

# PlantFUNCO: integrative functional genomics database reveals clues into duplicates divergence evolution

Víctor Roces<sup>1</sup>, Sara Guerrero<sup>1</sup>, Ana Álvarez<sup>1</sup>, Jesús Pascual<sup>1\*</sup>, Mónica  
Meijón<sup>1\*</sup>

<sup>5</sup> <sup>1</sup> Plant Physiology, Department of Organisms and Systems Biology, Faculty of Biology  
<sup>6</sup> and Biotechnology Institute of Asturias, University of Oviedo, Asturias, Spain

<sup>7</sup> \* Co-correspondence: [meijonmonica@uniovi.es](mailto:meijonmonica@uniovi.es), [pascualjesus@uniovi.es](mailto:pascualjesus@uniovi.es)

8 # Contributed equally

9 Abstract

10 Evolutionary epigenomics and, more generally, evolutionary functional-genomics, are  
11 emerging fields that study how non-DNA-encoded alterations in gene expression  
12 regulation are an important form of plasticity and adaptation. Previous evidence  
13 analysing plants' comparative functional genomics has mostly focused on comparing  
14 same assay-matched experiments, missing the power of heterogeneous datasets for  
15 conservation inference. To fill this gap, we developed PlantFUN(ctional)CO(nsevation)  
16 database, which is constituted by several tools and two main resources: inter-species  
17 chromatin states and functional genomics conservation scores, presented and analysed  
18 in this work for three well-established plant models (*Arabidopsis thaliana*, *Oryza sativa*  
19 and *Zea mays*). Overall, PlantFUNCO elucidated evolutionary information in terms of  
20 cross-species functional agreement. Therefore, providing a new complementary  
21 comparative-genomics source for assessing evolutionary studies. To illustrate the  
22 potential applications of this database, we replicated two previously published models  
23 predicting genetic redundancy in *A. thaliana* and found that chromatin states are a  
24 determinant of paralogs degree of functional divergence. These predictions were  
25 validated based on the phenotypes of mitochondrial alternative oxidase knockout  
26 mutants under two different stressors. Taking all the above into account, PlantFUNCO  
27 aim to leverage data diversity and extrapolate molecular mechanisms findings from  
28 different model organisms to determine the extent of functional conservation, thus,  
29 deepening our understanding of how plants epigenome and functional non-coding  
30 genome have evolved. PlantFUNCO is available at  
31 <https://rocesv.github.io/PlantFUNCO>.

32 **Keywords:** evolutionary epigenomics, functional genomics, integrative approach,  
33 database, paralogs.

34 **Introduction**

35 A fundamental question in biology is how complex patterns of gene expression are  
36 determined to explain different phenotypes (Schmitz, Grotewold, and Stam, 2022;  
37 Marand et al., 2023). Today, it is largely known that genome function is dynamically  
38 regulated in part by chromatin organisation, which consists of histones, non-histone  
39 proteins and RNA molecules that package DNA (Ho et al., 2014). In this sense, the  
40 generation of comprehensive chromatin state (CS) maps, defined as the homogeneous  
41 coexistence of multiple chromatin modifications at the whole genome level, provides  
42 valuable information for annotating coding and non-coding genome features, including  
43 the identification of various types of regulatory elements. Chromatin states can facilitate  
44 our understanding of regulatory elements and variants associated with core life  
45 processes, such as development, and disease and stress responses (Liu et al., 2018).  
46 Great efforts have been made by the plant research community to contribute to the  
47 comprehension of chromatin mechanisms using different models (Zhao et al., 2020;  
48 Jamge et al., 2023); nevertheless, universal annotation allowing the extrapolation and  
49 unification of earlier conclusions across species/conditions still needs to be addressed.

50 Evolutionary theory has been dominated by the idea that selection proceeds by changes  
51 in allele frequencies within and between populations and mutations that occur randomly  
52 with respect to their consequences. The last theoretical and experimental advances in  
53 the field point to phenotypic plasticity as an adaptative trait subjected to natural selection,  
54 therefore, similar genotypes that differently develop appropriate phenotypes without  
55 sequence changes are equally responsible for evolutionary changes (Ashe, Colot, and  
56 Oldroyd, 2021; Monroe et al., 2022). This brings us to evolutionary epigenomics, and,  
57 more generally, evolutionary functional genomics, which are emerging fields evaluating  
58 how alterations in the conservation of epigenome regulators and cytosine methylation  
59 over multiple generations represent a crucial form of plasticity and epigenetic adaptation.  
60 Regulatory elements states have begun to be regarded as major targets of evolution,  
61 given that their diversity plays a critical role in phenotypic variance across all organisms,  
62 enabling them to adapt to various environmental niches (Yocca and Edger, 2022).  
63 Although relevant research in plants has lagged behind animal species (Schmitz et al.,  
64 2022), some of the most controversial findings in evolutionary biology use plants as  
65 model species, for example, mutations occur less often in functionally constrained  
66 regions, and epimutations are located in hotspots with specific chromatin features  
67 (Hazarika et al., 2022; Monroe et al., 2022). These findings support the clear importance  
68 of the plant kingdom in evolutionary functional genomics. Plants present a series of  
69 interesting molecular features that allow same sequence different function scenarios; for

70 instance, cytosine methylation is more easily transgenerationally transmitted due to soft  
71 epigenetic reset during meiosis and early development, epialleles are quite common and  
72 a relative high rate of duplication events, thus, multiple original exact gene copies with  
73 distinct selection pressures in response to the environment may exist (Ashe et al., 2021;  
74 Cusack et al., 2021). Many comparative genomics studies interrogate sequence-  
75 conserved loci of interest across a wide range of species, and their functions are  
76 determined by perturbing their homologous in a single model organism. In this context,  
77 a maze of opportunities and challenges appears to systematically and confidently  
78 determine the extent of conservation at the functional genomics level between model  
79 species (Kwon and Ernst, 2021).

80 Previous evidence analysing comparative functional genomics has mostly focused on  
81 comparing same assay-matched experiments (Maher et al., 2018; Lu et al., 2019). These  
82 works have been crucial for the in-depth study of molecular machinery but lack the power  
83 of diverse datasets for conservation inference. In contrast to this narrow but deep  
84 knowledge bottleneck, we adopted a broad but shallow approach using heterogeneous  
85 functional genomics to search directly simple large-scale answers that we would never  
86 have contemplated asking based on our understanding of single assay/species  
87 information (Kliebenstein, 2019). In the current Earth Biogenome era, an increasing  
88 number of genomes and functional tracks are becoming available (Expósito-Alonso et  
89 al., 2020), thus highlighting the urge to use integrative tools that consider the vast  
90 diversity of biological strategies and enable wide genomic element characterisation.  
91 Considering the abovementioned knowledge trade-off, in the present study, we  
92 introduced PlantFUN(ctional)CO(nservation), an integrative functional-genomics  
93 database constituted by several tools and two main resources, inter-species chromatin  
94 states and functional genomics conservation scores, for the well-known plant models  
95 *Arabidopsis thaliana*, *Oryza sativa* and *Zea mays*. To illustrate how the results derived  
96 from the generated resources could be functionally relevant, we developed an  
97 application of the database and found that chromatin state information improved the  
98 paralogous degree of functional divergence predictions. Lastly, we validated the  
99 redundancy predictions based on the phenotypic effects of alternative oxidase (AOX)  
100 gene knockout mutants under several stressors and provided insights into the evolution  
101 of these genes.

102    **Results**

103    Characterisation of shared and species-specific chromatin states

104    We generated a universal CS map annotation from 10 common chromatin modifications  
105    (greatest number of tracks found simultaneously available) (**supplementary fig. S1**)  
106    using hiHMM software for three widely-studied model plant species: *A. thaliana*, *O. sativa*  
107    and *Z. mays*. We focused our analysis on a model with 16 CSs (see **Methods**). In turn,  
108    the states were divided into 5 functional groups (bivalent, active, divergent, repressive  
109    and quiescent/no-signal), with different levels of genome coverage, transposable  
110    element (TE) enrichment and overlap with other genomic features (**fig. 1**).

