed between both categories
389 (**Figure 5B**).

390 Most known -Fe regulators in Arabidopsis are bHLH TFs (FIT, subgroup Ib and IVc
391 bHLH proteins, PYE, e.g. (Jakoby et al., 2004; Wang et al., 2007; Long et al., 2010; Palmer et
392 al., 2013; Zhang et al., 2015)). bHLH and MYB TF families were identified, and even with
393 higher proportion in GS-enriched clusters than in non-enriched clusters, which is indeed
394 consistent with their role in -Fe response regulation. Other matching TF families over-
395 represented in GS-enriched clusters were bZIP (FET, $q<0.05$), B3, TCP and NAC. Although a
396 B3 TF (ABI3VP1 subfamily; IDEF1) and a NAC TF (IDEF2) are important regulators of Strategy
397 II Fe acquisition in grasses (Kobayashi et al., 2007; Ogo et al., 2008), the role for these TF
398 families in Strategy I non-grass plant species has not yet been described. In contrast, ARID,
399 WRKY (both FET, $q<0.05$), Homeobox, and CAMTA TF families were matched more in non-
400 enriched than GS-enriched clusters, pointing towards roles during -Fe stress other than Fe
401 uptake or homeostasis, in which GS genes are mostly involved.

402 Next, freq-pCRE-TFBM matches from GS-enriched clusters served to infer specific
403 upstream regulators of -Fe-responsive modules. More than 50% freq-pCREs (60 out of 115)
404 matched TFBMs of a specific TF (**Figure 5A**). Of those, 29 freq-pCREs shared perfect
405 sequence similarity (PCC=1) to the TFBM, which were then of particular interest. From these
406 perfect matches, 23 pCREs were unique for GS-enriched clusters. Example TF candidates for
407 these 23 cases were FUS3 (an ABI3VP1/B3 TF), bHLH104, bZIP3, 16 and 42, TCP13 (PTF1),
408 and FAR1 (**Supplemental Table S4A**). While the DAP-seq and CIS-BP TFBM databases
409 contain binding information for many TFs, they are far from exhaustive. For example, out of
410 162 known Arabidopsis bHLHs (Bailey et al., 2003), only 46 were available to be included in
411 the analysis. Therefore, some TF families were likely under-represented in our analysis and
412 some top match TFBMs may not accurately reflect the binding partner for certain pCREs.
413 Consequently, some important pCRE-TFBM matches might not be detectable at this time.
414 However, as new experimental TF binding data is collected, we might gain more biological
415 insight into our -Fe pCREs.

416 **Inferring upstream regulators of the -Fe response**

417 Because we believe our genome-wide approach for identifying regulatory elements
418 may shed light on areas of -Fe response that are less well understood, we next put our findings
419 in context with open questions in the field. For example, the ABI3VP1/B3-type TF IDEF1, a
420 key regulatory factor of -Fe responses in rice and barley roots, recognizes the CATGC core of
421 IDE1 (Kobayashi et al., 2003; Kobayashi et al., 2005; Kobayashi et al., 2007; Kobayashi et al.,
422 2009; Kobayashi et al., 2010). With ten of our freq-pCREs having an IDE1 CATGC (or GCATG)
423 core and matching ABI3VP1/B3 family TFBMs, IDE1-likes were fairly dominant among the
424 freq-pCREs and unique to GS-enriched clusters (**Supplemental Table S3**). This strongly
425 suggests an important function for IDE1-like motifs in Arabidopsis. Arabidopsis AFLs (B3 family
426 TFs ABI3/FUS3/LEC2), are the closest homologs of the rice IDEF1, and may bind to the IDE1-
427 likes. In fact, ABI3 and FUS3 bind to RY-like elements (CATGCA), regulating FeS cluster
428 subunit formation during seed maturation (Roschzttardtz et al., 2009). However, ABI3 or FUS3
429 functions during later developmental stages, particularly in the root during -Fe response,
430 remain to be elucidated. Since the FUS3 TFBM matched our top most abundant IDE1-like
431 (CATGCC; **Supplemental Table S4A**), and because *FUS3* is expressed in the root epidermis
432 and in lateral root primordia during later developmental stages (Boulard et al., 2017; Tang et
433 al., 2017), FUS3 might be an IDEF1 homolog in Strategy I plants.

434 Another -Fe response-related TF in rice and barley, IDEF2, belongs to the NAC family
435 and binds to the CA(A/C)G(T/C)(T/C/A)(T/C/A) core in IDE2 (Ogo et al., 2008). Although we
436 did not have a perfect (PCC=1) pCRE-NAC TFBM match, we found matched NAC TFBMs
437 slightly over-represented in GS-enriched clusters (**Figure 5B**). Furthermore, two of the top ten
438 most abundant freq-pCREs unique to GS-enriched clusters matched NAC TFBMs (PCC>0.9),

439 with one freq-pCRE being highly similar to the IDE2 core (CACGCC). This indicates that IDE2-
440 like motifs might also play a role during Arabidopsis -Fe responses.

441 One freq-pCRE (CGTGCC) perfectly matched to a bHLH104 TFBM (**Supplemental**
442 **Table S4A**). bHLH104 binds to the promoters of subgroup Ib *BHLH* genes
443 *BHLH38/39/100/101* (Zhang et al., 2015; Li et al., 2016), positively regulating Fe uptake.
444 Consistently, we found CGTGCC in clusters containing *BHLH101* (*AT5G04150*;
445 **Supplemental Table S2**).

446 Other freq-pCREs matched to known TFBMs from TFs with unknown roles in -Fe
447 response. For example, bZIP TFBMs were significantly over-represented in GS-enriched
448 clusters, but have no known direct roles in -Fe response. However, bZIP TFs are known
449 regulators of the Zn deficiency response, which, together with the fact that one GS gene, *ZIP9*,
450 is also responsive to Zn deficiency, could indicate an interdependency of Zn and Fe
451 homeostasis (Assunção et al., 2010; Sinclair et al., 2018). Furthermore, two matched TFs,
452 bZIP3 and bZIP16, are involved in ABA signaling, which is connected to -Fe response amongst
453 others by modulating root growth (Séguéla et al., 2008; Matioli et al., 2011; Hsieh et al., 2012).
454 Possible functions of bZIP TFs in response to -Fe stress should be explored in the future.

455 TCP13 (PTF1) and FAR1 TFBMs are two more examples for perfect freq-pCRE
456 matches with yet unknown specific roles of the TFs during -Fe, although their specificity to GS-
457 enriched clusters points towards important roles in regulating GS genes. TCPs are involved in
458 plant development, but also act in signaling of hormones that influence -Fe responses (Davière
459 et al., 2014; Brumbarova et al., 2015; Resentini et al., 2015; Nicolas and Cubas, 2016). For
460 example, TCP20 was reported to bind to the *BHLH39* promoter (Andriankaja et al., 2014),
461 indicating a possible connection of TCPs and -Fe responses during plant development. TCP13
462 is involved in regulating responses to light shade signals through *PHYTOCHROME*
463 *INTERACTING FACTORS* (*PIFs*) (Zhou et al., 2018). Interestingly, FAR1 and its homolog
464 FHY3 also act in phytochrome-PIF signaling (Wang and Wang, 2015). Together, this suggests
465 a connection of light perception and -Fe responses mediated through these TFs, which is
466 consistent with the known diurnal influence on Fe uptake (Vert et al., 2003; Santi and Schmidt,
467 2009; Hong et al., 2013; Salomé et al., 2013). In addition, FAR1/FHY3 act in the regulation of
468 phosphate starvation response, together with ethylene regulator EIN3 (Liu et al., 2017), which
469 also binds FIT to promote Fe uptake (Lingam et al., 2011). Therefore, FAR1 might also regulate
470 Fe acquisition via the ethylene pathway.

471 Finally, a perfect freq-pCRE-WRKY11 match indicates that WRKY TFs, although
472 significantly over-represented in non-enriched clusters, are also important for regulating GS-
473 enriched clusters. WRKY11 is involved in abiotic stress tolerance in Arabidopsis (Ali et al.,
474 2018), with no specific role known during -Fe response yet. However, WRKYs in general have
475 already been connected to -Fe, for example as putative regulators of the coumarin transporter

476 gene *PDR9* (Ito and Gray, 2006) and of *PYE* (Koryachko et al., 2015). Furthermore, WRKY46
477 negatively regulates the vacuolar Fe importer gene *VTL1/VITL1* (Gollhofer et al., 2011;
478 Gollhofer et al., 2014; Yan et al., 2016). We found a WRKY TFBM (GTCAAC) in several non-
479 enriched clusters containing down-regulated Fe-responsive genes, including the *VTL1*
480 homolog *VTL5* (*AT3G25190*; **Supplemental Table S2**), indicating that some of the TFs
481 matching non-enriched cluster pCREs might act as repressors of Fe excess genes.

482 In summary, many pCREs commonly found among GS-enriched clusters shared
483 significant sequence similarity with known -Fe CREs, such as IDE1, or with binding sites of
484 known -Fe-associated TF families, such as ABI3VP1/B3, NAC, MYB and bHLH. Notably, we
485 found evidence for IDE1-like motifs being relevant not only in Strategy II plants, but also in the
486 Strategy I plant Arabidopsis. Our results also suggest novel associations, such as the role of
487 bZIPs or TCPs in -Fe responses. We assessed in the next paragraph in which specific -Fe
488 response processes pCREs of particular interest, such as IDE1-likes, might be involved in.

489 **Associating important pCREs with FIT-dependent or FIT-independent functions**

490 After identifying novel potential regulators in the -Fe response, we pinpointed some of
491 those which could best explain models of -Fe-responsive up-regulation and explored their
492 potential functions.

493 More than 1,500 pCREs were identified in total in GS-enriched clusters, raising the
494 question of a core set of important pCREs needed to robustly predict -Fe response in each
495 cluster. Using pCRE abundance (freq-pCREs) among GS-enriched clusters as the only criteria
496 for selecting informative motifs for those clusters could result in missing motifs simply due to
497 the fact that some co-expression clusters were more unique than others. This is supported by
498 the fact that rare pCREs still can have a high importance rank (**Figure 4A**), meaning that those
499 pCREs were not included in the set of freq-pCREs although they seem to be important for
500 regulating individual GS-enriched clusters. We identified the most important pCREs (defined
501 as the minimum set of pCREs; min-pCREs) by building RF models iteratively with successively
502 deleting the least important pCREs in each round (**Supplemental Figure S6**). Applying this
503 approach to the 159 GS-enriched clusters resulted in a collective set of 615 min-pCREs. They
504 were part of the minimum sets of between 1 (0.6%) and 48 (30%) GS-enriched clusters, with
505 the IDE1-like CATGCC being the top most abundant min-pCRE. Together with CATGCC, two
506 more IDE1-like motifs, TCATGC and CCATGC, were among the top ten most abundant min-
507 pCREs (**Supplemental Table S3**). This supports a previous computational analysis of rice
508 promoters, in which IDE1-like was among the top scoring motifs (Takei et al., 2013). Together
509 with our previous finding typing IDE1-like ABI3VP1/B3 TFBMs to GS-enriched clusters, it
510 suggests an important, yet unknown, function of IDE1-like motifs in Arabidopsis -Fe response
511 regulation. Min-pCREs matching a bHLH (CGTGAC), a MYB (TAACTA), and the IDE2-like
512 NAC TFBM (CACGCC; all **Supplemental Table S4A**) were also among the top ten most

513 abundant min-pCREs (**Supplemental Table S3**), further demonstrating the utility of our
514 approach.

515 To determine in which processes min-pCREs might function during -Fe, we tested if
516 min-pCREs were more likely to be found in FIT-dependent or FIT-independent co-expression
517 clusters. More than 60% of the 159 GS-enriched clusters were classified as either FIT-
518 dependent with a likely function in root iron acquisition (35% out of 159) or FIT-independent
519 with either a function in internal Fe homeostasis in shoots and roots or in -Fe response
520 regulation (28% out of 159; **Figure 6A**). We then calculated the proportion of min-pCREs
521 (present in ≥ 5 GS-enriched clusters) in each cluster category (**Supplemental Table S5**).
522 Interestingly, the two IDE1-like motifs, CATGCC, CCATGC and the related ABI3VP1/B3 TFBM
523 matched ATGCAT, were predominantly identified in FIT-dependent clusters, but IDE2-like
524 CACGCC had no preference for either FIT-dependent or FIT-independent clusters (**Figure**
525 **6B**). This suggests that the IDE1-like pCREs tend to be more important for FIT-dependent root
526 Fe acquisition rather than FIT-independent Fe sensing, signaling and distribution. This is also
527 consistent with the role of grass IDE1 in Fe uptake (Kobayashi et al., 2003; Kobayashi et al.,
528 2005). Two ARF TFBM matched min-pCREs (AACGTA/ARF16, GTCGGA/ARF2) were also
529 preferentially found in FIT-dependent clusters. ARFs are involved in auxin signaling, thereby
530 controlling - among other functions - root hair elongation (Pitts et al., 1998; Mangano et al.,
531 2017; Choi et al., 2018). Different studies reported that -Fe responses can be accompanied by
532 an increase of root hair number, elongation of root hairs, deformed or short root hairs (Schmidt
533 et al., 2000; Müller and Schmidt, 2004; Dinneny et al., 2008). ARF2 and ARF16 TFs are root
534 hair growth repressors (Choi et al., 2018), which would be consistent with a short root hair
535 phenotype and down-regulation of respective GO terms under -Fe (Dinneny et al., 2008)
536 (**Supplemental Figure S1**). During -Fe, several root hair-acting genes are co-expressed in a
537 regulon which also contains *IRT2* (Ivanov et al., 2012), indicating a possible connection of ARF
538 TFBM matched min-pCREs with these root hair processes.

539 In contrast, three bZIP TFBM matched min-pCREs, GTGGCA, CACGTC and CACTAC,
540 were predominantly identified in FIT-independent clusters. As described in the previous
541 section, bZIPs are involved in ABA signaling. ABA negatively regulates FIT-dependent Fe³⁺
542 reductase gene *FRO2* and Fe²⁺ importer gene *IRT1* (Séguéla et al., 2008). However, ABA
543 signaling also leads to enhanced apoplastic and vacuolar Fe utilization and root to shoot
544 transport under -Fe (Séguéla et al., 2008; Lei et al., 2014). (Lei et al., 2014) propose that ABA-
545 responsive gene regulation and Fe remobilization and transport are connected through bZIPs,
546 which is consistent with our results that bZIP TFBMs are preferentially found in FIT-
547 independent co-expression clusters of genes involved in Fe mobilization and translocation.
548 Similarly, two TCP TFBMs, GACCAC and ACCCAC, were identified almost exclusively in FIT-
549 independent clusters, which is in agreement with TCP20 regulating *BHLH39* in a FIT-

550 independent manner (Andriankaja et al., 2014). Finally, bHLH TFBMs were identified in all
551 cluster categories with a preference for mixed clusters. This matches the ubiquitous nature of
552 bHLH target motifs (E-/G-boxes), which act at many levels in the -Fe bHLH cascade. In
553 summary, we can propose plausible roles for pCREs as TFBMs in FIT-dependent and FIT-
554 independent processes.