111    The co-occurrence of chromatin modification pairs exists between these species, but  
112    there are clearly specific patterns in both CSs and correlation analyses (**fig. 1**;  
113    **supplementary fig. S2**). Despite the diversity of data, we found some conserved  
114    chromatin definitions, such as Bivalent TSS/Promoter CS1, which is strongly linked to all  
115    active marks with very low enrichment in H3K27me3 and without the clear presence of  
116    heavy repressive marks, such as 5mC and H3K9me2; and Active CS6, which is  
117    established in gene bodies and mainly constituted by H3K36me3, H3K4me2, H3K4me3  
118    and H3K9ac in the three species. However, many CS definitions exhibit species-specific  
119    nuances at different levels, which could actually reflect how epigenomic complexity has  
120    evolved in plants. The various degrees of CS divergence were determined based on CS  
121    chromatin modifications composition (**fig. 1, top panel**) and genomic distribution (**fig. 1,**  
122    **bottom panel**). Ranging from less to more divergent: 1) states which shared genomic  
123    distribution and were constituted by chromatin modifications with the same roles but  
124    covered with different chromatin modifications, such as Heretochromatin 1 strong CS11  
125    and Heterochromatin 2 weak CS12 (**fig. 1**). Repressive modifications, which were also  
126    pinpointed in the correlation analysis with the highest inter-species variance  
127    (**supplementary fig. S2**), suggested two distinct types of heterochromatin across  
128    species, requiring H3K27me3 for strong and H3K9me2 for weak definitions in *A. thaliana*.  
129    However, they were not necessary in *O. sativa* or *Z. mays*. 2) Landscapes whose  
130    chromatin modifications and genomic distribution gradually transitioned between  
131    species. A good case representing this could be Active weak TSS > TES CS8, mainly  
132    dominated by H3K36me3 deposition in gene bodies and TSS in *A. thaliana*, while in the  
133    two remaining species H3K4me2 is added and the distribution changed towards the TES.  
134    3) Ultimately, the divergent region CS10 had totally different chromatin modifications and  
135    genomic distribution profiles. CS10 corresponded to heterochromatic, bivalent and active  
136    states in *A. thaliana*, *O. sativa* and *Z. mays*, respectively.

137 We next performed additional annotation analyses based on non-common chromatin-  
138 binding proteins and histone marks tracks for all species under study to test our states  
139 definitions (**fig. 2**). There was evidence supporting our interpretation of the states for  
140 each species under study. For example, RNA polymerase II (Pol2) was significantly  
141 located in all active and several bivalent states, and there was enrichment of the well-  
142 known H3K9-demethylase (IBM1) and transposon-methylase (CMT3) over  
143 heterochromatic states in *A. thaliana*. Most of the transcription factors (TFs) observed in  
144 heterochromatin states were related to flowering, an organ missed in our collection, and  
145 cell cycle/division functions, which have been previously described as present in  
146 chromatin barriers and strictly under control, with low expression levels (Feng and  
147 Michaels, 2015; Velay, Méteignier, and Laloi, 2022). Essentially, all non-common active  
148 and repressive histone marks/variants evaluated were enriched in active/bivalent and  
149 heterochromatic states, respectively, with only two exceptions: H3K27me1 located in  
150 Bivalent Promoter CS2 in *A. thaliana*, which did not impact the state definition because  
151 this was already presented as bivalent due to the presence of H3K27me3; and  
152 H3K9me1/me3 in Active gradual bivalent flank > intergenic CS7 in *O. sativa*. Although  
153 the initial definition included gradual bivalent, this only alluded to *Z. mays*, as *O. sativa*  
154 CS7 was absent of any repressive mark; therefore, this could potentially increase the  
155 CS7 relationship between both Poaceae family members. We decided to be conservative  
156 and maintain our initial interpretation because H3K9me3 data were not available for all  
157 species.

158 Taking advantage of the inter-species approach, we further evaluated whether the states  
159 could involve evolutionary information. We observed a remarkable gradient across  
160 functional groups, excluding quiescent/no signal from the analysis due to the lack of  
161 epigenetic regulation (**fig. 3; supplementary figs. S3 and S4; supplementary table**  
162 **S1**). A decreasing trend in gene functional convergence (KO and GO) and the proportion  
163 of orthologous relationships was identified, following the order active > bivalent >  
164 heterochromatin, illustrated by CS6 > CS1 > CS11, respectively (the first state of each  
165 functional group was selected for representation). CS10 represented a divergent state  
166 corresponding to heterochromatic, bivalent and active states in *A. thaliana*, *O. sativa* and  
167 *Z. mays*, respectively. Additionally, most of the PhastCons elements' genomic overlaps  
168 were located in the active and bivalent states (**fig. 4**). Conserved non-coding elements  
169 (CNEs) localisation in the same states for *A. thaliana* and the greater number of CNE  
170 enriched states when comparing both species of monocots again showed how CS could  
171 reflect the closer distance between *O. sativa* and *Z. mays*. Even though most of the  
172 states enriched in conserved TF binding sites (BS) were active and bivalent in *A. thaliana*

173 and *O. sativa*, we did not observe a constrained pattern for the three species in TF motifs  
174 and genetic variability annotation modules (**fig. 4**). In opposition to conservation, these  
175 results could indicate that CS information is still useful because significant overlaps were  
176 detected, but it probably reflects species-specific features in genetic variability and TF  
177 motif contexts.

178 Taken together, these discoveries introduce a single plant inter-species CS annotation  
179 as a resource to provide conservation and diversity evolutionary epigenomic information  
180 for future research.

181 Chromatin state features improve predictions of paralogs functional  
182 divergence

183 To exemplify an application of the generated resource, we reproduced two previously  
184 published models predicting *A. thaliana* genetic redundancy (Cusack et al., 2021; Ezoe,  
185 Shirai, and Hanada, 2021), including CS information to determine which of the feature  
186 categories (such as evolutionary properties, gene expression patterns, protein sequence  
187 properties, epigenetic modification, and CS) could be relevant regulators of paralogs'  
188 functional divergence. To the best of our knowledge, *A. thaliana* is the only organism  
189 under study with an experimentally validated set of mutants for paralogous gene pairs,  
190 which allowed the development of these models. Under the initial hypothesis that two  
191 paralogs covered by different state profiles are more likely to have divergent functions,  
192 we computed similarity and distance metrics between both CS profiles and fed these  
193 data to the abovementioned models (**fig. 5A**; see **Methods**).

194 For the models developed by Ezoe, Shirai, and Hanada, 2021 (**fig. 5B-E**), we first  
195 checked whether the custom chromatin state metric (CCSM; see **Methods**) proposed  
196 could be a determinant of functional divergence using the same paralogous gene pairs  
197 as the original article (**fig. 5B**). High and low CCSM values were significantly associated  
198 with high and low diversified pairs, respectively (P-value = 3.4e-15, two tailed Wilcoxon  
199 rank sum test). Despite the epigenomic features tested in the reference did not pass this  
200 threshold, our CS metric even joined the two best explanatory variables Ka/Ks (protein  
201 divergence rate) and Re/Ks (gene expression similarity rate) in terms of relative  
202 importance (**fig. 5C**; see **Methods**). These results indicate out the need to use integrative  
203 metrics when predicting genome elements. Logistic regression models (see **Methods**)  
204 using different sets of features were compared by calculating the area under the curve-  
205 receiver operating characteristic (AUC-ROC) and the area under-precision recall curve  
206 (AU-PRC) values (**fig. 5D**). Models including CS information had higher AUC-ROC and  
207 AU-PRC values and slightly improved the performance of the best final model reported

208 in the original article ( $Ka/Ks+Re/Ks$ ). This improvement was more obvious in the reduced  
209 formula ( $Ka/Ks+Re/Ks+CCSM$ ) and the small range of improvement between full  
210 ( $Ka/Ks+Re/Ks+CCSM+FD+PPI+GO$ ) and reduced formulas also agreed with the  
211 information reported by the main article. The degree of functional divergence (DFD) can  
212 be inferred from the best formula by logistic regression analysis. DFD values close to 0  
213 and 1 reflected low ( $<0.5$ ) and high ( $>0.5$ ) functional divergence, respectively. To enable  
214 the potential validation of paralogous pairs DFD in upcoming studies and to minimise the  
215 erroneous assignment of high and low diversified duplicates, we calculated 5% FDR as  
216 a threshold. DFD stringent thresholds were 0.93 and 0.46 for high and low diversified  
217 pairs, respectively (**fig. 5E**). **Supplementary table S3** contains labelled genome-wide  
218 predictions with additional filters to assist paralog redundancy experimental verification  
219 (see **Methods**).

220 In contrast, the models developed by Cusack et al. (2021) (**fig. 5F-I**) categorised  
221 redundancy into different definitions, covering a plethora of features with distinct  
222 transformations. Consequently, we opted to incorporate all CS metrics to model  
223 redundancy for each definition, resulting in four different sets: RD4 (extreme redundancy,  
224 where single mutants have no abnormal phenotype, and the double mutant is lethal;  
225 without CS information), RD4C (with CS information), RD9 (inclusive redundancy,  
226 general definition that also included RD4 gene pairs; without CS information) and RD9C  
227 (with CS information). Analysis of models without CS information (RD4 and RD9)  
228 revealed that the number of variables and the relative importance of the six feature  
229 categories largely corroborated the discoveries in the reference (**fig. 5F**). In summary,  
230 the ranking from best to worst, based on median importance ranks in those categories  
231 for RD4/RD9-based models (without CS information), was functional annotation (37/16)  
232 > network properties (57.5/64.5) > evolutionary properties (76/110) > gene expression  
233 (104/105) > protein properties (145/88) > epigenetic modifications (121/127), with gene  
234 expression being the category with the highest number of variables in both cases. These  
235 findings validated the reproducibility of the models and ensured rigorous interpretation  
236 of subsequent results. Considering RD4C/RD9C-based models (with CS information),  
237 the CS feature category was sixth/second in importance rankings and emerged as the  
238 first in terms of the number of variables for both cases. This suggests that CS information  
239 is more valuable when predicting general (RD9 definition gene pairs) than extreme  
240 redundancy (RD4 definition gene pairs). This notion was further verified when comparing  
241 SVM models (see **Methods**) with different sets using AUC-ROC and AU-PRC values  
242 (**fig. 5G,H**). While CS data notably improved predictions for general redundancy (RD9C  
243 vs RD9, AUC-ROC = 0.665 vs 0.634, AU-PRC = 0.651 vs 0.603), it also reduced the

244 values for the extreme definition (RD4C vs RD4, AUC-ROC = 0.807 vs 0.842, AU-PRC  
245 = 0.795 vs 0.825). Finally, we observed that the intersection with the highest number of  
246 features was common to all sets suggesting that the core predicting power remained  
247 constant for all models, thereby ensuring accurate comparisons between all mentioned  
248 models (**fig. 5I**).

249 Collectively, we revealed that CS information could give clues into duplicates' general  
250 functional divergence, corroborated by the replication of two independent previously  
251 published models.

252 **Defining functional genomics conservation scores and the database**  
253 Evolutionary functional (epi)genomics is an emerging field of study with a growing body  
254 of literature reporting the massive generation of functional genomics data, however, the  
255 determinants underlying these processes are still not well understood due to a lack of a  
256 holistic point of view. To fill this gap, we adopted an integrative approach and expanded  
257 the resource generated with functional genomics conservation scores computed by the  
258 LECIF algorithm (Kwon and Ernst, 2021). LECIF was applied to integrate epigenomic,  
259 CSs, whole genome alignments and transcriptomic information for all pairwise  
260 comparisons between the species. By querying the LECIF scores, we sought to identify  
261 genomic regions with a high degree of functional tracks convergence and, therefore,  
262 similar phenotypic properties (**fig. 6A**).

263 To research elements highlighted by LECIF, we characterised the genome distribution  
264 of the scores over genetic variability, chromatin states and conservation modules. In all  
265 comparisons, the LECIF score density decreased in centromeres due to the lower  
266 number of alignments in these regions (**supplementary fig. S5**). As mentioned  
267 previously, we did not find a constrained pattern in the genetic variability module.  
268 Although both *Z. mays* contrasts and *O. sativa* vs *Z. mays* GWAS significant SNPs were  
269 enriched in regions with high functional conservation, neither *A. thaliana* contrast  
270 reflected any enrichment and *O. sativa* vs *A. thaliana* was enriched in regions with low  
271 LECIF scores (**fig. 6B-D; bar plots**). This could be explained by a balanced significant  
272 SNP distribution in the *A. thaliana* genome due to its architecture and higher number of  
273 GWAS, more similarity in the traits studied between the monocots and/or *O. sativa* only  
274 being able to retain functional conservation information related to the closest species.

275 In the CS module, genome-wide distributions were shifted to the left because of the  
276 higher weights of negative (only aligned) vs positive (aligned and functionally conserved)  
277 samples to ensure that only regions with strong functional evidence were underlined (**fig.**  
278 **6E-G; histograms**). To validate that the LECIF score displays the expected cross-

279 species similarity in functional genomics features, we examined it in relation to CS  
280 annotation. In each of the six query vs target comparisons, CS linked to strong regulatory  
281 or transcription activity tended to have a higher mean LECIF-score than the other states  
282 (**fig. 6E-G; violin plots**). We investigated cross-species CS similarity for different ranges  
283 of the LECIF score (**fig. 6E-G; line plots**). As the LECIF score increased, cross-species  
284 CS agreement was gradually higher in the active, bivalent and heterochromatin  
285 functional groups. This pattern was not fulfilled for divergent and quies/no-signal states  
286 because similarity was not expected by definition and the absence of epigenetic  
287 regulation, respectively. To provide further proof, we analysed CS annotations in regions  
288 where functional genomics (LECIF) and comparative genomics (PhyloP) scores  
289 disagreed (**fig. 6E-G; horizontal grouped bar plots**). Specifically, for pairs of regions  
290 where the LECIF score was high (percentile rank>60) and the PhyloP score was low  
291 (percentile rank<40), we computed CS similarity. We appreciate that such pairs were  
292 more likely to exhibit convergent states for all groups and vice versa.

293 Next, we evaluated the relationships between functional/comparative-genomics scores  
294 and annotations more deeply (**fig. 6H-J; box plots**). As we studied distantly related  
295 species, the scores of annotations with a high coverage percentage in aligning regions,  
296 such as PhastCons/PhyloP (Tian, Yang, Meng, Jin, and Gao, 2020) sequence-based  
297 conservation, would be influenced by the high negative:positive weights ratio. We found  
298 that regions overlapping the PhastCons elements did not have a greater average LECIF  
299 score compared to the genome-wide distribution, and the LECIF-score was not  
300 correlated with the PhyloP score (min-max range: 0.04-0.119 and 0.005-0.118 for PCC  
301 and SCC, respectively). Interestingly, CNEs followed the same trend as PhastCons  
302 elements except for Poaceae members vs *A. thaliana* pairs, which had higher LECIF  
303 scores. This is reasonable since CNEs preserved during longer timescales are more  
304 likely to be functionally conserved.

305 In summary, these reports suggest that plant LECIF scores can capture functional  
306 conservation without being correlated with other comparative genomics and sequence  
307 constraint scores. We expect the LECIF score and inter-species CS to be useful tools  
308 for unifying and extrapolating molecular mechanism discoveries using different model  
309 systems, thus, we developed an integrated hub called PlantFUN(ctional)CO(nservation)  
310 to provide interactive user-friendly functionalities for further requests (**fig. 6A**; see  
311 **Methods**). The PlantFUNCO database is available at  
312 <https://rocesv.github.io/PlantFUNCO/>.

313 Experimental validation of potential divergent duplicates  
314 To illustrate that the functional uses of the database could be translated into solutions  
315 for complex biological problems, we focused on the experimental validation of  
316 mitochondrial alternative oxidase (AOX) redundancy in *A. thaliana*. Although these pairs  
317 did not pass the stringent threshold ( $>0.93/<0.46$ ; **fig. 5E**), they presented high enough  
318 DFD values to be considered high divergent paralogs (AOX1A-AOX1C: 0.77, AOX1A-  
319 AOX1D: 0.72, AOX1C-AOX1D: 0.89; **fig. 7**). We assessed AOX redundancy by  
320 monitoring root phenotypes under two stressors, considering the previously described  
321 roles of these genes in response and retrograde signalling (Fuchs et al., 2022). Two out  
322 of five paralogs are not root expressed (Papatheodorou et al., 2020), simplifying the  
323 system and evaluating the seedling stages. The DFD of duplicates can be inferred based  
324 on the phenotypes of knockout plants. When single knockouts exhibit abnormal  
325 phenotypes related to the wild-type (WT, Col-0) under a specific condition, the duplicates  
326 are not compensated by the other gene copies; thus, they are assumed to be functionally  
327 divergent (Ezoe, Shirai, and Hanada, 2021).

328 Seedling phenotypes followed the same pattern for the control and mock conditions,  
329 there were significant differences for all AOX genotypes in root length  
330 (WT>*aox1c*>*aox1a*>*aox1d*), hypocotyl length (*aox1c*>*aox1d*>*aox1a*>WT) and  
331 root:hypocotyl ratio (WT>*aox1a*/*aox1c*>*aox1d*) (**fig. 7**). Under drought/heat (PEGxHeat)  
332 stress, significant differences were also observed, with two exceptions: *aox1c* root length  
333 and *aox1a* hypocotyl length. We established an additional stress assay using Antimycin  
334 A (AA), a mitochondrial complex III inhibitor that can be tolerated in plants due to electron  
335 bypass via AOX, but not when the activity of these genes is suppressed/diminished  
336 (Strodtkotter et al., 2009). Only root length was monitored because of the small size of  
337 *aox1a* seedlings. Again, significant changes were found for all AOX genotypes measured  
338 in root length and root:hypocotyl ratio. The greater p-values for hypocotyl length in  
339 drought/heat and no significance in AA suggest a general stress hypocotyl elongation  
340 mechanism in these mutants. In view of the roles of the AOX genes in the redox state,  
341 DAB staining quantification was performed to measure hydrogen peroxide levels.  
342 Although both stressors agreed in the WT, *aox1d* relevant increase of hydrogen peroxide  
343 levels, and *aox1c* was not significant; *aox1a* trends were not congruent. In *aox1a*,  
344 hydrogen peroxide content change was not meaningful for drought/heat, while a  
345 significant increase was detected during AA. Finally, in terms of functional genomics, the  
346 dominant isoform AOX1A seems to be the most crucial because it was covered by active  
347 CS and marked with high LECIF scores compared to *O. sativa*.