555 **Distribution and conservation of min-pCREs in co-expression cluster gene promoters**

556 Next, to explore if min-pCREs displayed significant positional bias in the promoter
557 regions of co-expressed genes, we compared the observed min-pCRE frequencies in 100 bp
558 bins of -1000 to +500 bp and of -500 to +1000 bp flanking regions adjacent to the transcription
559 start site (TSS) and transcription termination site (TTS), respectively, with the expected
560 frequencies from shuffled pCRE sequences (according to Uygun et al., 2017) for all 615 min-
561 pCREs (**Supplemental Figure S9**). Furthermore, we examined the non-coding as well as the
562 coding sequences of the transcribed regions. Overall, the distributions of min-pCREs revealed
563 a slight positional bias in the promoter regions (**top, Figure 6C**). We investigated the
564 distribution plots separately for ten selected min-pCREs: FIT-dependent ABI3VP1/B3 TFBM-
565 matched min-pCREs (containing the two IDE1-likes), the IDE2-like (NAC TFBM match), FIT-
566 independent bZIP TFBM matches, and bHLH TFBM matches (**Figure 6C**). These min-pCREs
567 had significant location bias in the putative promoters up to 1000 bp upstream of the TSS.
568 Because known CREs often exhibit positional bias (Zou et al., 2011; Heyndrickx et al., 2014;
569 Yu et al., 2016), this provides additional support for these pCREs having regulatory functions
570 in Fe uptake and homeostasis. Interestingly, some of the pCREs common to mixed clusters
571 did not show position bias in any of the genomic regions tested (e.g. ABI3VP1/CTTATA and
572 MYB/TAACTA; **bottom, Figure 6C**), indicating that genes of such clusters are less likely to be
573 transcriptionally co-regulated.

574 Next, we sought to pinpoint the specific processes of Fe homeostasis, which these ten
575 min-pCREs might regulate. To assess this, we determined the genes containing the respective
576 min-pCRE and counted the number of incidents in which these genes were likely regulated by
577 the min-pCRE (**Figure 6D; Supplemental Table S6A**). For example, CATGCC was found in
578 48 GS-enriched clusters, and 36 of those clusters (75%) included the min-pCRE-containing
579 gene *MYB10*, while only four of those clusters (8%) included *MAPKKK16* (**second row right,**
580 **Figure 6D**). We inferred that CATGCC might be regulating predominantly processes in which
581 MYB10 is required. As an additional line of evidence for pCRE functionality, we determined if
582 min-pCREs overlapped with conserved noncoding sequences (CNS) of the Brassicaceae
583 family (Haudry et al., 2013) (**Supplemental Table S6B**). As expected, bHLH TFBM matched
584 min-pCREs were identified in many gene promoters, including *BHLH39/BHLH101* (direct
585 targets of bHLH IVc TFs, (Zhang et al., 2015)), *NAS4* (direct PYE target, (Long et al., 2010)),
586 and *IRT1*, *AT3G07720* and *GRF11* (FIT targets; **Figure 6D**; (Sivitz et al., 2012; Yang et al.,

587 2013)). In *BHLH39/BHLH101*, *IRT1*, and *GRF11*, the respective min-pCREs overlapped with
588 CNS, further supporting the importance of these motifs. Interestingly, bHLH matched min-
589 pCREs were also located in CNS of *BTS* and *BTSL1*, two genes that negatively regulate Fe
590 uptake by marking positive regulators (e.g. bHLH IVc TFs) for degradation (Selote et al., 2015).
591 If *BTS* were to be regulated by bHLH proteins from the same regulatory cascade, this may
592 indicate a negative feedback loop.

593 FIT-dependent IDE1-likes CATGCC and CCATGC were located in the *IRT1* promoter,
594 overlapping with CNS (**second row, Figure 6D**). Interestingly, both IDE1-likes were found in
595 several genes encoding enzymes and TFs involved in coumarin biosynthesis (*CYP82C4*, *S8H*,
596 *F6'H1*, *MYB72/MYB10*; (Kai et al., 2008; Murgia et al., 2011; Fourcroy et al., 2014; Schmid et
597 al., 2014; Zamioudis et al., 2014; Rajniak et al., 2018; Siwinska et al., 2018)). We propose that
598 IDE1 is important for synthesis of Fe chelators in response to low Fe conditions in both
599 monocots and dicots (see Kobayashi et al., 2003; Kobayashi et al., 2005). The IDE2-like min-
600 pCRE was located amongst others in *BTSL1* (within a CNS), *BHLH39*, *ORG1*, and a number
601 of uncharacterized genes. While the role of IDE2 in the Strategy II Fe response has not been
602 comprehensively explored, it is known to regulate expression of *OsYSL2*, a phloem Fe²⁺-
603 nicotianamine transporter (Kobayashi et al., 2003; Ogo et al., 2008). This putative involvement
604 in phloem translocation of Fe suggests that the IDE2-like might preferably associate with FIT-
605 independent gene functions. However, we could not assess the relationship between IDE2 and
606 YSLs in this analysis because *YSL1/2/3* were not expressed above the log₂FC>1 threshold.

607 Next, we explored the bZIP TFBM-matched min-pCREs, since they had the strongest
608 FIT-independent preference. These min-pCREs were located in a number of genes involved
609 in translocation of Fe/Fe chelates or the synthesis of Fe chelators (e.g. *NRAMP4*, *ZIF1*, *OPT3*
610 (Lanquar et al., 2005; Haydon et al., 2012; Mendoza-Cózatl et al., 2014; Zhai et al., 2014)) or
611 in Fe sensing and signaling (e.g. *OPT3*, *BTS*, (Mendoza-Cózatl et al., 2014; Zhai et al., 2014;
612 Selote et al., 2015; Khan et al., 2018)). Furthermore, GTGGCA (matched to bZIP TFBM)
613 overlapped with a CNS of *CGLD27* (**third row, Figure 6D**), which has been associated with
614 photoprotection in leaves during -Fe (Ruiz-Sola and Rodríguez-Concepción, 2012; Rodríguez-
615 Celma et al., 2013). However its function in roots remains elusive. Taken together, our findings
616 suggest diverse roles for bZIP TFBMs, including Fe transport and adjustment of the plastid
617 proteome.

618 In summary, we identified more than 100 -Fe pCREs which, in addition to sharing
619 significant sequence similarity to known TFBMs, were also part of the core sets of pCREs
620 needed for robust prediction of -Fe responses of GS-enriched clusters (min-pCREs).
621 Furthermore, they were preferentially located in promoter regions upstream of the TSS, and
622 even in CNS' of some genes. Together, these findings indicate that these pCREs might be
623 authentic -Fe CREs. From the biological context of the genes which are likely regulated by

624 some of the pCREs, we were able to greatly improve our understanding of -Fe response
625 regulation in Arabidopsis. For example, our work highlighted that in addition to the bHLH
626 TFs, IDE1-like motifs and bZIP TFs are likely involved in different responses to -Fe and
627 should be considered of high interest for future work.

628

629 **Conclusion**

630 We identified 5,639 pCREs enriched in promoters of co-expressed -Fe-responsive
631 genes that were used as features to predict -Fe-responsive regulation of root-expressed genes
632 on a genome-wide scale. Of those, 173 reliably predicted -Fe response genes of >5% of our
633 defined co-expression clusters (freq-pCREs). Because most of those pCREs were either
634 unique to co-expression clusters enriched for our gold standard Fe acquisition and
635 homeostasis genes, or unique to co-expression clusters lacking those genes, we conclude that
636 our approach had captured motifs specifically regulating different responses during -Fe. To
637 take advantage of the publicly available *in vitro* TF binding information, we compared the freq-
638 pCREs to TFs from two studies (Weirauch et al., 2014; O'Malley et al., 2016), and found
639 that our approach had captured known Strategy I -Fe recognition motifs for bHLH and MYB
640 proteins. Our approach also led to novel regulatory connections of bZIP, B3, NAC, and TCP
641 families to Strategy I -Fe response regulation. While bZIP and bHLH TF families are also
642 associated with high salinity stress response (Uygun et al., 2017), other high salinity stress
643 response associated TF families (e.g. WRKY and AP2) were not common among our -Fe
644 pCREs regulating GS-enriched co-expression clusters, highlighting the usefulness of this
645 approach to pinpoint regulators specific to a stress condition.

646 We inferred possible functions of pCREs which were most important for modeling -Fe
647 responses (min-pCREs) from their enrichment in FIT-dependent or FIT-independent co-
648 expression clusters and their location bias in promoters of particular -Fe-responsive genes
649 (**Figure 6B, D**). Our results provide evidence that B3 TF pCREs containing the IDE1 core
650 motif CATGC are linked to coumarin synthesis, indicating that the function of IDE1-like motifs
651 to ensure supply of Fe-chelating compounds for Fe acquisition could be an evolutionarily
652 conserved function at least among flowering plants. While our results highlight the importance
653 of IDE1-like motifs for Fe acquisition, it was not the only prominent -Fe pCRE. This is in contrast
654 to Zn deficiency, where ZDRE seems to be singularly associated with multiple Zn deficiency
655 responses (Assunção et al., 2010), and indicates that despite of overlaps of Zn and Fe
656 homeostasis control (Briat et al., 2015), their transcriptional regulation must follow different
657 mechanisms.

658 Our results support a concept in which -Fe is not regulated by only one or few regulatory
659 elements. Of the many important pCREs for -Fe response, many share significant similarity
660 with TFs of TF families known to undergo hetero-dimerization and protein interaction across

661 families, such as bHLH, MYB, bZIP, TCP, and ABI3VP1 (Bemer et al., 2017). A combinatorial
662 mechanism would dramatically increase the flexibility of transcriptional responses driven by a
663 set of few TFs. It might be that some pCREs not as important in our prediction models, would
664 become informative in combination, as suggested for high salinity stress response (Zou et al.,
665 2011; Uygun et al., 2017). A next step would therefore be to build -Fe prediction models that
666 explicitly account for interactions between pCREs. A limitation of our approach is that our co-
667 expression clusters were based on ATH1 chip microarray data, the only comprehensive -Fe
668 time course transcriptome set available to date. Some important -Fe marker genes (e.g. *FRO2*)
669 are not represented on the chip and others might not have passed our significance threshold
670 because of sensitivity issues with the microarray technology. Additionally, we restricted our
671 analysis to the promoter region 1000 bp upstream of the TSS. While this is expected to cover
672 most important *cis*-acting elements and reduce the occurrence of promoters overlapping with
673 adjacent genes, introns as well as more distal promoter regions are known to harbor *cis*-acting
674 elements (Rose et al., 2008; Rose et al., 2016).

675 The large number of TFs known to be involved in -Fe-induced up-regulation points
676 towards the importance of transcriptional regulation. However -Fe responses are also heavily
677 controlled at the post-transcriptional and post-translational level (Lingam et al., 2011; Meiser
678 et al., 2011; Sivitz et al., 2011; Selote et al., 2015; Zhang et al., 2015; Gratz et al., 2019).
679 Naturally, our approach cannot cover such regulatory aspects. However, it allows us to predict
680 TF families, that may act upstream of the known -Fe-responsive genes. We suggest TFs of
681 the bZIP, ABI3VP1/B3, NAC, and TCP families as upstream regulators of -Fe response in the
682 root. Because major Fe sinks are located in the shoot, a systemic shoot-to-root signal must
683 exist for proper Fe supply (Vert et al., 2003; Garcia et al., 2013). Integrating shoot
684 transcriptomic data would expand our knowledge on how responses to -Fe stress are
685 coordinated at the whole-plant level.

686 In conclusion, we demonstrate that our machine learning-based approach can identify
687 pCREs for -Fe-induced gene up-regulation. This strategy can be applied to various stresses
688 and developmental conditions to elucidate regulatory mechanisms, especially when *cis*- and/or
689 *trans*-acting elements were previously elusive (Zou et al., 2011; Uygun et al., 2017). We
690 provide a comprehensive source of potential -Fe response *cis*-regulators for a wide range
691 of -Fe-responsive genes. Because the identified pCREs are potentially involved in enhancing
692 Fe uptake and translocation, they generate potential for future applications in engineering
693 plants with improved plant performance traits, e.g. higher nutritional value because of better
694 Fe allocation and coping with unfavorable soil conditions.

695

696 **Methods**

697 **Expression data processing and generation of multiple expression data combinations**

698 Expression data (Affymetrix ATH1) from an -Fe treatment time course experiment with
699 six time points and of four -Fe treated root zones (both Dinneny et al., 2008) were downloaded
700 from Gene Expression Omnibus (GEO, <https://www.ncbi.nlm.nih.gov/geo/>, GSE10502,
701 GSE10497) in CEL format, preprocessed, normalized and contrasted as described below.
702 AtGenExpress expression data (Affymetrix ATH1) of abiotic stresses ((Kilian et al., 2007),
703 GSE5620-5628 or TAIR-ME00325-330, only data of root samples were used), hormone
704 treatment ((Goda et al., 2008), GSE39384 or TAIR-ME00333-340, ME00343-344, ME00350-
705 352, ME00356) and plant development ((Schmid et al., 2005), GSE5629-5634 or TAIR-
706 ME00319) were downloaded from The Arabidopsis Information Resource (TAIR;
707 <https://www.arabidopsis.org/portals/expression/microarray/ATGenExpress.jsp>) preprocessed,
708 normalized and contrasted by S. Uygun (Uygun et al., 2016) as described below. Background
709 correction and quantile normalization of CEL files were performed with Robust Multi-Array
710 Average expression measure (RMA) using the Bioconductor affy package (Gautier et al.,
711 2004). The \log_2 fold-change (\log_2FC) in expression was calculated for all data sets except
712 developmental data by pairwise comparison of treatment and control experiments for each
713 treatment and time point. Contrast matrices and linear model fits were created using R and the
714 Bioconductor LIMMA package (Ritchie et al., 2015; Phipson et al., 2016). Because
715 developmental stages have no control treatment, absolute normalized fluorescence intensity
716 values were used. The p -values for \log_2FC or fluorescence intensities were corrected for
717 multiple testing (adjusted p -values= q) using the BH method (Benjamini and Hochberg, 1995).
718 Genes were regarded as -Fe responsive if $\text{abs}(\log_2FC) \geq 1$, and $q < 0.05$ at least at one -Fe
719 treatment time point or in at least one -Fe treated root zone. -Fe deficiency time course data
720 was combined with -Fe root zone expression data or ATGenExpress datasets in different
721 combinations (**Figure 1C**) either including only genes up-regulated (“up”) or all genes up- or
722 down- regulated (“up & down”) in ≥ 1 -Fe time point or root zone. This resulted in 12 different
723 expression data combinations.