348 In brief, these findings validated our high divergence predictions and set a scenario in  
349 which AOX1A appeared to retain the ancestral function, allowing the understanding of  
350 the remaining AOX gene redundancy in relation to this reference.

351 **Discussion**

352 We introduced PlantFUNCO, a database that allows for further inspection of the crosstalk  
353 between evolution and epigenome/functional non-coding genome. This database is  
354 derived from two resources presented and analysed in this work for three well-  
355 established plant models. We generated inter-species CS using hiHMM (**fig. 1**). While  
356 this flexible framework provides a consistent definition of CS across multiple genomes,  
357 making the extrapolation of intra-species analyses between them easier, the stack  
358 approach allows for an understanding of the potential epigenomic regulation over several  
359 tissues/conditions, such as differentiating constitutively active/repressive regions (Vu  
360 and Ernst, 2022). CS links with different types of evolutionary information set a  
361 foundation for the epigenomics inter-species perspective (**figs. 3 and 4; supplementary**  
362 **fig. S3 and S4**). All the approaches have trade-offs; thus, this resource should be  
363 considered complementary to and not a replacement for other single-species/condition  
364 annotations. We obtained functional genomic conservation scores using LECIF. In  
365 accordance with the abovementioned framework, LECIF can handle very diverse  
366 datasets and take advantage of them to quantify functional conservation. Plant ECIF  
367 score elucidated functional genomic cross-species agreement without being correlated  
368 with other comparative genomics sources (**fig. 6**). This probably reflects a  
369 complementary side of evolution. Despite the greater divergence between plants models  
370 compared to metazoans (Ho et al., 2014; Kwon and Ernst, 2021), both resources results  
371 are congruent with a higher plant epigenomic/functional complexity probed by more  
372 states with species-specific features and lower LECIF scores.

373 A major focus of this study was to illustrate the application of the generated resources.  
374 Due to the holistic approach adopted and exploiting that our inter-species CS could differ  
375 between constitutively active/repressive regions, we replicated two previously published  
376 models predicting paralogous functional divergence in Arabidopsis (Cusack et al., 2021;  
377 Ezoe et al., 2021), including our CS information. We determined whether CS similarity  
378 could be a determinant of duplicates' degree of functional divergence under the initial  
379 hypothesis that two paralogs covered by different state profiles are more likely to present  
380 distinct functions. Although the models are far from perfect, useful information about  
381 gene features can be extrapolated. These models independently reported CS  
382 information as relevant, and including this type of data improved general redundancy

383 predictions (fig. 5). This shows an example of how PlantFUNCO's integrative resources  
384 can be effectively employed to predict genomic elements.

385 An important goal of a database is to functionally translate applications into solutions to  
386 explain complex biological mechanisms; thus, we decided to check the redundancy  
387 predictions of AOX genes. DFD values were high enough to be considered, and earlier  
388 AOX research made their context of high biological interest. Briefly, past reports mainly  
389 focused on the dominant isoform AOX1A (Giraud et al., 2008) which has a partial  
390 redundancy relation described with AOX1D (Strodtkotter et al., 2009), but current  
391 literature is not congruent with the use of single *aox1a* or double *aox1a-aox1d* mutants  
392 to discover causal drivers of retrograde-signalling/metabolism/stress-response (Giraud  
393 et al., 2009; Clercq et al., 2013; Oh Khim et al., 2022; Oh Khim et al., 2023). Additionally,  
394 more AOX isoforms exist, but their relationships were still not addressed. The abnormal  
395 seedling growth observed in control and mock conditions for all tested single mutants  
396 (*aox1a*, *aox1c*, *aox1d*) (fig. 7) validated the high functional divergence predicted by  
397 PlantFUNCO since in case of redundancy, other duplicates could rescue these  
398 phenotypes (Ezoe, Shirai, and Hanada, 2021). Our findings suggest that the dominant  
399 isoform AOX1A could retain the ancestral AOX function because it is marked as  
400 functionally conserved with the distantly related *O. sativa* and is the only one covered by  
401 an active CS; thus, all redundancy relationships can potentially be compared to this  
402 gene. Considering that oxidative stress was more severe than drought/heat conditions,  
403 we found putative evidence of a probable stress-dependent partial non-mutual  
404 redundancy of AOX1D to AOX1A. Although AOX1D could partially alleviate *aox1a* raw  
405 hydrogen peroxide content under drought/heat (no significance), during more severe  
406 oxidative conditions, AOX1D would not be enough to supply the AOX1A function  
407 (significant) (Strodtkotter et al., 2009). It is defined as a potential non-mutual relationship  
408 because, in all cases, *aox1d* phenotypes remained significant. Finally, non-meaningful  
409 differences in raw hydrogen peroxide content for both stressors and WT-like root lengths  
410 under drought/heat in *aox1c* would indicate that AOX1C as a non-stress-responsive  
411 gene. This could agree with the previously described AOX1C AA expression insensitivity  
412 (Yoshida and Noguchi, 2009), but we still found significant differences in root length in  
413 our severe oxidative assay. Compared to other genotypes, the p-value was close to not  
414 significant; thus, AOX1C may only be related to stress under severe conditions and could  
415 probably be defined as almost non-stress-responsive. In summary, stress seems to be  
416 a crucial evolutionary force driving sub-/neofunctionalisation (Panchy, Lehti-shiu, and  
417 Shiu, 2016) in AOX genes, and we characterised the unknown AOX1C as almost stress-  
418 insensitive during the seedling stages. Furthermore, extra attention should be taken

419 when using double AOX mutants to identify the causal determinants of biological  
420 processes because all AOX genes evaluated appeared to be functionally divergent  
421 during early development.

422 While we expect PlantFUNCO to be useful, we acknowledge certain limitations. Owing  
423 to our data collection design, the main goal of inter-species CS resources is to conduct  
424 intra-species analyses while leveraging the advantage of having additional layers of  
425 interpretation, including direct correspondence between CS and  
426 conservation/divergence relationships established across species. Direct cross-species  
427 comparisons of equivalent loci or CSs should be undertaken only in conjunction with  
428 plants' LECIF scores, as this algorithm is explicitly designed to handle highly diverse  
429 datasets. There may be states/regions that are functionally conserved, but have low  
430 scores/agreement in the database since the evidence was not present in our collection.  
431 While the interpretation of the resources generated is less ambiguous due to the broad-  
432 shallow perspective adopted, we also perceived that PlantFUNCO is limited by the input  
433 functional genomics resolution and does not provide direct information about which  
434 particular tracks/conditions supported the evidence. The results promoted the potential  
435 application of PlantFUNCO to further test new hypotheses in the context of duplicate  
436 evolution and other genomic elements prediction. For example, as CSs are determinants  
437 of paralogs' functional divergence and LECIF scores highlight regions with high  
438 phenotypic similarity, it could be possible to identify genes that are more likely to retain  
439 ancestral functions if high scores are found between orthologs in distantly related  
440 species (**fig. 6A**). Here, we focused on *A. thaliana*, *O. sativa* and *Z. mays*, which are  
441 widely used models in plant science research with substantial high-quality publicly  
442 available data. Given the increasing availability of epigenomics and functional genomics  
443 datasets, the utility of PlantFUNCO will continue to grow and serve as an additional  
444 resource to simplify functional conservation annotations for a more diverse set of species  
445 such as *Chlamydomonas reinhardtii*, *Marchantia polymorpha* and *Solanum lycopersicum*. Overall, PlantFUNCO aims to leverage data diversity and extrapolate  
446 findings from different models to determine the extent of molecular conservation, thus  
447 deepening our understanding of how plants epigenome and functional non-coding  
448 genome have fascinatingly evolved.

## 450 Methods

451 An overview of the methods workflow used in this study is shown in **supplementary fig.**  
452 **S1.**

453 Data collection

454 We collected epigenomic (ChIP-, MeDIP-, ATAC- and DNase-seq) and transcriptomic  
455 (RNA-seq) data from three plant model species: *Arabidopsis thaliana*, *Oryza sativa* and  
456 *Zea mays*.

457 For the epigenomic data, we used the previously published collection from the PCSD (Y.  
458 Liu et al., 2018) to ensure high-quality data. Then, we expanded the abovementioned list  
459 to include new common chromatin modifications published in recent years  
460 (**supplementary table S1**).

461 For the transcriptomic data, we used the baseline collection of the manually curated  
462 database EBI-ATLAS (Papatheodorou et al., 2020). We filtered this list to include only  
463 studies that covered multiple tissues/organs (**supplementary table S2**).