724 **Co-expression clustering using k -means**

725 To cluster genes with similar expression pattern, k -means clustering (Hartigan and
726 Wong, 1979) was applied using the Euclidean distance as the similarity measure. Because k -
727 means returns a local optimum solution depending on the number of clusters (k) created and
728 the random selection of genes as initial “means”, the outcome varies with run (i.e. non-
729 deterministic). Therefore, different k (25, 30, 35, 40, 50, 70, 80, 100) were tested and the
730 clustering was repeated up to four times. We build machine learning models (see below) with
731 all clusters generated from expression data combinations (DC) 1, 2, 3 and 5 (**Figure 1C**). To

732 prevent confusion, we point out that the total number of *k*-means-generated clusters used to
733 build models represents several repeated clustering events of always the same two sets of -
734 Fe-responsive genes (up; up & down, see above). The clustering events differ in the DC which
735 was used and in the *k*. Two DC were excluded from the analysis: DC 4 produced identical
736 clusters as DC 1, which were therefore not considered. DC 6 contained different measuring
737 units (log₂FC and absolute normalized fluorescence intensity), and could not be handled by
738 the *k*-means algorithm.

739 **Co-expression clustering by correlation with GS genes**

740 To generate co-expression clusters based on gold standard (GS) genes
741 (**Supplemental Table S1**), each GS gene was used as a query to identify genes with similar
742 expression patterns. Briefly, for each expression data combination (DC; **Figure 1C**), the PCC
743 was calculated between the query gene and each gene in DC using SciPy
744 (<http://www.scipy.org>, (Jones et al., 2001)). Similar to (Uygun et al., 2016), a random
745 background PCC was calculated representing the null distribution of expression correlation by
746 calculating the PCC of 10,000 randomly selected gene pairs in DC and the 95th percentile of
747 PCCs was used as the threshold for classifying a pair of genes as significantly correlated. For
748 some DC (mostly those containing only up-regulated -Fe responsive genes), we allowed a
749 significance threshold below 90% down to 45%, because the PCC between random Fe
750 responsive genes was already very high. On the other hand, when >50 genes were considered
751 significantly correlated, the threshold was raised above 95% to 99% to hone in on genes most
752 likely to be co-regulated. In addition, we generated a second version of clusters with >50
753 genes, containing only the 10 genes with highest PCC. We build machine learning models
754 (see below) with both versions and further used the results from the better performing version
755 only. Percentiles used for each PCC-generated cluster are given in **Supplemental Table S2**.
756 Two DC were excluded from the analysis: DC 4 (up), because the resulting clusters were
757 identical to those generated from DC 1 (up), and DC 6 (up & down), because developmental
758 data seemed to have a disproportional influence on the PCC with the result that even -Fe up-
759 and down-regulated gene pairs were identified as strongly correlating. As in the *k*-means
760 clustering, the total number of PCC-generated clusters used for modeling represents repeated
761 clustering events of the same two sets of -Fe-responsive genes (as described above).

762 **The co-expression clusters: GO and GS/FIT-dependent/FIT-independent gene** 763 **enrichment and GO/gene content similarity**

764 Gene ontology (GO) associations for *A. thaliana* were downloaded from TAIR
765 (ftp://ftp.arabidopsis.org/home/tair/Ontologies/Gene_Ontology/ (Berardini et al., 2004)).
766 Biological process (BP) GO annotations were downloaded from GO
767 (<http://purl.obolibrary.org/obo/go.obo>) and parsed for BP information. Enrichment of GO terms

768 in genes that were significantly differentially regulated ($q < 0.05$, $\text{abs}(\log_2\text{FC}) \geq 1$) in the -Fe time
769 course data set was determined with a Fisher's exact test (FET, <http://www.scipy.org>, (Jones
770 et al., 2001)), and p -values were corrected for multiple testing ($=q$) using the "qvalue" function
771 in R (Storey, 2002) (**Supplemental Table S7**).

772 All co-expression clusters were tested for enrichment of GO terms as described above.
773 The similarity of enriched GOs between co-expression clusters was assessed using the
774 Jaccard Index (JI), or the intersection of GOs divided by the union of the GOs, where $\text{JI} = 1$ if
775 the exact same GOs were enriched in both co-expression clusters. Co-expression clusters
776 were grouped by hierarchical clustering using the JI with the UPGMA method in the R cluster
777 package (Maechler et al., 2017). Groups containing >20 co-expression clusters and having a
778 within-mean JI that was significantly higher than the mean JI of all clusters were defined as
779 "superclusters". Biological functions of superclusters were defined through GOs shared by
780 $\geq 75\%$ (k -means clustering) or $\geq 90\%$ (GS gene correlation; PCC) of the clusters. Similarly,
781 FET with p -value correction for multiple testing was used to identify co-expression clusters
782 enriched for (A) -Fe GS genes, (B) FIT-dependent genes, and/or (C) FIT-independent genes.

783 **K-mer enrichment and identification of pCREs predictive of -Fe response using Random** 784 **Forest (pCRE identification pipeline)**

785 Promoter sequences 1 kb upstream from the transcription start site (TSS) were
786 downloaded from TAIR
787 ([ftp://ftp.arabidopsis.org/home/tair/Sequences/blast_datasets/TAIR10_blastsets/upstream_](ftp://ftp.arabidopsis.org/home/tair/Sequences/blast_datasets/TAIR10_blastsets/upstream_sequences/TAIR10_upstream_1000_20101104)
788 [sequences/TAIR10_upstream_1000_20101104](ftp://ftp.arabidopsis.org/home/tair/Sequences/blast_datasets/TAIR10_blastsets/upstream_sequences/TAIR10_upstream_1000_20101104)). A list of all possible 6-mers of A, T, C, G was
789 generated with the Python `itertools` function and the Biopython `Bio.Seq` module
790 (<http://biopython.org/wiki/Biopython>, (Cock et al., 2009)). Only one 6-mer for each reverse
791 complement pair was kept (resulting in 2,080 6-mers). Genes were considered -Fe non-
792 responsive if they were not significantly differentially expressed ($\text{abs}(\log_2 \text{FC}) < 0.4$) during any
793 time point during the -Fe time course experiment or in any -Fe treated root zone or in four
794 additional -Fe treatment experiments ((Li and Schmidt, 2010): GSE16964, (Long et al., 2010):
795 GSE21443, (Schuler et al., 2011): GSE24348, (Sivitz et al., 2012): GSE40076). The four
796 additional data sets were downloaded in CEL format from GEO and processed as described
797 in the first section of the **Methods** part.

798 Potentially meaningful *cis*-regulatory elements for -Fe response were identified in two
799 steps, where we first looked for enriched k -mers in the promoters of -Fe responsive genes and
800 then determined how well the enriched k -mers predicted -Fe response using machine learning.
801 The code for this analysis is available on GitHub (<https://github.com/ShiuLab/MotifDiscovery>,
802 https://github.com/ShiuLab/ML_Pipeline). For the first step, the promoter sequences of the
803 genes in co-expression clusters (positive set) were searched for enriched 6-mers in
804 comparison to promoter sequences of non-responsive genes (negative set). These enriched

805 6-mers were elongated by one base and tested again for enrichment. This process was
806 repeated until no longer k -mer was more enriched than the shorter k -mer. Enrichment was
807 calculated using a one-sided FET ($p < 0.01$).

808 For the second step, to determine which sets of enriched k -mers were predictive of -Fe
809 response, we generated features based on presence or absence of each enriched k -mer and
810 used these features to build machine learning models using the Random Forest (RF) algorithm
811 (Pedregosa et al., 2011). To avoid building biased models, 50 models were generated for each
812 co-expression cluster by randomly drawing from the negative set to generate balanced (i.e.
813 size positive set equals size of negative set) input datasets. A 10-fold cross-validation
814 approach was used to train and test the models. Briefly, the balanced datasets were divided
815 randomly into ten even groups with a 1:1 ratio of positive to negative class genes. The model
816 was trained on the 1-9 folds and applied to the 10th (and successively trained on 1-8+10 and
817 applied to the 9th, etc.). This cross-validation scheme was repeated ten times. Each RF model
818 was made up of 500 decision trees, each trained on a random subset of enriched k -mers and
819 of training set genes. The final model performance is represented by the mean F1 score (i.e.
820 F-measure) across all 50 balanced models. The F1 score is the harmonic mean of precision
821 ($P = TP / (FP + TP)$) and recall ($R = TP / (FN + TP)$), where TP=true positive, FP=false positive, and
822 FN=false negative. Only co-expression clusters for which the enriched k -mers were deemed
823 as good predictors ($F1 \geq 0.7$) were used in the downstream analysis.

824 Predictive k -mers (then referred to as putative *cis*-regulatory elements, pCREs) were
825 ranked by importance. The importance score is based on the Gini Index, which is a measure
826 of node purity, where important pCREs separate positive from negative class genes well and
827 low ranked pCREs are less informative. To determine how well the models predicted specific
828 -Fe responsive genes, we calculated the percent of times each gene was correctly predicted
829 (TP) out of the 50 balanced replicates.

830 **pCRE sequence similarity**

831 To assess sequence similarities between the 173 pCREs that were frequently identified
832 (in >5%; freq-pCREs) in GS-enriched or non-enriched co-expression clusters with good model
833 performance ($F1 \geq 0.7$), sequence dissimilarity of pCRE position weight matrices (PWMs) was
834 calculated by pair-wise PCC distance and a distance matrix was generated using the TAMO
835 package (Gordon et al., 2005). Freq-pCREs were grouped by hierarchical clustering of the
836 PCC distance matrix using the UPGMA method in the R cluster package (Maechler et al.,
837 2017), and visualized in a dendrogram (**Supplemental Figure S4E**). Due to group-wise
838 averaging of PCC distances during hierarchical clustering, the algorithm produced skewed
839 PCC distances of some similar pCRE pairs. Therefore, freq-pCRE clusters were additionally
840 visualized as a network, in which freq-pCREs with PCC distance=0 (identical freq-pCREs or
841 subsets of each other) were connected with black bold edges and freq-pCREs with PCC

842 distance \leq 0.22 were connected with light gray edges (**Figure 4C**). Highly interconnected nodes
843 were arranged in groups. The network was created using the Cytoscape software (Shannon
844 et al., 2003). To show a consensus of freq-pCREs within a network group, freq-pCRE
845 sequences were aligned using ClustalX (Larkin et al., 2007) with default parameters and a
846 sequence logo was created with weblogo (<https://weblogo.berkeley.edu/logo.cgi>).

847 **Identification of most informative pCREs (min-pCREs) by non-linear regression**

848 The most informative pCREs of a co-expression cluster were defined as the minimum
849 set of pCREs (min-pCREs) needed for RF models without sacrificing performance. To identify
850 min-pCREs, for each GS-enriched co-expression cluster, the pCREs used as features were
851 step-wise reduced, with the least important pCREs deleted at each step. First, for pCREs that
852 were subsets of each other (PCC distance=0), the lower ranked one was removed. Then, from
853 this list of pCREs and for successively shorter lists of pCREs (n=40, 30, 25, 20, 15, 12, 10, 8,
854 6, 5, 4, 3, 2, 1), 10 replicates of RF models were trained on balanced datasets. F1 scores were
855 plotted against the number of pCREs (x) and a non-linear regression curve was fitted to the
856 data points. An exponential recovery function

$$857 \quad F1(x) = a(1 - e^{-nx})$$

858 was found to best describe the data behavior. Starting values for variables a and n were
859 approximated by fitting a linear model to the logarithmic transformation of the function. The set
860 of pCREs with the highest F1 closest to the inflection point of the regression curve was defined
861 as min-pCRE set (example in **Supplemental Figure S6**).

862 **pCRE similarity to TFBMs**

863 *In vitro* binding data of Arabidopsis TFs to genomic DNA (DNA Affinity Purification
864 Sequencing, DAP-seq, (O'Malley et al., 2016)) and TF binding data based on protein-binding
865 microarray data or the TRANSFAC data base (Catalog of Inferred Sequence Binding
866 Preferences, CIS-BP, (Weirauch et al., 2014)) were used (**Supplemental Table S4B**), with
867 DAP-seq TFBMs used over CIS-BP TFBMs when the TF was present in both databases.
868 PWMs of pCREs were compared to PWMs of TFBMs using PCC and the pCREs were
869 classified as similar to (A) a specific TF, (B) a TF family, or (C) to TFs generally, based on the
870 degree of similarity to their best matching TFBM (Uygun et al., 2017). A pCRE was similar to
871 a specific TFBM (A) if the PCC between the pCRE and the TFBM was \geq 95th percentile of
872 PCCs between that TFBM and TFBMs from the same TF family. Alternatively a pCRE was
873 similar to TFBMs from a TF family (B) if the PCC between the pCRE and a TFBM from that
874 family was \geq 95th percentile of PCCs between TFBMs from that family and TFBMs from other
875 TF families. Finally, a pCRE was similar to TFBMs (C) if the PCC between the pCRE and any
876 known TFBM was \geq 95th percentile of PCCs between TFBMs and randomly generated 6-mers.
877 For 95th percentile PCC thresholds see **Supplemental Table S4C**. To determine if pCREs

878 similar to specific TF families were enriched in GS-enriched versus non-enriched co-
879 expression clusters, the percentage of pCREs similar to TFBMs (significance level A or B) from
880 each TF family was calculated for each co-expression cluster category. Then, FET with
881 multiple testing correction ($q \leq 0.05$) was used to determine if GS-enriched co-expression
882 clusters were enriched for TF families compared to non-GS-enriched co-expression clusters
883 and vice versa.

884 **Positional distribution of pCREs**

885 To determine the positional distribution of the min-pCREs for each GS-enriched co-
886 expression cluster, min-pCREs were converted to PWMs adjusted to the Arabidopsis
887 background AT (0.33) and CG (0.17) content using the TAMO package (Gordon et al., 2005)
888 and mapped to the promoter sequences ranging from 1000 bp upstream to 500 bp downstream
889 of the transcription start site (1000-TSS-500), using Motility (<http://cartwheel.caltech.edu>). For
890 comparison, min-pCRE PWMs were also mapped to exons and introns, respectively, and to
891 the region 500 bp upstream and 1000 bp downstream of the transcription termination site (500-
892 TTS-1000). Arabidopsis sequences were downloaded from TAIR
893 (ftp://ftp.arabidopsis.org/Sequences/blast_datasets/TAIR10_blastsets/). Positional
894 distributions were calculated as described in (Uygun et al., 2017). In brief, min-pCREs were
895 mapped to 100 bp bins of 1000-TSS-500 and 500-TTS-1000 and to whole exons and introns.
896 For comparison, min-pCREs were mapped to randomized versions of the sequences.
897 Randomization was performed within each 100 bp bin and in each exon or intron, respectively,
898 in order to maintain nucleotide composition and therefore GC content. Positional distribution
899 was calculated as \log_2FC of number of observed mappings divided by number of randomly
900 expected mappings ($\log_2FC(\text{observed}/\text{expected})$).

901 **pCRE coordinate overlap with CNS coordinates**

902 All 615 min-pCRE PWMs were mapped to the putative promoter region (1 kb upstream
903 of TSS) of -Fe response genes (described above). The min-pCRE coordinates were then
904 compared to the coordinates reported as conserved non-coding sequences (CNS) across nine
905 species in the Brassicaceae family (Haudry et al., 2013) downloaded from the UCSC
906 Genomics Bioinformatics website
907 (http://mustang.biol.mcgill.ca:8885/download/A.thaliana/gff/AT_CNS.gff; **Supplemental**
908 **Table S6B**).

909

910

911 **Supplemental Material**

912 The following supporting material is available as three supplemental PDF files (1:
913 **Supplemental_Figures_S1-S7_Supplemental_Table_S1_Supplemental_literature**; 2:
914 **Supplemental_Figure_S8**; 3: **Supplemental_Figure_S9**), and as supplemental Excel
915 spreadsheet (**Supplemental_Table_S2-S7_spreadsheet**).

916

917 **Supplemental Figure S1.** Complete GO enrichment analysis of -Fe-responsive genes.

918 **Supplemental Figure S2.** GO terms and -Fe GS gene enrichments of the defined co-
919 expression clusters containing up-/down-regulated genes, and mean similarity within
920 designated superclusters of up-regulated and up-/down-regulated genes.

921 **Supplemental Figure S3.** Co-expression cluster RF model performance of cluster category
922 (GS-enriched and non-enriched) and cluster size.

923 **Supplemental Figure S4.** Comparison of the pCRE abundance and importance in GS-
924 enriched clusters vs. non-enriched clusters and hierarchical clustering of freq-pCRE
925 sequences.

926 **Supplemental Figure S5.** Significance of sequence similarity for freq-pCRE from non-
927 enriched clusters and the best matching known TFBM.

928 **Supplemental Figure S6.** Example of a non-linear regression curve to determine the minimum
929 set of pCREs for a co-expression cluster.

930 **Supplemental Figure S7.** High-resolution image of **Figure 6**.

931 **Supplemental Figure S8.** Expression plots of all well-performing co-expression clusters.

932 **Supplemental Figure S9.** Positional distribution plots of all 615 min-pCREs.

933 **Supplemental Table S1.** Robust -Fe-responsive GS genes (FIT-dependent/FIT-
934 independent).

935 **Supplemental Table S2.** Detailed information of all generated co-expression clusters: input
936 expression data combinations, algorithm and parameters used for clustering, enrichment of
937 GS genes, FIT-dependent/FIT-independent/both genes, F1 score, gene content, all identified
938 pCREs, min-pCREs.

939 **Supplemental Table S3.** List of all pCREs (n=5,639) identified in well-performing clusters (GS-
940 enriched and non-enriched).

941 **Supplemental Table S4. A:** List of most relevant pCREs (freq-pCREs and min-pCREs) and
942 their similarity to DAP-seq and CIS-BP TFBMs. **B:** DAP-seq and CIS-BP TFBMs used in this
943 study. **C:** TF family 95th percentiles of within, between and random PCC thresholds determining
944 the pCRE-TFBM similarity.

945 **Supplemental Table S5.** Association of min-pCREs to FIT-dependent, FIT-independent or
946 both cluster categories.

947 **Supplemental Table S6. A:** List of all 615 min-pCREs with total counts of GS-enriched
948 clusters and genes having the min-pCRE. **B:** Overlap of min-pCRE coordinates with
949 Brassicaceae conserved non-coding sequences (CNS).

950 **Supplemental Table S7.** p - and q -values of complete GO enrichment analysis of -Fe-
951 responsive genes (**Supplemental Figure S1**).

952

953 **Acknowledgements**

954 We thank Sahra Uygun, Bethany Moore, Nicholas Panchy, and Peipei Wang for providing
955 scripts and help with programming. B.S. is a member of the international graduate school
956 iGRAD-Plant, Düsseldorf. Funding from the German Research Foundation through the DFG
957 International Research Training group 1525 to P.B. is greatly acknowledged. C.B.A. was
958 supported in part by an NSF Graduate Research Fellowship. This work was partly supported
959 by grants to S.-H.S. from the US National Science Foundation IOS-1546617, DEB-1655386,
960 and DGE-1828149, and the US Department of Energy (DOE) Great Lakes Bioenergy
961 Research Center (DOE Office of Science BER DE-SC0018409). This work received funding
962 from Germany's Excellence Strategy, EXC 2048/1, Project ID: 390686111.

963

964 **Figure Legends**

965 **Figure 1. -Fe pCRE identification workflow and transcriptomic data.**

966 **A:** pCRE identification workflow. **B:** Heatmap of enrichment (FET, $q < 0.05$) of selected GO
967 terms in genes that were significantly up- (red) or down-regulated (blue) ($q < 0.05$) at ≥ 1 of 6
968 time points in -Fe-treated roots of 6 d-old seedlings (Dinneny et al., 2008). Differential
969 regulation was defined as \log_2 fold-change (\log_2FC) > 1 or < -1 (treatment vs. control). GOs are
970 sorted by category, and expression patterns of genes corresponding to Fe-related GOs are
971 shown below the heatmap. Yellow genes indicate -Fe GS genes. **C:** Transcriptomic data
972 combinations which were used for clustering of co-expressed genes. Gray filled boxes in
973 columns depict (**top**) expression data used in the combination and (**bottom**) if up (up-regulated
974 only) or up & down (up- and down-regulated) genes were included. ¹(Dinneny et al., 2008),
975 ²(Kilian et al., 2007), ³(Goda et al., 2008), ⁴(Schmid et al., 2005), *Tested with (5a) and without
976 (5b) genotoxic stress data, (5b) input only up-regulated genes.

977

978 **Figure 2. Characterization of the defined co-expression clusters by GO terms, -Fe GS** 979 **gene content and FIT-dependent/FIT-independent gene content.**

980 **A:** Heatmap of GO similarity between co-expression clusters from k -means clustering (**top**,
981 $n=985$ clusters) and GS gene correlation (PCC; **bottom**, $n=238$), containing up-regulated

982 genes (**Figure 1C**; up- and down-regulated genes: **Supplemental Figure S2A**). Clusters were
983 grouped by hierarchical clustering and superclusters (A-F) were defined as groups of >20
984 clusters that have a within-mean Jaccard Index significantly higher than the mean Jaccard
985 Index of all clusters. Enriched GO terms shared by $\geq 75\%$ (*k*-means) and $\geq 90\%$ (PCC) of the
986 clusters in each supercluster are shown (**left**). Co-expression clusters enriched for -Fe GS
987 genes are designated (yellow, **right**). **B**: Proportions of all *k*-means (**top left**) and PCC (**top**
988 **right**) co-expression clusters in which -Fe GS genes are significantly over-represented
989 (yellow). Of those (**bottom**), the proportion enriched for FIT-dependent genes (FIT, blue), FIT-
990 independent genes (non-FIT, red) or for both (mixed, gray) was calculated. **C**: Proportion of
991 FIT, non-FIT, and mixed clusters found using each expression data combination (as in **Figure**
992 **1C**). 1: -Fe time course, 2: time course + root zones, 3: time course + abiotic stresses, 4: time
993 course + hormone treatments, 5a: time course + abiotic stresses + hormones, 5b: as 5a,
994 genotoxic stress deleted, 6: time course + abiotic stresses + developmental data. *PCC
995 clusters only. All enrichment analyses: FET, $q < 0.05$.

996

997 **Figure 3. Performance of -Fe response RF prediction models.**

998 **A**: F1 scores of all GS-enriched clusters (n=495). **Inset**: Proportions of well-performing
999 ($F1 \geq 0.7$) and poorly performing clusters among the GS-enriched clusters. **B**: F1 scores of all
1000 non-enriched clusters (n=1,240). **Inset**: Proportions of well-performing and poorly performing
1001 clusters among the non-enriched clusters. **C**: Mean F1 score distributions of all GS-enriched
1002 clusters (yellow) and non-enriched clusters (gray). Statistical analysis: Mann-Whitney U (****
1003 $p < 2.358e-09$). **D-F**: Example GS-enriched co-expression clusters with good (**D**, **E**; cluster IDs:
1004 493, 823) and bad (**F**; cluster ID: 1297) model performance. (**Left**) Expression (\log_2 fold-
1005 change: \log_2FC) profile of all genes in the co-expression cluster. (**Center**) Percent of times
1006 across RF replicates each gene was correctly predicted as -Fe-responsive (true positive (TP);
1007 black=100%, white=0%). (**Right**) pCREs sorted by importance rank (top ranked pCRE on the
1008 left) with heatmap designating when pCRE was present (gray) or absent (white) in a gene's
1009 promoter. T: -Fe treatment time course. R: -Fe-treated root zones 1-4. F1 score: harmonic
1010 mean of precision and recall, with 1=perfect prediction and 0.5=random guessing. Cluster IDs
1011 and details: **Supplemental Table S2**.

1012

1013 **Figure 4. Analysis of pCREs predictive of -Fe co-expression clusters.**

1014 **A**: Proportion of GS-enriched (yellow) and non-enriched (gray) clusters in which each pCRE
1015 (total n=5,639) was identified (y-axis) and mean importance rank (1=most important) of that
1016 pCRE in those clusters (x-axis). **Inset**: Numbers of unique and shared pCREs of GS-enriched
1017 and non-enriched cluster categories. **Upper**: all 5,639 pCREs. **Lower**: pCREs identified in >5%
1018 of GS-enriched or non-enriched clusters (n=173; freq-pCREs). **B**: Frequency of normalized

1019 mean importance ranks across all pCREs in GS-enriched (yellow) and non-enriched (gray)
1020 clusters. **C:** Cytoscape network of the 173 freq-pCREs based on sequence similarity, where
1021 similar pCREs (nodes) are connected by edges representing pair-wise correlation (PCC)
1022 distance of freq-pCRE PWMs. Bold black edges: distance=0. Light gray edges: distance
1023 ≤ 0.22 . Highly interconnected freq-pCREs were arranged in groups and numbered.
1024 Hierarchical clustering representation of PCC distances: **Supplemental Figure S5E**. Yellow
1025 filled: freq-pCRE unique for GS-enriched clusters, gray filled: freq-pCRE unique for non-
1026 enriched clusters, not filled: shared freq-pCRE. **D:** PWMs of merged freq-pCREs from the
1027 same group (as in 4C).

1028

1029 **Figure 5. Similarity of freq-pCREs to *in vitro* TFBMs.**

1030 **A:** Significance of sequence similarity for freq-pCREs from GS-enriched clusters and the best
1031 matching known TFBM. Bars represent 95th percentile (PCC) significance thresholds for within
1032 TF family (red, pCRE sequence is more similar to a specific TFBM than other TFBMs from the
1033 same family), between TF families (light blue, pCRE sequence is more similar to a TFBM in a
1034 TF family than TFBMs from other TF families), or random (dark blue, pCRE sequence is more
1035 similar to a TFBM from a family than random 6-mers). Similarity of freq-pCREs from non-
1036 enriched clusters to TFBMs: **Supplemental Figure S5**. **B:** Proportion of TF family TFBMs
1037 (representing freq-pCRE matches meeting at least “between” threshold) in GS-enriched
1038 clusters (x-axis) and non-enriched clusters (y-axis). TFBM matches significantly over-
1039 represented (FET, $q < 0.05$) in the GS-enriched or non-enriched cluster category are depicted
1040 in red and marked with “X”. Dashed line marks theoretical position for TF family TFBMs with
1041 the same proportion in both categories.

1042

1043 **Figure 6. Characteristics of the most informative pCREs (min-pCREs).**

1044 **A:** Proportion of GS-enriched co-expression clusters enriched for FIT-dependent genes (FIT,
1045 blue), FIT-independent genes (non-FIT, red) or both (mixed, gray). **B:** Ternary plot including
1046 min-pCREs identified in >3% (n=5) GS-enriched clusters. Position of the min-pCREs
1047 corresponds to the normalized proportions of FIT, non-FIT, and mixed clusters in which the
1048 min-pCRE was identified. Bubble size corresponds to the overall proportion of GS-enriched
1049 clusters with min-pCRE. Labeled min-pCREs are shown in 6C, D or mentioned in the main
1050 text. **C:** Positional bias of all (mean with standard deviation; **top**) and selected min-pCREs
1051 (**below**) in the putative promoter region (**1st column**), all introns (In) and all exons (Ex) (**2nd**
1052 **column**; mean with standard deviation), and in the putative non-coding region (**3rd column**).
1053 1st and 3rd column: position distributions in all co-expression clusters with min-pCRE (gray
1054 areas) with mean distribution (red line). TFBM matches (PCC) for each min-pCRE are shown
1055 (**4th column**) and min-pCREs are sorted by TF family. $\log_2(\text{obs/exp})$: \log_2 of the number of

1056 observed (obs) min-pCRE occurrences divided by the number of min-pCRE occurrences in
1057 randomized sequences (expected, exp). **D**: Genes which might be regulated by the selected
1058 min-pCREs. Count: number of GS-enriched clusters in which the min-pCRE was identified and
1059 which included the respective gene having the min-pCRE in its promoter. Dashed line: total
1060 number of GS-enriched clusters with the min-pCRE. Genes in which the min-pCRE overlaps
1061 with a CNS are designated with black bars. A high-resolution image is available as
1062 **Supplemental Figure S7**.