464 Epigenomic data processing

465 Raw reads were trimmed and adapters were removed using trim\_galore v.0.6.6 as an  
466 interface to CutAdapt (Martin, 2011). The remaining reads were aligned to the reference  
467 genome (*A. thaliana*: TAIR10, *O. sativa*: IRGSP-1.0, *Z. mays*: RefGen v4) using the  
468 bowtie2 algorithm (Langmead and Salzberg, 2012). Mapped reads with a MAPQ > 30  
469 were used to secure the optimal quality of the data. Aligned reads were sorted using  
470 SAMtools v.1.9, and duplicate reads were removed using Picard v.2.26  
471 (<https://github.com/broadinstitute/picard>). For all subsequent analyses we performed  
472 peak calling (narrow and broad), signal track building, correlation, and formatting with  
473 MACS2 and deepTools (Zhang et al., 2008; Ram et al., 2016). Briefly, the –g argument  
474 was changed for each species (*A. thaliana*: 91254070, *O. sativa*: 215463918, *Z. mays*:  
475 1975365725), FDR < 0.1 was used for broad peak calling, and the arguments --nomodel/  
476 --shift -75 --extsize 150 were added for ATAC- and DNase-seq file processing. Additional  
477 information detailing intra-species correlations and variance can be found in  
478 **supplementary table S1**. To guarantee the reproducibility of the analysis, a docker was  
479 created and it is available at <https://hub.docker.com/r/rocesv/plantina-chiplike>.

480 Inter-species chromatin states definition and annotation

481 We applied hiHMM (Sohn et al., 2015) to jointly infer multiple species chromatin states  
482 (CS) using common chromatin modifications signal tracks from several tissues as input.  
483 Signal tracks consisted of scaled log2 (fold enrichment + 0.5) values averaged in 200 bp  
484 bins in all three species, as described in the original application (Ho et al., 2014). The  
485 analysis was restricted to nuclear chromosomes. hiHMM can handle an unbounded  
486 number of hidden states; thus, the number of states is learned from the training data  
487 instead of a pre-specified value by the user. The model inferred a total of 15 CSs with

488 unmappable regions added *a posteriori* as the 16th state to avoid any bias in the  
489 segmentation. We defined CSs based on the colocalisation of chromatin modifications  
490 and overlap enrichments of different genomic features using ChromHMM (Ernst and  
491 Kellis, 2017).

492 To further improve the interpretability of the states, additional annotations and  
493 descriptions were performed. The annotation was based on significant overlap  
494 enrichments using the LOLA package (Sheffield and Bock, 2016) and was divided as  
495 follows: 1) assessment of the presence of other epigenomic features employing non-  
496 common liftover information in PCSD; 2) conservation covered by PhastCons elements  
497 in PlantRegMap and pairwise CNEs; 3) transcription factor binding motifs collected in  
498 PlantRegMap (Tian et al., 2020); 4) genetic variability represented by significant SNPs  
499 compiled in GWAS-ATLAS and AraGWAS (Togninalli et al., 2020; Liu et al., 2023). The  
500 description involved KEGG-Orthology(KO)/Gene-Ontology(GO) enrichments using  
501 clusterProfiler/REVIGO, respectively, and gene biotype-orthology correspondence using  
502 inParanoid information stored in Phytozome (Goodstein et al., 2012).

### 503 Modelling paralogs' degrees of functional divergence

504 We reproduced two published models that predict genetic redundancy in *A. thaliana*  
505 paralogs (Cusack et al., 2021; Ezoe et al., 2021) including our inter-species CS distance  
506 metrics. To define state distance metrics, we first binned different genomic features  
507 (promoters and genes) into a fixed number of windows and computed both presence (1  
508 = present; 0 = absent) and frequency (% of bp covered in a window) vectors for each  
509 state and gene. Additionally, we included a third type of vector, with each element having  
510 the frequency of a particular state over a non-binned genomic feature. Lastly, distinct  
511 distance metrics were calculated between genes of the same paralog pair, comparing  
512 equivalent vectors using the philientropy package (Drost, 2018).

513 To reproduce both studies, we followed the workflow originally established for the best  
514 performing model. In brief, for the model described by Ezoe, Shirai and Hanada (2021)  
515 feature selection was executed by two-tailed Wilcoxon rank sum test p-values between  
516 pairs labelled as redundant or divergent, followed by logistic regression relative  
517 importance to examine the explanatory weights of the best variables. Since this model  
518 is designed to perform genome-wide predictions and only some of the distance state  
519 metrics could be informative, a small number of features are desirable. We combined the  
520 information of the best-scored features into a single metric defined as the custom  
521 chromatin state metric (CCSM) (**supplementary table S3**). To compare the performance  
522 of logistic regression models using different sets of features, we calculated the AUC-

523 ROC and AU-PRC values. All the analyses were conducted in the R software  
524 environment ([Team R Development Core 2013](#)).

525 However, in the model developed by Cusack et al. (2021) multiple transformations and  
526 interpretations of the same feature were included; thus, all the distance state metrics  
527 were considered. Only the available extreme (RD4) and inclusive (RD9) redundancy  
528 gene pair sets were analysed, deleting variables identified as mispredictors in the main  
529 article. Non-redundant gene pairs were randomly downsampled to generate balanced  
530 cross-validation sets. Feature selection was executed using random forest top 200 best  
531 transformed variables (determined by feature importance) for sets without (RD4-RD9)  
532 and with (RD4C-RD9C) chromatin information. The C value for the SVM algorithm was  
533 set as a hyperparameter during the tuning. To measure SVM performance using different  
534 feature sets, we calculated AUC-ROC and AU-PRC values. All analyses were conducted  
535 using the pipeline implemented and developed by the authors  
536 (<https://github.com/ShiuLab/ML-Pipeline>).

### 537 Genome-wide redundancy predictions

538 To generate genome-wide predictions, we used the best performing model from the first  
539 pipeline described above. The stringent threshold for identifying high and low diversified  
540 pairs with the logistic regression formula (DFD=degree of functional divergence) was  
541 defined by a 100 cross-validation test where the FDR was under 5 %. As a result,  
542 high/low divergent pairs have  $>0.5/<0.5$  and  $>0.93/<0.46$  DFD values with relaxed and  
543 stringent thresholds, respectively. *Arabidopsis thaliana* genes (longest sequence) were  
544 used as queries to search for self-match homologues with DIAMOND v2 (E-value=1e-  
545 04) (Buchfink, Reuter, and Drost, 2021). We only focused on pairs with the best hits,  $>$   
546 30% identity and  $>$  50% coverage. We identified 7852 pairs, of which 1444/6898 were  
547 predicted as high and 723/954 as low diversified duplicates with strict/relaxed thresholds,  
548 respectively. Ka/Ks (number of nonsynonymous/synonymous substitutions per  
549 nonsynonymous/synonymous site) and the similarity of expression patterns (Re) were  
550 calculated as described by Ezoe, Shirai and Hanada (2021). An additional table is  
551 provided with filters, such as the same second closest paralog and expression under  
552 stress and in the seedling stages, to assist experimental validation in future studies  
553 (**supplementary table S3**).

### 554 Experimental validation of potential divergent paralogs

555 The *A. thaliana* T-DNA insertion line *aox1a* (SALK\_084897) was previously described  
556 as a knockout and validated by genotyping before use (Fuchs et al., 2022). We  
557 characterised the *aox1c* (Sail\_420\_A04) and *aox1d* (SM\_3\_24421) insertion lines as

558 homozygous and knockout by genotyping and RT-PCR analysis, respectively. Briefly,  
559 RNA was extracted as described by Valledor et al. (2014) and quantified by a Navi  
560 UV/Vis Nano Spectrophotometer, integrity was evaluated by agarose gel  
561 electrophoresis. cDNA was obtained from 500 ng of RNA using the RevertAid kit  
562 (ThermoFisherScientific), where random hexamers were used as primers following the  
563 manufacturer's instructions. RT-PCR analysis reported these lines as knockouts  
564 because no amplification was detected in the mutants (all primers are available in  
565 **supplementary table S3**).