1063

1064 **Literature Cited**

1065 Ali MA, Azeem F, Nawaz MA, Acet T, Abbas A, Imran QM, Shah KH, Rehman HM, Chung G,
1066 Yang SH, Bohlmann H (2018) Transcription factors WRKY11 and WRKY17 are
1067 involved in abiotic stress responses in Arabidopsis. *J Plant Physiol* 226: 12-21

1068 Andriankaja ME, Danisman S, Mignolet-Spruyt LF, Claeys H, Kochanke I, Vermeersch M, De
1069 Milde L, De Bodt S, Storme V, Skiryicz A, Maurer F, Bauer P, Mühlenbock P, Van
1070 Breusegem F, Angenent GC, Immink RG, Inzé D (2014) Transcriptional coordination
1071 between leaf cell differentiation and chloroplast development established by TCP20
1072 and the subgroup Ib bHLH transcription factors. *Plant Mol Biol* 85: 233-245

1073 Assunção AG, Herrero E, Lin YF, Huettel B, Talukdar S, Smaczniak C, Immink RG, van Eldik
1074 M, Fiers M, Schat H, Aarts MG (2010) Arabidopsis thaliana transcription factors bZIP19
1075 and bZIP23 regulate the adaptation to zinc deficiency. *Proc Natl Acad Sci U S A* 107:
1076 10296-10301

1077 Bailey PC, Martin C, Toledo-Ortiz G, Quail PH, Huq E, Heim MA, Jakoby M, Werber M,
1078 Weisshaar B (2003) Update on the basic helix-loop-helix transcription factor gene
1079 family in Arabidopsis thaliana. *Plant Cell* 15: 2497-2502

1080 Balk J, Schaedler TA (2014) Iron cofactor assembly in plants. *Annu Rev Plant Biol* 65: 125-
1081 153

1082 Bemer M, van Dijk ADJ, Immink RGH, Angenent GC (2017) Cross-family transcription factor
1083 interactions: an additional layer of gene regulation. *Trends Plant Sci* 22: 66-80

1084 Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and powerful
1085 approach to multiple testing. *Journal of the Royal Statistical Society. Series B*
1086 (Methodological) 57: 12

- 1087 Berardini TZ, Mundodi S, Reiser L, Huala E, Garcia-Hernandez M, Zhang P, Mueller LA, Yoon
1088 J, Doyle A, Lander G, Moseyko N, Yoo D, Xu I, Zoeckler B, Montoya M, Miller N,
1089 Weems D, Rhee SY (2004) Functional annotation of the Arabidopsis genome using
1090 controlled vocabularies. *Plant Physiol* 135: 745-755
- 1091 Blaby-Haas C, Merchant S (2013) Iron sparing and recycling in a compartmentalized cell.
1092 *Current Opinion in Microbiology* 16: 677-685
- 1093 Boulard C, Fatihi A, Lepiniec L, Dubreucq B (2017) Regulation and evolution of the interaction
1094 of the seed B3 transcription factors with NF-Y subunits. *Biochim Biophys Acta Gene
1095 Regul Mech* 1860: 1069-1078
- 1096 Briat JF, Rouached H, Tissot N, Gaymard F, Dubos C (2015) Integration of P, S, Fe, and Zn
1097 nutrition signals in *Arabidopsis thaliana*: potential involvement of PHOSPHATE
1098 STARVATION RESPONSE 1 (PHR1). *Frontiers in Plant Science* 6
- 1099 Brumbarova T, Bauer P, Ivanov R (2015) Molecular mechanisms governing Arabidopsis iron
1100 uptake. *Trends Plant Sci* 20: 124-133
- 1101 Choi HS, Seo M, Cho HT (2018) Two TPL-binding motifs of ARF2 are involved in repression
1102 of auxin responses. *Front Plant Sci* 9: 372
- 1103 Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T,
1104 Kauff F, Wilczynski B, de Hoon MJL (2009) Biopython: freely available Python tools for
1105 computational molecular biology and bioinformatics. *Bioinformatics* 25: 1422-1423
- 1106 Colangelo EP, Guerinot ML (2004) The essential basic helix-loop-helix protein FIT1 is required
1107 for the iron deficiency response. *Plant Cell* 16: 3400-3412
- 1108 Curie C, Mari S (2017) New routes for plant iron mining. *New Phytologist* 214: 521-525
- 1109 Davière JM, Wild M, Regnault T, Baumberger N, Eisler H, Genschik P, Achard P (2014) Class
1110 I TCP-DELLA interactions in inflorescence shoot apex determine plant height. *Curr Biol*
1111 24: 1923-1928
- 1112 Deng W, Zhang K, Busov V, Wei H (2017) Recursive random forest algorithm for constructing
1113 multilayered hierarchical gene regulatory networks that govern biological pathways.
1114 *PLoS One* 12: e0171532

- 1115 Dinneny JR, Long TA, Wang JY, Jung JW, Mace D, Pointer S, Barron C, Brady SM,
1116 Schiefelbein J, Benfey PN (2008) Cell identity mediates the response of Arabidopsis
1117 roots to abiotic stress. *Science* 320: 942-945
- 1118 Fourcroy P, Sisó-Terraza P, Sudre D, Savirón M, Reyt G, Gaymard F, Abadía A, Abadia J,
1119 Alvarez-Fernández A, Briat JF (2014) Involvement of the ABCG37 transporter in
1120 secretion of scopoletin and derivatives by Arabidopsis roots in response to iron
1121 deficiency. *New Phytol* 201: 155-167
- 1122 Garcia MJ, Romera FJ, Stacey MG, Stacey G, Villar E, Alcantara E, Perez-Vicente R (2013)
1123 Shoot to root communication is necessary to control the expression of iron-acquisition
1124 genes in Strategy I plants. *Planta* 237: 65-75
- 1125 Gautier L, Cope L, Bolstad BM, Irizarry RA (2004) affy--analysis of Affymetrix GeneChip data
1126 at the probe level. *Bioinformatics* 20: 307-315
- 1127 Goda H, Sasaki E, Akiyama K, Maruyama-Nakashita A, Nakabayashi K, Li W, Ogawa M,
1128 Yamauchi Y, Preston J, Aoki K, Kiba T, Takatsuto S, Fujioka S, Asami T, Nakano T,
1129 Kato H, Mizuno T, Sakakibara H, Yamaguchi S, Nambara E, Kamiya Y, Takahashi H,
1130 Hirai MY, Sakurai T, Shinozaki K, Saito K, Yoshida S, Shimada Y (2008) The
1131 AtGenExpress hormone and chemical treatment data set: experimental design, data
1132 evaluation, model data analysis and data access. *Plant J* 55: 526-542
- 1133 Gollhofer J, Schläwicke C, Jungnick N, Schmidt W, Buckhout TJ (2011) Members of a small
1134 family of nodulin-like genes are regulated under iron deficiency in roots of Arabidopsis
1135 thaliana. *Plant Physiol Biochem* 49: 557-564
- 1136 Gollhofer J, Timofeev R, Lan P, Schmidt W, Buckhout TJ (2014) Vacuolar-Iron-Transporter1-
1137 Like proteins mediate iron homeostasis in Arabidopsis. *PLoS One* 9: e110468
- 1138 Gordon DB, Nekludova L, McCallum S, Fraenkel E (2005) TAMO: a flexible, object-oriented
1139 framework for analyzing transcriptional regulation using DNA-sequence motifs.
1140 *Bioinformatics* 21: 3164-3165
- 1141 Gratz R, Manishankar P, Ivanov R, Köster P, Mohr I, Trofimov K, Steinhorst L, Meiser J, Mai
1142 HJ, Drerup M, Arendt S, Holtkamp M, Karst U, Kudla J, Bauer P, Brumbarova T (2019)
1143 CIPK11-dependent phosphorylation modulates FIT transcription factor activity,
1144 promoting Arabidopsis iron acquisition in response to calcium signaling. *Developmental*
1145 *Cell* 48

- 1146 Hantzis LJ, Kroh GE, Jahn CE, Cantrell M, Peers G, Pilon M, Ravet K (2018) A program for
1147 iron economy during deficiency targets specific Fe proteins. *Plant Physiol* 176: 596-
1148 610
- 1149 Hartigan JA, Wong MA (1979) A k-means clustering algorithm. *Journal of the Royal Statistical*
1150 *Society. Series C (Applied Statistics)* 28: 9
- 1151 Haudry A, Platts AE, Vello E, Hoen DR, Leclercq M, Williamson RJ, Forczek E, Joly-Lopez Z,
1152 Steffen JG, Hazzouri KM, Dewar K, Stinchcombe JR, Schoen DJ, Wang XW, Schmutz
1153 J, Town CD, Edger PP, Pires JC, Schumaker KS, Jarvis DE, Mandáková T, Lysak MA,
1154 van den Bergh E, Schranz ME, Harrison PM, Moses AM, Bureau TE, Wright SI,
1155 Blanchette M (2013) An atlas of over 90,000 conserved noncoding sequences provides
1156 insight into crucifer regulatory regions. *Nature Genetics* 45: 891-U228
- 1157 Haydon MJ, Kawachi M, Wirtz M, Hillmer S, Hell R, Krämer U (2012) Vacuolar nicotianamine
1158 has critical and distinct roles under iron deficiency and for zinc sequestration in
1159 *Arabidopsis*. *Plant Cell* 24: 724-737
- 1160 Heyndrickx KS, Van de Velde J, Wang CM, Weigei D, Vandepoele K (2014) A functional and
1161 evolutionary perspective on transcription factor binding in *Arabidopsis thaliana*. *Plant*
1162 *Cell* 26: 3894-3910
- 1163 Hong S, Kim SA, Guerinot ML, McClung CR (2013) Reciprocal interaction of the circadian
1164 clock with the iron homeostasis network in *Arabidopsis*. *Plant Physiology* 161: 893-903
- 1165 Hsieh WP, Hsieh HL, Wu SH (2012) *Arabidopsis* bZIP16 transcription factor integrates light
1166 and hormone signaling pathways to regulate early seedling development. *Plant Cell*
1167 24: 3997-4011
- 1168 Ito H, Gray WM (2006) A gain-of-function mutation in the *Arabidopsis* pleiotropic drug
1169 resistance transporter PDR9 confers resistance to auxinic herbicides. *Plant Physiol*
1170 142: 63-74
- 1171 Ivanov R, Brumbarova T, Bauer P (2012) Fitting into the harsh reality: regulation of iron-
1172 deficiency responses in dicotyledonous plants. *Mol Plant* 5: 27-42
- 1173 Jakoby M, Wang HY, Reidt W, Weisshaar B, Bauer P (2004) FRU (BHLH029) is required for
1174 induction of iron mobilization genes in *Arabidopsis thaliana*. *FEBS Lett* 577: 528-534
- 1175 Jeong J, Merkovich A, Clyne M, Connolly EL (2017) Directing iron transport in dicots: regulation
1176 of iron acquisition and translocation. *Curr Opin Plant Biol* 39: 106-113

- 1177 Jones E, Oliphant T, Peterson P, al. e (2001) SciPy: Open source scientific tools for Python.
1178 In,
- 1179 Kai K, Mizutani M, Kawamura N, Yamamoto R, Tamai M, Yamaguchi H, Sakata K, Shimizu B
1180 (2008) Scopoletin is biosynthesized via ortho-hydroxylation of feruloyl CoA by a 2-
1181 oxoglutarate-dependent dioxygenase in Arabidopsis thaliana. Plant J 55: 989-999
- 1182 Kakei Y, Ogo Y, Itai RN, Kobayashi T, Yamakawa T, Nakanishi H, Nishizawa NK (2013)
1183 Development of a novel prediction method of cis-elements to hypothesize collaborative
1184 functions of cis-element pairs in iron-deficient rice. Rice (N Y) 6: 22
- 1185 Khan MA, Castro-Guerrero NA, McInturf SA, Nguyen NT, Dame AN, Wang J, Bindbeutel RK,
1186 Joshi T, Jurisson SS, Nusinow DA, Mendoza-Cozatl DG (2018) Changes in iron
1187 availability in Arabidopsis are rapidly sensed in the leaf vasculature and impaired
1188 sensing leads to opposite transcriptional programs in leaves and roots. Plant Cell
1189 Environ
- 1190 Kilian J, Whitehead D, Horak J, Wanke D, Weinl S, Batistic O, D'Angelo C, Bornberg-Bauer E,
1191 Kudla J, Harter K (2007) The AtGenExpress global stress expression data set:
1192 protocols, evaluation and model data analysis of UV-B light, drought and cold stress
1193 responses. Plant J 50: 347-363
- 1194 Kobayashi T, Itai RN, Ogo Y, Kakei Y, Nakanishi H, Takahashi M, Nishizawa NK (2009) The
1195 rice transcription factor IDEF1 is essential for the early response to iron deficiency, and
1196 induces vegetative expression of late embryogenesis abundant genes. Plant J 60: 948-
1197 961
- 1198 Kobayashi T, Nakayama Y, Itai RN, Nakanishi H, Yoshihara T, Mori S, Nishizawa NK (2003)
1199 Identification of novel cis-acting elements, IDE1 and IDE2, of the barley IDS2 gene
1200 promoter conferring iron-deficiency-inducible, root-specific expression in
1201 heterogeneous tobacco plants. Plant J 36: 780-793
- 1202 Kobayashi T, Ogo Y, Aung MS, Nozoye T, Itai RN, Nakanishi H, Yamakawa T, Nishizawa NK
1203 (2010) The spatial expression and regulation of transcription factors IDEF1 and IDEF2.
1204 Ann Bot 105: 1109-1117
- 1205 Kobayashi T, Ogo Y, Itai RN, Nakanishi H, Takahashi M, Mori S, Nishizawa NK (2007) The
1206 transcription factor IDEF1 regulates the response to and tolerance of iron deficiency in
1207 plants. Proc Natl Acad Sci U S A 104: 19150-19155

- 1208 Kobayashi T, Suzuki M, Inoue H, Itai RN, Takahashi M, Nakanishi H, Mori S, Nishizawa NK
1209 (2005) Expression of iron-acquisition-related genes in iron-deficient rice is co-ordinately
1210 induced by partially conserved iron-deficiency-responsive elements. *J Exp Bot* 56:
1211 1305-1316
- 1212 Koryachko A, Matthiadis A, Muhammad D, Foret J, Brady SM, Ducoste JJ, Tuck J, Long TA,
1213 Williams C (2015) Clustering and differential alignment algorithm: identification of early
1214 stage regulators in the *Arabidopsis thaliana* iron deficiency response. *PLoS One* 10:
1215 e0136591
- 1216 Lanquar V, Lelièvre F, Bolte S, Hamès C, Alcon C, Neumann D, Vansuyt G, Curie C, Schröder
1217 A, Krämer U, Barbier-Brygoo H, Thomine S (2005) Mobilization of vacuolar iron by
1218 AtNRAMP3 and AtNRAMP4 is essential for seed germination on low iron. *EMBO J* 24:
1219 4041-4051
- 1220 Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, Valentin F,
1221 Wallace IM, Wilm A, Lopez R, Thompson JD, Gibson TJ, Higgins DG (2007) Clustal W
1222 and clustal X version 2.0. *Bioinformatics* 23: 2947-2948
- 1223 Lei GJ, Zhu XF, Wang ZW, Dong F, Dong NY, Zheng SJ (2014) Abscisic acid alleviates iron
1224 deficiency by promoting root iron reutilization and transport from root to shoot in
1225 *Arabidopsis*. *Plant Cell Environ* 37: 852-863
- 1226 Li WF, Lan P (2017) The understanding of the plant iron deficiency responses in Strategy I
1227 plants and the role of ethylene in this process by omic approaches. *Frontiers in Plant*
1228 *Science* 8
- 1229 Li WF, Schmidt W (2010) A lysine-63-linked ubiquitin chain-forming conjugase, UBC13,
1230 promotes the developmental responses to iron deficiency in *Arabidopsis* roots. *Plant*
1231 *Journal* 62: 330-343
- 1232 Li X, Zhang H, Ai Q, Liang G, Yu D (2016) Two bHLH transcription factors, bHLH34 and
1233 bHLH104, regulate iron homeostasis in *Arabidopsis thaliana*. *Plant Physiol* 170: 2478-
1234 2493
- 1235 Liang G, Zhang H, Li X, Ai Q, Yu D (2017) bHLH transcription factor bHLH115 regulates iron
1236 homeostasis in *Arabidopsis thaliana*. *J Exp Bot* 68: 1743-1755
- 1237 Lingam S, Mohrbacher J, Brumbarova T, Potuschak T, Fink-Straube C, Blondet E, Genschik
1238 P, Bauer P (2011) Interaction between the bHLH transcription factor FIT and

- 1239 ETHYLENE INSENSITIVE3/ETHYLENE INSENSITIVE3-LIKE1 reveals molecular
1240 linkage between the regulation of iron acquisition and ethylene signaling in Arabidopsis.
1241 Plant Cell 23: 1815-1829
- 1242 Liu Y, Xie Y, Wang H, Ma X, Yao W (2017) Light and ethylene coordinately regulate the
1243 phosphate starvation response through transcriptional regulation of PHOSPHATE
1244 STARVATION RESPONSE1. Plant Cell 29: 2269-2284
- 1245 Long TA, Tsukagoshi H, Busch W, Lahner B, Salt DE, Benfey PN (2010) The bHLH
1246 transcription factor POPEYE regulates response to iron deficiency in Arabidopsis roots.
1247 Plant Cell 22: 2219-2236
- 1248 Ma C, Zhang HH, Wang X (2014) Machine learning for Big Data analytics in plants. Trends
1249 Plant Sci 19: 798-808
- 1250 Maechler M, Rousseeuw P, Struyf A, Hubert M, Hornik K (2017) cluster: Cluster Analysis
1251 Basics and Extensions. In, Ed R package version 2.0.6
- 1252 Mai HJ, Pateyron S, Bauer P (2016) Iron homeostasis in Arabidopsis thaliana: transcriptomic
1253 analyses reveal novel FIT-regulated genes, iron deficiency marker genes and
1254 functional gene networks. BMC Plant Biol 16: 211
- 1255 Mangano S, Denita-Juarez SP, Choi HS, Marzol E, Hwang Y, Ranocha P, Velasquez SM,
1256 Borassi C, Barberini ML, Aptekmann AA, Muschietti JP, Nadra AD, Dunand C, Cho HT,
1257 Estevez JM (2017) Molecular link between auxin and ROS-mediated polar growth. Proc
1258 Natl Acad Sci U S A 114: 5289-5294
- 1259 Marschner H, Römheld V (1994) Strategies of plants for acquisition of iron. Plant and Soil 165:
1260 261-274
- 1261 Matioli CC, Tomaz JP, Duarte GT, Prado FM, Del Bem LE, Silveira AB, Gauer L, Corrêa LG,
1262 Drumond RD, Viana AJ, Di Mascio P, Meyer C, Vincentz M (2011) The Arabidopsis
1263 bZIP gene AtbZIP63 is a sensitive integrator of transient abscisic acid and glucose
1264 signals. Plant Physiol 157: 692-705
- 1265 Meiser J, Lingam S, Bauer P (2011) Posttranslational regulation of the iron deficiency basic
1266 helix-loop-helix transcription factor FIT is affected by iron and nitric oxide. Plant Physiol
1267 157: 2154-2166
- 1268 Mendoza-Cózatl DG, Xie Q, Akmajian GZ, Jobe TO, Patel A, Stacey MG, Song L, Demoin
1269 DW, Jurisson SS, Stacey G, Schroeder JI (2014) OPT3 is a component of the iron-

- 1270 signaling network between leaves and roots and misregulation of OPT3 leads to an
1271 over-accumulation of cadmium in seeds. *Mol Plant* 7: 1455-1469
- 1272 Murgia I, Tarantino D, Soave C, Morandini P (2011) Arabidopsis CYP82C4 expression is
1273 dependent on Fe availability and circadian rhythm, and correlates with genes involved
1274 in the early Fe deficiency response. *J Plant Physiol* 168: 894-902
- 1275 Müller M, Schmidt W (2004) Environmentally induced plasticity of root hair development in
1276 Arabidopsis. *Plant Physiol* 134: 409-419
- 1277 Nicolas M, Cubas P (2016) TCP factors: new kids on the signaling block. *Curr Opin Plant Biol*
1278 33: 33-41
- 1279 O'Malley RC, Huang SSC, Song L, Lewsey MG, Bartlett A, Nery JR, Galli M, Gallavotti A,
1280 Ecker JR (2016) Cistrome and epicistrome features shape the regulatory DNA
1281 landscape. *Cell* 165: 1280-1292
- 1282 Ogo Y, Itai RN, Nakanishi H, Kobayashi T, Takahashi M, Mori S, Nishizawa NK (2007) The
1283 rice bHLH protein OsIRO2 is an essential regulator of the genes involved in Fe uptake
1284 under Fe-deficient conditions. *Plant J* 51: 366-377
- 1285 Ogo Y, Kobayashi T, Nakanishi Itai R, Nakanishi H, Kakei Y, Takahashi M, Toki S, Mori S,
1286 Nishizawa NK (2008) A novel NAC transcription factor, IDEF2, that recognizes the iron
1287 deficiency-responsive element 2 regulates the genes involved in iron homeostasis in
1288 plants. *J Biol Chem* 283: 13407-13417
- 1289 Palmer CM, Hindt MN, Schmidt H, Clemens S, Guerinot ML (2013) MYB10 and MYB72 are
1290 required for growth under iron-limiting conditions. *PLoS Genet* 9: e1003953
- 1291 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer
1292 P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M,
1293 Duchesnay E (2011) Scikit-learn: Machine Learning in Python. *Journal of Machine*
1294 *Learning Research* 12: 2825-2830
- 1295 Phipson B, Lee S, Majewski IJ, Alexander WS, Smyth GK (2016) Robust hyperparameter
1296 estimation protects against hypervariable genes and improves power to detect
1297 differential expression. *Ann Appl Stat* 10: 946-963
- 1298 Pitts RJ, Cernac A, Estelle M (1998) Auxin and ethylene promote root hair elongation in
1299 Arabidopsis. *Plant J* 16: 553-560

- 1300 Rajniak J, Giehl RFH, Chang E, Murgia I, von Wirén N, Sattely ES (2018) Biosynthesis of
1301 redox-active metabolites in response to iron deficiency in plants. *Nat Chem Biol* 14:
1302 442-450
- 1303 Resentini F, Felipe-Benavent A, Colombo L, Blázquez MA, Alabadí D, Masiero S (2015)
1304 TCP14 and TCP15 mediate the promotion of seed germination by gibberellins in
1305 *Arabidopsis thaliana*. *Mol Plant* 8: 482-485
- 1306 Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK (2015) limma powers
1307 differential expression analyses for RNA-sequencing and microarray studies. *Nucleic*
1308 *Acids Res* 43: e47
- 1309 Rodríguez-Celma J, Pan IC, Li W, Lan P, Buckhout TJ, Schmidt W (2013) The transcriptional
1310 response of *Arabidopsis* leaves to Fe deficiency. *Front Plant Sci* 4: 276
- 1311 Roschzttardtz H, Fuentes I, Vásquez M, Corvalán C, León G, Gómez I, Araya A, Holuigue L,
1312 Vicente-Carbajosa J, Jordana X (2009) A nuclear gene encoding the iron-sulfur subunit
1313 of mitochondrial complex II is regulated by B3 domain transcription factors during seed
1314 development in *Arabidopsis*. *Plant Physiol* 150: 84-95
- 1315 Rose AB, Carter A, Korf I, Kojima N (2016) Intron sequences that stimulate gene expression
1316 in *Arabidopsis*. *Plant Mol Biol* 92: 337-346
- 1317 Rose AB, Elfersi T, Parra G, Korf I (2008) Promoter-proximal introns in *Arabidopsis thaliana*
1318 are enriched in dispersed signals that elevate gene expression. *Plant Cell* 20: 543-551
- 1319 Ruiz-Sola M, Rodríguez-Concepción M (2012) Carotenoid biosynthesis in *Arabidopsis*: a
1320 colorful pathway. *Arabidopsis Book* 10: e0158
- 1321 Salomé PA, Oliva M, Weigel D, Krämer U (2013) Circadian clock adjustment to plant iron status
1322 depends on chloroplast and phytochrome function. *Embo Journal* 32: 511-523
- 1323 Santi S, Schmidt W (2009) Dissecting iron deficiency-induced proton extrusion in *Arabidopsis*
1324 roots. *New Phytol* 183: 1072-1084
- 1325 Schmid M, Davison TS, Henz SR, Pape UJ, Demar M, Vingron M, Schölkopf B, Weigel D,
1326 Lohmann JU (2005) A gene expression map of *Arabidopsis thaliana* development. *Nat*
1327 *Genet* 37: 501-506

- 1328 Schmid NB, Giehl RF, Döll S, Mock HP, Strehmel N, Scheel D, Kong X, Hider RC, von Wirén
1329 N (2014) Feruloyl-CoA 6'-Hydroxylase1-dependent coumarins mediate iron acquisition
1330 from alkaline substrates in Arabidopsis. *Plant Physiol* 164: 160-172
- 1331 Schmidt W, Tittel J, Schikora A (2000) Role of hormones in the induction of iron deficiency
1332 responses in Arabidopsis roots. *Plant Physiol* 122: 1109-1118
- 1333 Schuler M, Keller A, Backes C, Philippar K, Lenhof HP, Bauer P (2011) Transcriptome analysis
1334 by GeneTrail revealed regulation of functional categories in response to alterations of
1335 iron homeostasis in Arabidopsis thaliana. *Bmc Plant Biology* 11
- 1336 Selote D, Samira R, Matthiadis A, Gillikin JW, Long TA (2015) Iron-binding E3 ligase mediates
1337 iron response in plants by targeting basic helix-loop-helix transcription factors. *Plant*
1338 *Physiol* 167: 273-286
- 1339 Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B,
1340 Ideker T (2003) Cytoscape: A software environment for integrated models of
1341 biomolecular interaction networks. *Genome Research* 13: 2498-2504
- 1342 Sinclair SA, Senger T, Talke IN, Cobbett CS, Haydon MJ, Krämer U (2018) Systemic
1343 upregulation of MTP2- and HMA2-mediated Zn partitioning to the shoot supplements
1344 local Zn deficiency responses. *Plant Cell* 30: 2463-2479
- 1345 Sivitz A, Grinvalds C, Barberon M, Curie C, Vert G (2011) Proteasome-mediated turnover of
1346 the transcriptional activator FIT is required for plant iron-deficiency responses. *Plant*
1347 *Journal* 66: 1044-1052
- 1348 Sivitz AB, Hermand V, Curie C, Vert G (2012) Arabidopsis bHLH100 and bHLH101 control iron
1349 homeostasis via a FIT-independent pathway. *PLoS One* 7: e44843
- 1350 Siwinska J, Siatkowska K, Olry A, Grosjean J, Hehn A, Bourgaud F, Meharg AA, Carey M,
1351 Lojkowska E, Ichnatowicz A (2018) Scopoletin 8-hydroxylase: a novel enzyme involved
1352 in coumarin biosynthesis and iron-deficiency responses in Arabidopsis. *J Exp Bot* 69:
1353 1735-1748
- 1354 Storey JD (2002) A direct approach to false discovery rates. *Journal of the Royal Statistical*
1355 *Society Series B-Statistical Methodology* 64: 479-498
- 1356 Séguéla M, Briat JF, Vert G, Curie C (2008) Cytokinins negatively regulate the root iron uptake
1357 machinery in Arabidopsis through a growth-dependent pathway. *Plant J* 55: 289-300

- 1358 Tang LP, Zhou C, Wang SS, Yuan J, Zhang XS, Su YH (2017) FUSCA3 interacting with LEAFY
1359 COTYLEDON2 controls lateral root formation through regulating YUCCA4 gene
1360 expression in *Arabidopsis thaliana*. *New Phytol* 213: 1740-1754
- 1361 Uygun S, Peng C, Lehti-Shiu MD, Last RL, Shiu SH (2016) Utility and limitations of using gene
1362 expression data to identify functional associations. *PLoS Comput Biol* 12: e1005244
- 1363 Uygun S, Seddon AE, Azodi CB, Shiu SH (2017) Predictive models of spatial transcriptional
1364 response to high salinity. *Plant Physiol* 174: 450-464
- 1365 Vert GA, Briat JF, Curie C (2003) Dual regulation of the *Arabidopsis* high-affinity root iron
1366 uptake system by local and long-distance signals. *Plant Physiology* 132: 796-804
- 1367 Wang H, Wang HY (2015) Multifaceted roles of FHY3 and FAR1 in light signaling and beyond.
1368 *Trends in Plant Science* 20: 453-461
- 1369 Wang HY, Klatt M, Jakoby M, Bäumlein H, Weisshaar B, Bauer P (2007) Iron deficiency-
1370 mediated stress regulation of four subgroup Ib BHLH genes in *Arabidopsis thaliana*.
1371 *Planta* 226: 897-908
- 1372 Wang N, Cui Y, Liu Y, Fan H, Du J, Huang Z, Yuan Y, Wu H, Ling HQ (2013) Requirement
1373 and functional redundancy of Ib subgroup bHLH proteins for iron deficiency responses
1374 and uptake in *Arabidopsis thaliana*. *Mol Plant* 6: 503-513
- 1375 Weirauch MT, Yang A, Albu M, Cote AG, Montenegro-Montero A, Drewe P, Najafabadi HS,
1376 Lambert SA, Mann I, Cook K, Zheng H, Goity A, van Bakel H, Lozano JC, Galli M,
1377 Lewsey MG, Huang EY, Mukherjee T, Chen XT, Reece-Hoyes JS, Govindarajan S,
1378 Shaulsky G, Walhout AJM, Bouget FY, Ratsch G, Larrondo LF, Ecker JR, Hughes TR
1379 (2014) Determination and inference of eukaryotic transcription factor sequence
1380 specificity. *Cell* 158: 1431-1443
- 1381 Yan JY, Li CX, Sun L, Ren JY, Li GX, Ding ZJ, Zheng SJ (2016) A WRKY transcription factor
1382 regulates Fe translocation under Fe deficiency. *Plant Physiol* 171: 2017-2027
- 1383 Yang JL, Chen WW, Chen LQ, Qin C, Jin CW, Shi YZ, Zheng SJ (2013) The 14-3-3 protein
1384 GENERAL REGULATORY FACTOR11 (GRF11) acts downstream of nitric oxide to
1385 regulate iron acquisition in *Arabidopsis thaliana*. *New Phytol* 197: 815-824
- 1386 Yu CP, Lin JJ, Li WH (2016) Positional distribution of transcription factor binding sites in
1387 *Arabidopsis thaliana*. *Scientific Reports* 6