566 For stress evaluation, *aox1a*, *aox1c* and *aox1d* seeds were surface sterilised in 2.8%  
567 hypochlorite solution and washed several times with sterile water; they were stratified for  
568 3 days at 4° C in darkness. The in vitro culture of seeds was carried out in 12x12 plates  
569 (Greiner) containing 50 mL of MS medium, pH 5.8, 1% (w/v) sucrose and 0.8% (w/v)  
570 agar and they were vertically placed under a long-day photoperiod (16 h light 21° C, 8 h  
571 dark 18° C) for control conditions. To avoid a position effect, the four genotypes (Col-0  
572 as WT, *aox1a*, *aox1c* and *aox1d*) were located in every plate position by rotating sectors  
573 in different plates. For the combined drought/heat stress, 2.5% PEG8000  
574 (ThermoFisherScientific) was added to the initial plates and seedlings were subjected to  
575 37° C stress for 1 h every day at the same hour, gradually increasing and decreasing the  
576 temperature. For the antimycin A (AA) treatment, 50 µM AA (Sigma-Aldrich) was added  
577 to the initial plates; control conditions were set as a mock due to AA being dissolved in  
578 ethanol. Phenotypic monitoring was conducted 5 days after germination by scanning  
579 culture plates with high-resolution scans (EpsonPerfectionV600); hypocotyl and root  
580 lengths were measured with ImageJ software (Schneider, Rasband, and Eliceiri, 2012)  
581 in at least 12 biological replicates. Furthermore 3,3-Diaminobenzidine (DAB) staining  
582 (Sigma-Aldrich) was performed 5 days after germination for at least 3 biological  
583 replicates per treatment, following the protocol described by Daudi and O'Brien, (2012);  
584 DAB quantification was carried out using ImageJ.

## 585 RNA-seq data processing

586 The sequence quality of RNA-seq libraries was evaluated by FastQC and multiQC  
587 (Andrews, 2013; Ewels, Lundin, and Max, 2016). Raw reads were trimmed and adapters  
588 were removed using trim\_galore v.0.6.6. Cleaned reads were mapped using STAR  
589 v.2.7.10 (Dobin et al., 2013) changing the reference genome and minimum/maximum  
590 intron size according to species. Bigwig files were obtained using the *bamCoverage*  
591 command from deepTools (Ram et al., 2016).

592 Whole genome alignments and identification of conserved non-coding  
593 elements

594 Whole genome alignments (WGA) were computed for each pairwise comparison. In  
595 summary, *lastz* alignments with far (vs *A. thaliana*; >100 MYA according to TimeTree  
596 (Kumar et al., 2022)) and medium (*O. sativa* vs *Z. mays*; >15 and <100 MYA) *distance*  
597 arguments were performed using the CNEr package interface (Tan, Polychronopoulos,  
598 and Lenhard, 2019). This was followed by format conversion, chain building, and  
599 processing using *lavToPsl*, *maf-convert*, *axtChain* and *chainMergeSort*. RepeatFiller  
600 (Osipova, Hecker, and Hiller, 2019) was applied to the chains to improve the  
601 identification of conserved non-coding elements (CNEs). After RepeatFiller, we executed  
602 ChainCleaner (Suarez, Langer, Ladde, and Hiller, 2017) to improve alignment specificity  
603 and chains were then converted into alignment nets using Hillerlab *chainNet* and  
604 *netToAxt*. Finally, Axt files were used as input for the pairwise identification of CNEs  
605 using the CNEr package with 45-identity/50-length windows while considering the  
606 difference in whole genome duplication history between these species, as described by  
607 Ren et al. (2018).

608 To take advantage of previously processed epigenetic tracks in PCSD that are not  
609 included in our initial collection (not common for all species), we executed another WGA  
610 pipeline to lift over these files to the new reference assemblies. In summary, we used  
611 near as a *distance* argument, and skipped the RepeatFiller-ChainCleaner step because  
612 we aligned the same species, and liftover was carried out using CrossMap v.0.6.2 (Hao  
613 Zhao et al., 2014). To guarantee the reproducibility of the analysis, a docker was created;  
614 it is available at <https://hub.docker.com/r/rocesv/compncnes>.

## 615 Functional genomics conservation score

616 The LECIF algorithm (Kwon and Ernst, 2021) was applied to obtain a functional  
617 genomics conservation score between all possible pairwise comparisons, integrating  
618 whole genome alignments, epigenomics, CSs, and transcriptomic information. The  
619 negative to positive sample weight ratio was set to 10 because the species under study  
620 are distantly related, with a lower number of samples aligning but more likely to be  
621 functionally conserved. For the training and evaluation, we adopted the same approach  
622 as the authors based on odd and even chromosomes (**supplementary table S4**). LECIF  
623 downstream analyses were performed in the R software environment ([Team R](#)  
624 [Development Core 2013](#)).

625 **Database resource**  
626 We developed PlantFUN(ctional)CO(nservation) database to provide public availability  
627 of the functional integrative tracks generated in this work and to facilitate future research  
628 in evolutionary functional genomics. PlantFUNCO contains three main tools: 1) a search  
629 section with interactive tables to retrieve gene- or superenhancer-level (Zhao et al.,  
630 2022) functional and comparative genomics information; 2) a shiny-application to  
631 compute LOLA genomic overlap enrichments of user query bed files over CSs and  
632 LECIF/PhyloP binned scores; and 3) a JBrowse2 genome browser (Diesh et al., 2023).  
633 PlantFUNCO is available at <https://rocesv.github.io/PlantFUNCO>.

## 634 **Data availability**

635 All data generated in this study are available at the PlantFUNCO database  
636 <https://rocesv.github.io/PlantFUNCO> and <https://zenodo.org/record/7852329>. The code  
637 used in this work is available at [https://github.com/RocesV/PlantFUNCO\\_manuscript](https://github.com/RocesV/PlantFUNCO_manuscript).

## 638 **Acknowledgements**

639 This work was generously financed by the Spanish Ministry of Science, Innovation and  
640 Universities (PID2020-113896GB-I00). VR and AA were supported by a fellowship from  
641 the Spanish Ministry of Universities (FPU18/02953 and FPU19/01142, respectively). SG  
642 was supported by the Severo Ochoa Predoctoral Program (BP19-145). JP was  
643 supported by the Juan de la Cierva Incorporación Programme (IJC-2019-040330-I). We  
644 are grateful to Prof. James Whelan (Zhejian University) for kindly sharing the aox mutant  
645 lines used in this study.

## 646 **Conflict of interest**

647 The authors declare that there are no conflicts of interest.

## 648 **Author's contributions**

649 VR and MM conceived the study. VR designed the research. VR and AA collected the  
650 data and built the figures. SG performed all mutant generation, validation and stress  
651 experiments. VR performed computational analyses, analysed and interpreted the data,  
652 and wrote the manuscript. JP and MM supervised the study. All authors revised, read,  
653 and approved the final manuscript.

654      **References**

- 655      Andrews S. 2013. Babraham Bioinformatics -FastQC A Quality Control tool for High  
656      Throughput Sequence Data.
- 657      Ashe A, Colot V, Oldroyd BP. 2021. How does epigenetics influence the course of  
658      evolution ?         *Philosophical Transactions B*,         376(20200111).  
659      <https://doi.org/10.1098/rstb.2020.0111>
- 660      Buchfink B, Reuter K, Drost H. 2021. Sensitive protein alignments at tree-of-life scale  
661      using DIAMOND. *Nature Methods*, 18(April). <https://doi.org/10.1038/s41592-021-01101-x>
- 663      Clercq I De, Vermeirssen V, Aken O Van, Vandepoele K, Murcha M W, Law S R, Inzé  
664      A, Ng S, Ivanova A, Rombaut D, et al. 2013. The Membrane-Bound NAC Transcription  
665      Factor ANAC013 Functions in Mitochondrial Retrograde Regulation of the Oxidative  
666      Stress Response in Arabidopsis. *The Plant Cell*, 25(September), 3472–3490.  
667      <https://doi.org/10.1105/tpc.113.117168>
- 668      Cusack S A, Wang P, Lotreck S G, Moore B M, Meng F, Conner J K, Krysan PJ, Lehti-  
669      Shiu MD, Shiu-Han S. 2021. Predictive Models of Genetic Redundancy in Arabidopsis  
670      thaliana. *Molecular Biology and Evolution*, 38(8), 3397–3414.  
671      <https://doi.org/10.1093/molbev/msab111>
- 672      Daudi A, A O'Brien J. 2012. Detection of Hydrogen Peroxide by DAB Staining in  
673      Arabidopsis Leaves. *Bio Protoc.*, 2(18), 4–7.
- 674      Diesh C, Stevens G J, Xie P, Martinez T D J, Hershberg E A, Leung A, Guo E, Dider S,  
675      Zhang J, Bridge C, et al. 2023. JBrowse 2 : a modular genome browser with views of  
676      synteny and structural variation. *Genome Biology*, 1–21. <https://doi.org/10.1186/s13059-023-02914-z>
- 678      Dobin A, Davis C A, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M,  
679      Gingeras T R. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, 29(1),  
680      15–21. <https://doi.org/10.1093/bioinformatics/bts635>
- 681      Drost H. 2018. Philentropy : Information Theory and Distance Quantification with R. *The  
682      Journal of Open Source Software*, 1, 1–4. <https://doi.org/10.21105/joss.00765>
- 683      Ernst J, Kellis M. 2017. Chromatin-state discovery and genome annotation with  
684      ChromHMM. *Nature Publishing Group*, 12(12), 2478–2492.  
685      <https://doi.org/10.1038/nprot.2017.124>

- 686 Ewels P, Lundin S, Max K. 2016. MultiQC : summarize analysis results for multiple tools  
687 and samples in a single report. *Bioinformatics*, 32(June), 3047–3048.  
688 <https://doi.org/10.1093/bioinformatics/btw354>
- 689 Expósito-Alonso M, Drost H, Burbano H A, Weigel D. 2020. The Earth BioGenome  
690 project : opportunities and challenges for plant genomics and conservation. *Plant*  
691 *Journal*, 102, 222–229. <https://doi.org/10.1111/tpj.14631>
- 692 Ezoe A, Shirai K, Hanada K. 2021. Degree of Functional Divergence in Duplicates Is  
693 Associated with Distinct Roles in Plant Evolution. *Molecular Biology and Evolution*, 38(4),  
694 1447–1459. <https://doi.org/10.1093/molbev/msaa302>
- 695 Feng W, Michaels S D. 2015. Accessing the Inaccessible : The Organization ,  
696 Transcription , Replication , and Repair of Heterochromatin in Plants. *Annual Review of*  
697 *Genetics*, 49, 439–459. <https://doi.org/10.1146/annurev-genet-112414-055048>
- 698 Fuchs P, Bohle F, Lichtenauer S, Ugalde M, Araujo E F, Mansuroglu B, Ruberti C,  
699 Wagner S, Müller-Schüssle J, Meyer AJ, et al. 2022. Reductive stress triggers  
700 ANAC017-mediated retrograde signaling to safeguard the endoplasmic reticulum by  
701 boosting mitochondrial respiratory capacity. *The Plant Cell*, 34, 1375–1395.  
702 <https://doi.org/10.1093/plcell/koac017>
- 703 Giraud E, Ho L H M, Clifton R, Carroll A, Estavillo G, Tan Y, Howell KA, Ivanova A,  
704 Pogson BJ, Millar AH, et al. 2008. The Absence of ALTERNATIVE OXIDASE1a in  
705 Arabidopsis Results in Acute Sensitivity to Combined. *Plant Physiology*, 147(June), 595–  
706 610. <https://doi.org/10.1104/pp.107.115121>
- 707 Giraud E, Aken O Van, Ho L H M, Whelan J. 2009. The Transcription Factor ABI4 Is a  
708 Regulator of Mitochondrial Retrograde Expression of. *Plant Physiology*, 150(July), 1286–  
709 1296. <https://doi.org/10.1104/pp.109.139782>
- 710 Goodstein D M, Shu S, Howson R, Neupane R, Hayes R D, Fazo J, Mitros T, Dirks W,  
711 Hellsten U, Putnam N, et al. 2012. Phytozome : a comparative platform for green plant  
712 genomics. *Nucleic Acids Research*, 40(November 2011), 1178–1186.  
713 <https://doi.org/10.1093/nar/gkr944>
- 714 Hazarika R R, Serra M, Zhang Z, Zhang Y, Schmitz R J, Johannes F. 2022. Molecular  
715 properties of epimutation hotspots. *Nature Plants*, 8(February), 146–156.  
716 <https://doi.org/10.1038/s41477-021-01086-7>

- 717 Ho J W K, Jung Y L, Liu T, Alver B H, Lee S, Ikegami K, Sohn K, Minoda A, Tolstorukov  
718 MY, Appert A, et al. 2014. Comparative analysis of metazoan chromatin organization.  
719 *Nature*, 512(7515), 449–452. <https://doi.org/10.1038/nature13415>
- 720 Jamge B, Lorkovi Z J, Axelsson E, Osakabe A, Shukla V, Yelagandula R, Akimcheva S,  
721 Kuehn AL, Berger F. 2023. Histone variants shape chromatin states in *Arabidopsis*.  
722 *eLife*, 12(RP87714), 1–26. <https://doi.org/10.7554/eLife.87714.3>
- 723 Kliebenstein D J. 2019. Questionomics : Using Big Data to Ask and Answer. *The Plant*  
724 *Cell*, 31(July), 1404–1405. <https://doi.org/10.1105/tpc.19.00344>
- 725 Kumar S, Suleski M, Craig J M, Kasprowicz A E, Sanderford M, Li M, Li M, Stecher G,  
726 Hedges S B. 2022. TimeTree 5 : An Expanded Resource for Species Divergence Times.  
727 *Molecular Biology and Evolution*, 39(8), 1–6. <https://doi.org/10.1093/molbev/msac174>
- 728 Kwon S Bin, Ernst J. 2021. Learning a genome-wide score of human–mouse  
729 conservation at the functional genomics level. *Nature Communications*, 12, 2495.  
730 <https://doi.org/10.1038/s41467-021-22653-8>
- 731 Liu X, Tian D, Li C, Tang B, Wang Z, Zhang R, Pan Y, Wang Y, Zou D, Zhang Z, et al.  
732 2023. GWAS Atlas : an updated knowledgebase integrating more curated associations  
733 in plants and animals. *Nucleic Acids Research*, 51(October 2022), 969–976.  
734 <https://doi.org/10.1093/nar/gkac924>
- 735 Liu Y, Tian T, Zhang K, You Q, Yan H, Zhao N, Yi X, Xu W, Su Z. 2018. PCSD : a plant  
736 chromatin state database. *Nucleic Acids Research*, 46(October 2017), 1157–1167.  
737 <https://doi.org/10.1093/nar/gkx919>
- 738 Lu Z, Marand A P, Ricci W A, Ethridge C L, Zhang X, Schmitz R J. 2019. The prevalence,  
739 evolution and chromatin signatures of plant regulatory elements. *Nature Plants*,  
740 5(December), 1250–1259. <https://doi.org/10.1038/s41477-019-0548-z>
- 741 Maher K A, Bajic M, Kajala K, Reynoso M, Pauluzzi G, West D A, Zumstein K,  
742 Woodhouse M, Bubb K, Dorrrity MW, et al. 2018. Profiling of Accessible Chromatin  
743 Regions across Multiple Plant Species and Cell Types Reveals Common Gene  
744 Regulatory Principles and New Control Modules. *The Plant Cell*, 30(January), 15–36.  
745 <https://doi.org/10.1105/tpc.17.00581>
- 746 Marand A P, Eveland A L, Kaufmann K, Springer N M. 2023. cis -Regulatory Elements  
747 in Plant Development , Adaptation , and Evolution. *Annual Review of Plant Biology*, 74,  
748 111–137. <https://doi.org/10.1146/annurev-arplant-070122-030236>

- 749 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing  
750 reads. *EMBnet.Journal*, 17(1), 10–12.
- 751 Monroe J G, Srikanth T, Carbonell-bejerano P, Becker C, Lensink M, Exposito-alonso M,  
752 Klein M, Hildebrandt J, Neumann M, Kliebenstein D, et al. 2022. Mutation bias reflects  
753 natural selection in *Arabidopsis thaliana*. *Nature*, 602(3), 101–105.  
754 <https://doi.org/10.1038/s41586-021-04269-6>
- 755 Oh Khim G G, Leary B M O, Signorelli S, Millar A H. 2022. Alternative oxidase (AOX)  
756 1a and 1d limit proline-induced oxidative stress and aid salinity recovery in *Arabidopsis*.  
757 *Plant Physiology*, 188, 1521–1536. <https://doi.org/10.1093/plphys/kiab578>
- 758 Oh Khim G G, Kumari V, Millar A H, Leary B M O. 2023. Alternative oxidase 1a and 1d  
759 enable metabolic flexibility during Ala catabolism in *Arabidopsis* Research Article. *Plant*  
760 *Physiology*, 192(4), 2958–2970. <https://doi.org/10.1093/plphys/kiad233>
- 761 Osipova E, Hecker N, Hiller M. 2019. RepeatFiller newly identifies megabases of aligning  
762 repetitive sequences and improves annotations of conserved non-exonic elements.  
763 *GigaScience*, 8, 1–10. <https://doi.org/10.1093/gigascience/giz132>
- 764 Panchy N, Lehti-shiu M, Shiu S. 2016. Evolution of Gene Duplication in Plants. *Plant*  
765 *Physiology*, 171(August), 2294–2316. <https://doi.org/10.1104/pp.16.00523>
- 766 Papathodorou I, Moreno P, Manning J, George N, Fexova S, Fonseca N A, Füllgrabe  
767 A, Green M, Huang N, Huerta L, et al. 2020. Expression Atlas update : from tissues to  
768 single cells Anja F ullgrabe. *Nucleic Acids Research*, 48(October 2019), 77–83.  
769 <https://doi.org/10.1093/nar/gkz947>
- 770 Ram F, Ryan D P, Bhardwaj V, Kilpert F, Richter A S, Heyne S, Dündar F, Manke T.  
771 2016. deepTools2 : a next generation web server for deep-sequencing data analysis.  
772 *Nucleic Acids Research*, 44(April), 160–165. <https://doi.org/10.1093/nar/gkw257>
- 773 Ren R, Wang H, Guo C, Zhang N, Zeng L, Chen Y, Hong M, Qi J. 2018. Widespread  
774 Whole Genome Duplications Contribute to Genome Complexity and Species Diversity in  
775 Angiosperms. *Molecular Plant*, 11, 414–428. <https://doi.org/10.1016/j.molp.2018.01.002>
- 776 Schmitz R J, Grotewold E, Stam M. 2022. Cis-regulatory sequences in plants : Their  
777 importance , discovery , and future challenges. *The Plant Cell*, 34, 718–741.  
778 <https://doi.org/10.1093/plcell/koab281>
- 779 Schneider C A, Rasband W S, Eliceiri K W. 2012. NIH Image to ImageJ : 25 years of  
780 Image Analysis. *Nature Methods*, 9(7), 671–675. <https://doi.org/10.1038/nmeth.2089>