- 1388 Yuan Y, Wu H, Wang N, Li J, Zhao W, Du J, Wang D, Ling HQ (2008) FIT interacts with
1389 AtbHLH38 and AtbHLH39 in regulating iron uptake gene expression for iron
1390 homeostasis in Arabidopsis. *Cell Res* 18: 385-397
- 1391 Yáñez-Cuna JO, Kvon EZ, Stark A (2013) Deciphering the transcriptional cis-regulatory code.
1392 *Trends in Genetics* 29: 11-22
- 1393 Zamioudis C, Hanson J, Pieterse CM (2014) β -Glucosidase BGLU42 is a MYB72-dependent
1394 key regulator of rhizobacteria-induced systemic resistance and modulates iron
1395 deficiency responses in Arabidopsis roots. *New Phytol* 204: 368-379
- 1396 Zhai Z, Gayomba SR, Jung HI, Vimalakumari NK, Piñeros M, Craft E, Rutzke MA, Danku J,
1397 Lahner B, Punshon T, Guerinot ML, Salt DE, Kochian LV, Vatamaniuk OK (2014) OPT3
1398 is a phloem-specific iron transporter that is essential for systemic iron signaling and
1399 redistribution of iron and cadmium in Arabidopsis. *Plant Cell* 26: 2249-2264
- 1400 Zhang HM, Li Y, Yao XN, Liang G, Yu DQ (2017) POSITIVE REGULATOR OF IRON
1401 HOMEOSTASIS1, OsPRI1, facilitates iron homeostasis. *Plant Physiology* 175: 543-
1402 554
- 1403 Zhang J, Liu B, Li M, Feng D, Jin H, Wang P, Liu J, Xiong F, Wang J, Wang HB (2015) The
1404 bHLH transcription factor bHLH104 interacts with IAA-LEUCINE RESISTANT3 and
1405 modulates iron homeostasis in Arabidopsis. *Plant Cell* 27: 787-805
- 1406 Zheng L, Ying Y, Wang L, Wang F, Whelan J, Shou H (2010) Identification of a novel iron
1407 regulated basic helix-loop-helix protein involved in Fe homeostasis in *Oryza sativa*.
1408 *BMC Plant Biol* 10: 166
- 1409 Zhou Y, Zhang DZ, An JX, Yin HJ, Fang S, Chu JF, Zhao YD, Li J (2018) TCP transcription
1410 factors regulate shade avoidance via directly mediating the expression of both
1411 PHYTOCHROME INTERACTING FACTORS and auxin biosynthetic genes. *Plant*
1412 *Physiology* 176: 1850-1861
- 1413 Zou C, Sun K, Mackaluso JD, Seddon AE, Jin R, Thomashow MF, Shiu SH (2011) Cis-
1414 regulatory code of stress-responsive transcription in Arabidopsis thaliana. *Proc Natl*
1415 *Acad Sci U S A* 108: 14992-14997
- 1416

Figure 1

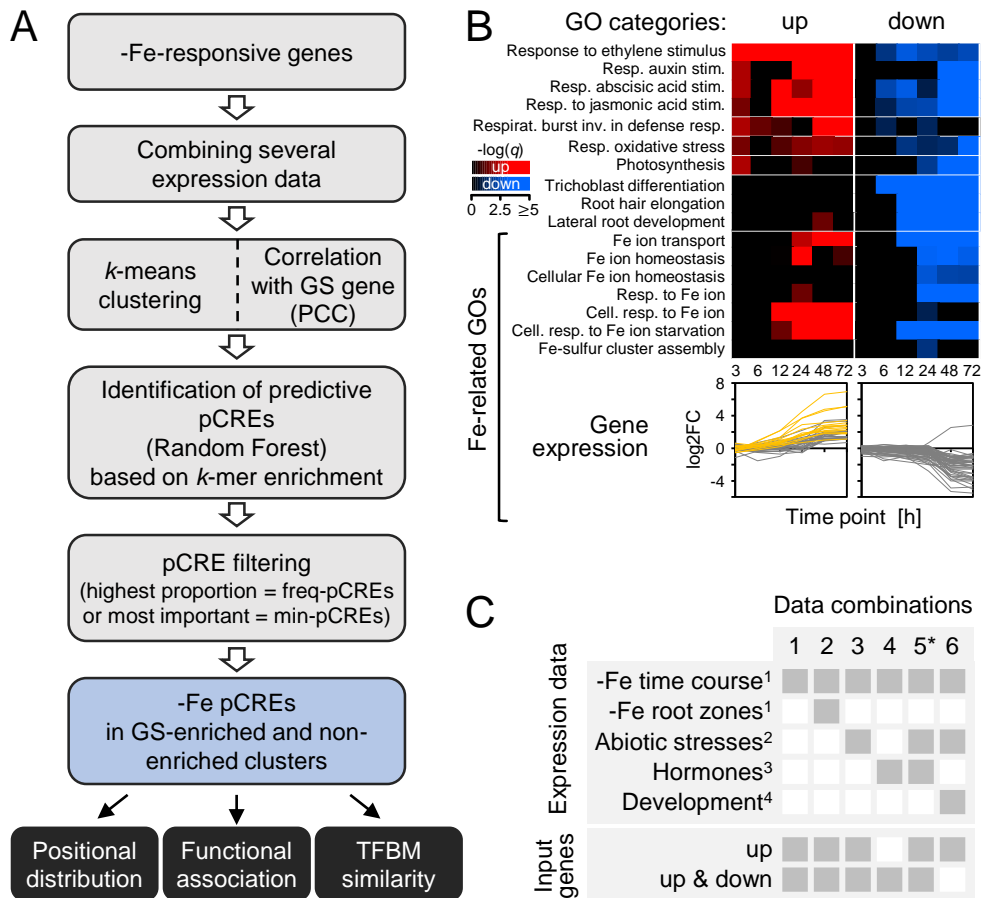


Figure 1. -Fe pCRE identification workflow and transcriptomic data.

A: pCRE identification workflow. **B:** Heatmap of enrichment (FET, $q < 0.05$) of selected GO terms in genes that were significantly up- (red) or down-regulated (blue) ($q < 0.05$) at ≥ 1 of 6 time points in -Fe-treated roots of 6 d-old seedlings (Dinneny et al., 2008). Differential regulation was defined as \log_2 fold-change (\log_2FC) > 1 or < -1 (treatment vs. control). GOs are sorted by category, and expression patterns of genes corresponding to Fe-related GOs are shown below the heatmap. Yellow genes indicate -Fe GS genes. **C:** Transcriptomic data combinations which were used for clustering of co-expressed genes. Gray filled boxes in columns depict (**top**) expression data used in the combination and (**bottom**) if up (up-regulated only) or up & down (up- and down-regulated) genes were included. ¹(Dinneny et al., 2008), ²(Kilian et al., 2007), ³(Goda et al., 2008), ⁴(Schmid et al., 2005), *Tested with (5a) and without (5b) genotoxic stress data, (5b) input only up-regulated genes.

Figure 2

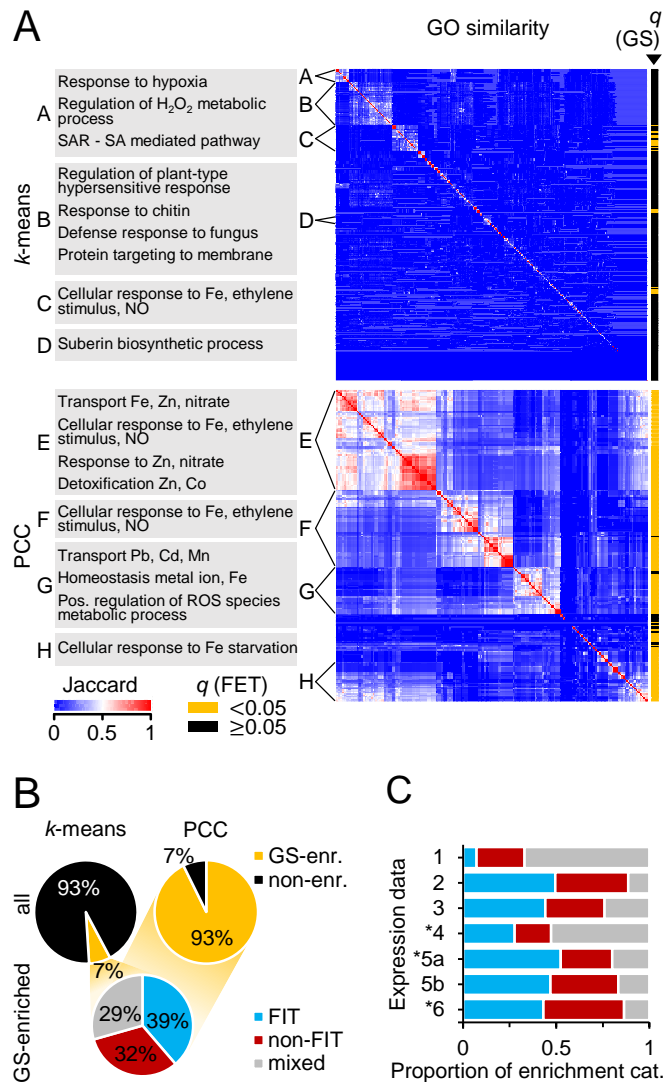


Figure 2. Characterization of the defined co-expression clusters by GO terms, -Fe GS gene content and FIT-dependent/FIT-independent gene content.

A: Heatmap of GO similarity between co-expression clusters from *k*-means clustering (**top**, $n=985$ clusters) and GS gene correlation (PCC; **bottom**, $n=238$), containing up-regulated genes (**Figure 1C**; up- and down-regulated genes: **Supplemental Figure S2A**). Clusters were grouped by hierarchical clustering and superclusters (A-F) were defined as groups of >20 clusters that have a within-mean Jaccard Index significantly higher than the mean Jaccard Index of all clusters. Enriched GO terms shared by $\geq 75\%$ (*k*-means) and $\geq 90\%$ (PCC) of the clusters in each supercluster are shown (**left**). Co-expression clusters enriched for -Fe GS genes are designated (yellow, **right**). **B:** Proportions of all *k*-means (**top left**) and PCC (**top right**) co-expression clusters in which -Fe GS genes are significantly over-represented (yellow). Of those (**bottom**), the proportion enriched for FIT-dependent genes (FIT, blue), FIT-independent genes (non-FIT, red) or for both (mixed, gray) was calculated. **C:** Proportion of enrichment categories FIT, non-FIT, and mixed clusters found using each expression data combination (as in **Figure 1C**). 1: -Fe time course, 2: time course + root zones, 3: time course + abiotic stresses, 4: time course + hormone treatments, 5a: time course + abiotic stresses + hormones, 5b: as 5a, genotoxic stress deleted, 6: time course + abiotic stresses + developmental data. *PCC clusters only. All enrichment analyses: FET, $q < 0.05$.

Figure 3

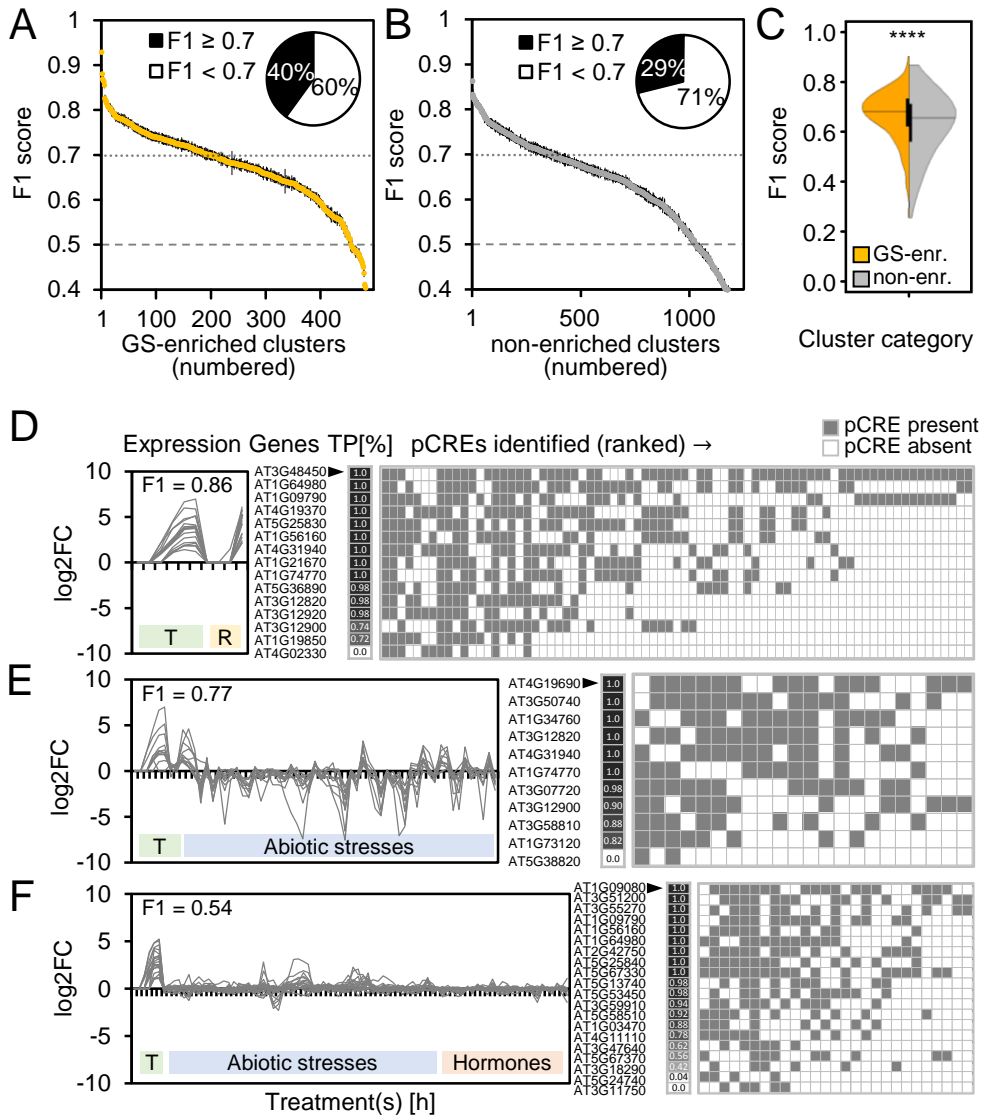


Figure 3. Performance of -Fe response RF prediction models.

A: F1 scores of all GS-enriched clusters ($n=495$). **Inset:** Proportions of well-performing ($F1 \geq 0.7$) and poorly performing clusters among the GS-enriched clusters. **B:** F1 scores of all non-enriched clusters ($n=1,240$). **Inset:** Proportions of well-performing and poorly performing clusters among the non-enriched clusters. **C:** Mean F1 score distributions of all GS-enriched clusters (yellow) and non-enriched clusters (gray). Statistical analysis: Mann-Whitney U (**** $p < 2.358e-09$). **D-F:** Example GS-enriched co-expression clusters with good (**D**, **E**; cluster IDs: 493, 823) and bad (**F**; cluster ID: 1297) model performance. **(Left)** Expression (\log_2 fold-change: \log_2FC) profile of all genes in the co-expression cluster. **(Center)** Percent of times across RF replicates each gene was correctly predicted as -Fe-responsive (true positive (TP); black=100%, white=0%). **(Right)** pCREs sorted by importance rank (top ranked pCRE on the left) with heatmap designating when pCRE was present (gray) or absent (white) in a gene's promoter. T: -Fe treatment time course. R: -Fe-treated root zones 1-4. F1 score: harmonic mean of precision and recall, with 1=perfect prediction and 0.5=random guessing. Cluster IDs and details: **Supplemental Table S2**.

Figure 4

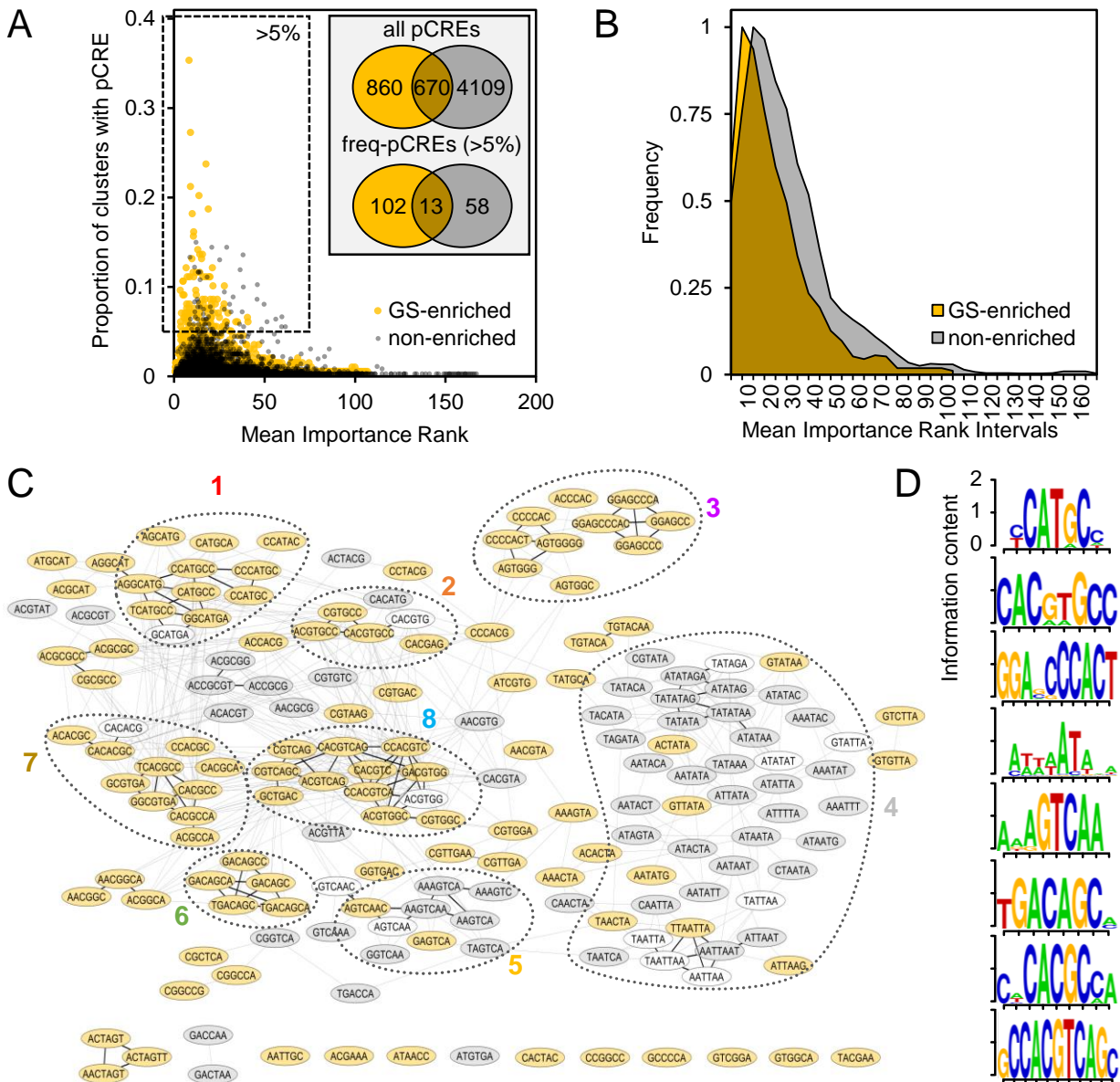


Figure 4. Analysis of pCREs predictive of -Fe co-expression clusters.

A: Proportion of GS-enriched (yellow) and non-enriched (gray) clusters in which each pCRE (total $n=5,639$) was identified (y-axis) and mean importance rank (1=most important) of that pCRE in those clusters (x-axis). **Inset:** Numbers of unique and shared pCREs of GS-enriched and non-enriched cluster categories. **Upper:** all 5,639 pCREs. **Lower:** pCREs identified in >5% of GS-enriched or non-enriched clusters ($n=173$; freq-pCREs). **B:** Frequency of normalized mean importance ranks across all pCREs in GS-enriched (yellow) and non-enriched (gray) clusters. **C:** Cytoscape network of the 173 freq-pCREs based on sequence similarity, where similar pCREs (nodes) are connected by edges representing pair-wise correlation (PCC) distance of freq-pCRE PWMs. Bold black edges: distance=0. Light gray edges: distance ≤ 0.22 . Highly interconnected freq-pCREs were arranged in groups and numbered. Hierarchical clustering representation of PCC distances: **Supplemental Figure S5E**. Yellow filled: freq-pCRE unique for GS-enriched clusters, gray filled: freq-pCRE unique for non-enriched clusters, not filled: shared freq-pCRE. **D:** PWMs of merged freq-pCREs from the same group (as in 4C).

Figure 5

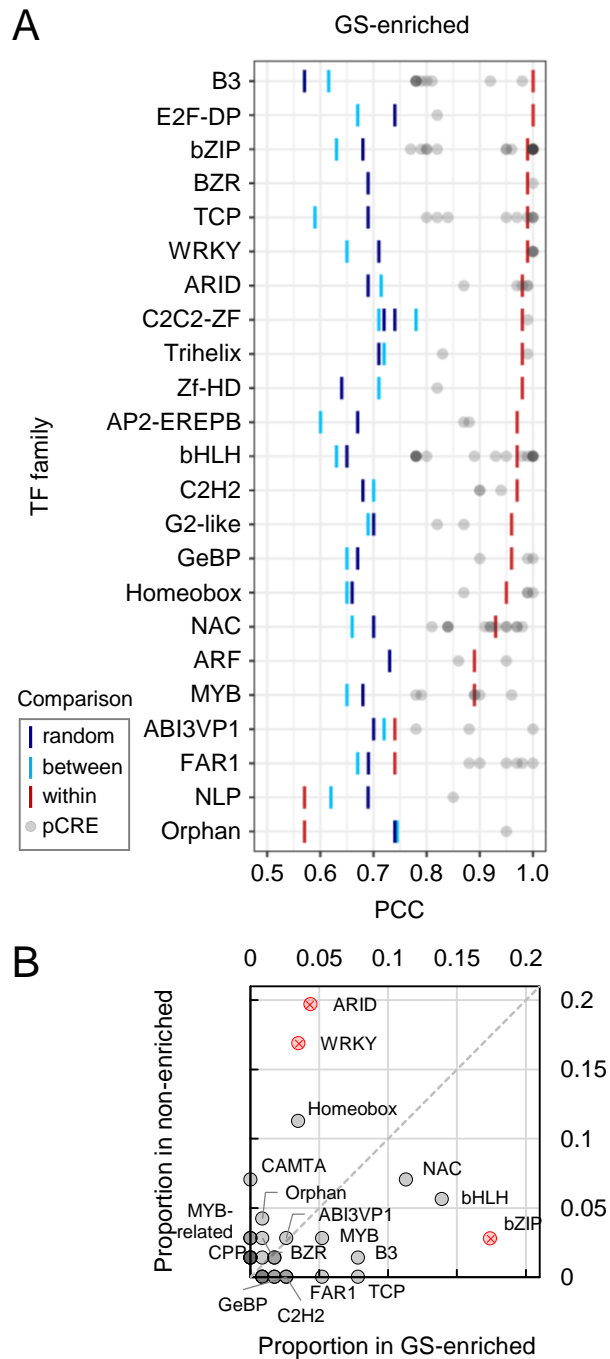


Figure 5. Similarity of freq-pCREs to *in vitro* TFBM.

A: Significance of sequence similarity for freq-pCREs from GS-enriched clusters and the best matching known TFBM. Bars represent 95th percentile (PCC) significance thresholds for within TF family (red, pCRE sequence is more similar to a specific TFBM than other TFBM from the same family), between TF families (light blue, pCRE sequence is more similar to a TFBM in a TF family than TFBM from other TF families), or random (dark blue, pCRE sequence is more similar to a TFBM from a family than random 6-mers). Similarity of freq-pCREs from non-enriched clusters to TFBM: **Supplemental Figure S5.** **B:** Proportion of TF family TFBM (representing freq-pCRE matches meeting at least “between” threshold) in GS-enriched clusters (x-axis) and non-enriched clusters (y-axis). TFBM matches significantly over-represented (FET, $q < 0.05$) in the GS-enriched or non-enriched cluster category are depicted in red and marked with “X”. Dashed line marks theoretical position for TF family TFBM with the same proportion in both categories.

Figure 6

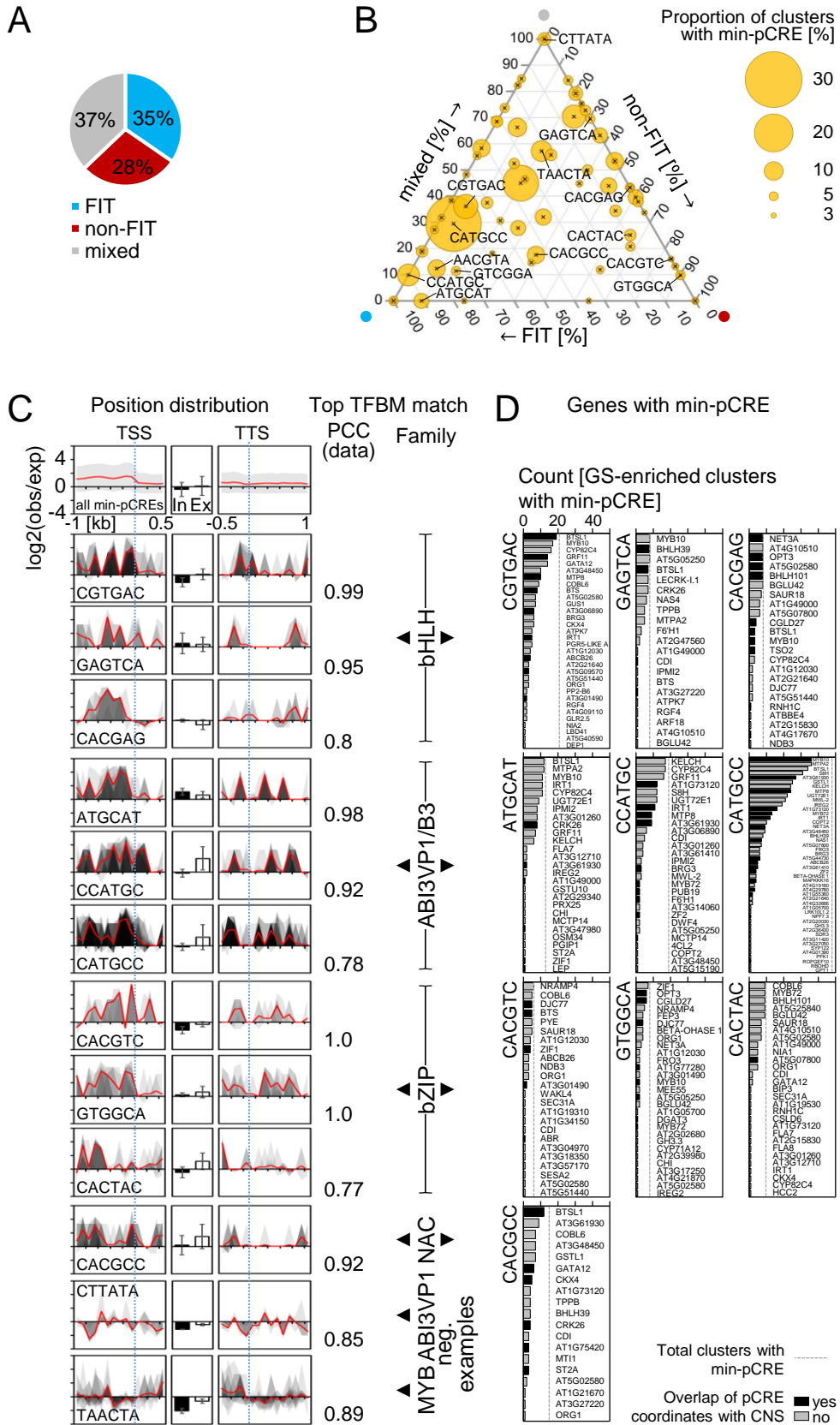


Figure 6

Figure 6. Characteristics of the most informative pCREs (min-pCREs).

A: Proportion of GS-enriched co-expression clusters enriched for FIT-dependent genes (FIT, blue), FIT-independent genes (non-FIT, red) or both (mixed, gray). **B:** Ternary plot including min-pCREs identified in >3% (n=5) GS-enriched clusters. Position of the min-pCREs corresponds to the normalized proportions of FIT, non-FIT, and mixed clusters in which the min-pCRE was identified. Bubble size corresponds to the overall proportion of GS-enriched clusters with min-pCRE. Labeled min-pCREs are shown in 6C, D or mentioned in the main text. **C:** Positional bias of all (mean with standard deviation; **top**) and selected min-pCREs (**below**) in the putative promoter region (**1st column**), all introns (In) and all exons (Ex) (**2nd column**; mean with standard deviation), and in the putative non-coding region (**3rd column**). **1st** and **3rd** column: position distributions in all co-expression clusters with min-pCRE (gray areas) with mean distribution (red line). TFBM matches (PCC) for each min-pCRE are shown (**4th column**) and min-pCREs are sorted by TF family. $\log_2(\text{obs}/\text{exp})$: \log_2 of the number of observed (obs) min-pCRE occurrences divided by the number of min-pCRE occurrences in randomized sequences (expected, exp). **D:** Genes which might be regulated by the selected min-pCREs. Count: number of GS-enriched clusters in which the min-pCRE was identified and which included the respective gene having the min-pCRE in its promoter. Dashed line: total number of GS-enriched clusters with the min-pCRE. Genes in which the min-pCRE overlaps with a CNS are designated with black bars. A high-resolution image is available as **Supplemental Figure S7**.