- 781 Sheffield N C, Bock C. 2016. LOLA : enrichment analysis for genomic region sets and  
782 regulatory elements in R and Bioconductor. *Bioinformatics*, 32(October 2015), 587–589.  
783 <https://doi.org/10.1093/bioinformatics/btv612>
- 784 Sohn K A, Ho J W K, Djordjevic D, Jeong H H, Park P J, Kim J H. 2015. HiHMM: Bayesian  
785 non-parametric joint inference of chromatin state maps. *Bioinformatics*, 31(13), 2066–  
786 2074. <https://doi.org/10.1093/bioinformatics/btv117>
- 787 Strodtkotter I, Padmasreea K, Dinakara C, Spetha B, Niazi P S, Wojtera J, Voss I, Do  
788 PT, Nunes-Nesi A, Fernie AR, et al. 2009. Induction of the AOX1D Isoform of Alternative  
789 Oxidase in *A. thaliana* T-DNA Insertion Lines Lacking Isoform AOX1A Is Insufficient to  
790 Optimize Photosynthesis when Treated with Antimycin A. *Molecular Plant*, 2(2).  
791 <https://doi.org/10.1093/mp/ssn089>
- 792 Suarez H G, Langer B E, Ladde P, Hiller M. 2017. ChainCleaner improves genome  
793 alignment specificity and sensitivity. *Bioinformatics*, 33(January), 1596–1603.  
794 <https://doi.org/10.1093/bioinformatics/btx024>
- 795 Tan G, Polychronopoulos D, Lenhard B. 2019. CNEr : A toolkit for exploring extreme  
796 noncoding conservation. *PLoS Computational Biology*, 15((8)), 1–16.  
797 <https://doi.org/10.1371/journal.pcbi.1006940>
- 798 Tian F, Yang D, Meng Y, Jin J, Gao G. 2020. PlantRegMap : charting functional  
799 regulatory maps in plants. *Nucleic Acids Research*, 48(November 2019), 1104–1113.  
800 <https://doi.org/10.1093/nar/gkz1020>
- 801 Togninalli M, Seren Ü, Freudenthal J A, Monroe J G, Meng D, Nordborg M, Weigel D,  
802 Borgwardt K, Korte A, Grimm GD. 2020. AraPheno and the AraGWAS Catalog 2020 : a  
803 major database update including RNA-Seq and knockout mutation data for *Arabidopsis*  
804 *thaliana*. *Nucleic Acids Research*, 48(October 2019), 1063–1068.  
805 <https://doi.org/10.1093/nar/gkz925>
- 806 Valledor L, Escandón M, Meijón M, Nukarinen E, Cañal M J, Weckwerth W. 2014. A  
807 universal protocol for the combined isolation of metabolites, DNA, long RNAs, small  
808 RNAs, and proteins from plants and microorganisms. *Plant Journal*, 79(1), 173–180.  
809 <https://doi.org/10.1111/tpj.12546>
- 810 Velay F, Méteignier L-V, Laloi C. 2022. You shall not pass ! A Chromatin barrier story in  
811 plants. *Frontiers in Plant Science*, 13(September), 1–9.  
812 <https://doi.org/10.3389/fpls.2022.888102>

- 813 Vu H, Ernst J. 2022. Universal annotation of the human genome through integration of  
814 over a thousand epigenomic datasets. *Genome Biology*, 23(9), 1–37.  
815 <https://doi.org/10.1186/s13059-021-02572-z>
- 816 Yocca A E, Edger P P. 2022. Current status and future perspectives on the evolution of  
817 cis -regulatory elements in plants. *Current Opinion in Plant Biology*, 65(102139).  
818 <https://doi.org/10.1016/j.pbi.2021.102139>
- 819 Yoshida K, Noguchi K. 2009. Differential Gene Expression Profiles of the Mitochondrial  
820 Respiratory Components in Illuminated *Arabidopsis* Leaves. *Plant and Cell Physiology*,  
821 50(8), 1449–1462. <https://doi.org/10.1093/pcp/pcp090>
- 822 Zhang Y, Liu T, Meyer C A, Eeckhoute J, Johnson D S, Bernstein B E, Nusbaum C,  
823 Myers RM, Brown M, Li W, et al. 2008. Open Access Model-based Analysis of ChIP-Seq  
824 ( MACS ). *Genome Biology*, R137(9). <https://doi.org/10.1186/gb-2008-9-9-r137>
- 825 Zhao H, Yang M, Bishop J, Teng Y, Cao Y, Beall B D, Li S, Liu T, Fang Q, Fang Q, et al.  
826 2022. Identification and functional validation of super-enhancers in *Arabidopsis thaliana*.  
827 *PNAS*, 119(48), 1–11. <https://doi.org/10.1073/pnas>.
- 828 Zhao H, Sun Z, Wang J, Huang H, Kocher J, Wang L. 2014. CrossMap : a versatile tool  
829 for coordinate conversion between genome assemblies. *Bioinformatics*, 30(7), 1006–  
830 1007. <https://doi.org/10.1093/bioinformatics/btt730>
- 831 Zhao L, Xie L, Zhang Q, Ouyang W, Deng L, Guan P, Ma M, Li Y, Zhang Y, Xiao Q, et  
832 al. 2020. Integrative analysis of reference epigenomes in 20 rice varieties. *Nature  
833 Communications*, 11(2658), 1–16. <https://doi.org/10.1038/s41467-020-16457-5>

## 834 Figure legends

835 **Fig. 1. Inter-species chromatin states definition.** **Top panel:** From left to right  
836 chromatin state definitions, abbreviation, species relation, track composition (emission  
837 probability) and genome coverage based on 10 common chromatin modifications.  
838 Chromatin states with “>” indicate definitions transitioning between species. Darkblue  
839 colors in relation heatmap highlight for which species the definition is similar and columns  
840 represent *A. thaliana* (*At*), *O. sativa* (*Os*) and *Z. mays* (*Zm*), respectively. **Bottom panel:**  
841 fold enrichments over different genomic features for each state and species.

842 **Fig. 2. Inter-species chromatin states annotation with non-common chromatin  
843 modifications.** Heatmaps depicting significant ( $p < 0.05$ ) genomic overlap-enrichment  
844 (odds ratio) of inter-species states with different annotation modules. From top to bottom:  
845 non-common chromatin-binding proteins and histone modifications/variants. Chromatin

846 states with “>” indicate definitions transitioning between species. Darkblue colors in  
847 relation heatmap higlight for which species the definition is similar and rows represent *A.*  
848 *thaliana* (*At*), *O. sativa* (*Os*) and *Z. mays* (*Zm*), respectively.

849 **Fig. 3. Inter-species chromatin states description.** Each chromatin functional group  
850 is exemplified by a module with a single state (CS1 – bivalent; CS6 – active; CS10 –  
851 divergent; CS11 – heterochromatin). Each module is constituted by three alluvial  
852 diagrams describing the distribution and correspondence between gene biotypes and  
853 orthologous for each species (*A. thaliana* (*At*), *O. sativa* (*Os*) and *Z. mays* (*Zm*)). Colors  
854 denote species. Minor gene biotypes are represented by different symbols.

855 **Fig. 4. Inter-species chromatin states annotation with conservation, genetic  
856 variability and transcription factor motifs modules.** Heatmaps depicting significant  
857 ( $p < 0.05$ ) genomic overlap-enrichment (odds ratio) of inter-species states with different  
858 annotation modules. From top to bottom: conservation covered by PhastCons elements  
859 and pairwise conserved non-coding elements (CNEs), transcription factor (TF) motifs  
860 illustrated by TF binding sites (BS) according to PlantRegMap categories and genetic  
861 variability represented by significant SNPs in GWAS. Chromatin states with “>” indicate  
862 definitions transitioning between species. Darkblue colors in relation heatmap higlight for  
863 which species the definition is similar and rows represent *A. thaliana* (*At*), *O. sativa* (*Os*)  
864 and *Z. mays* (*Zm*), respectively.

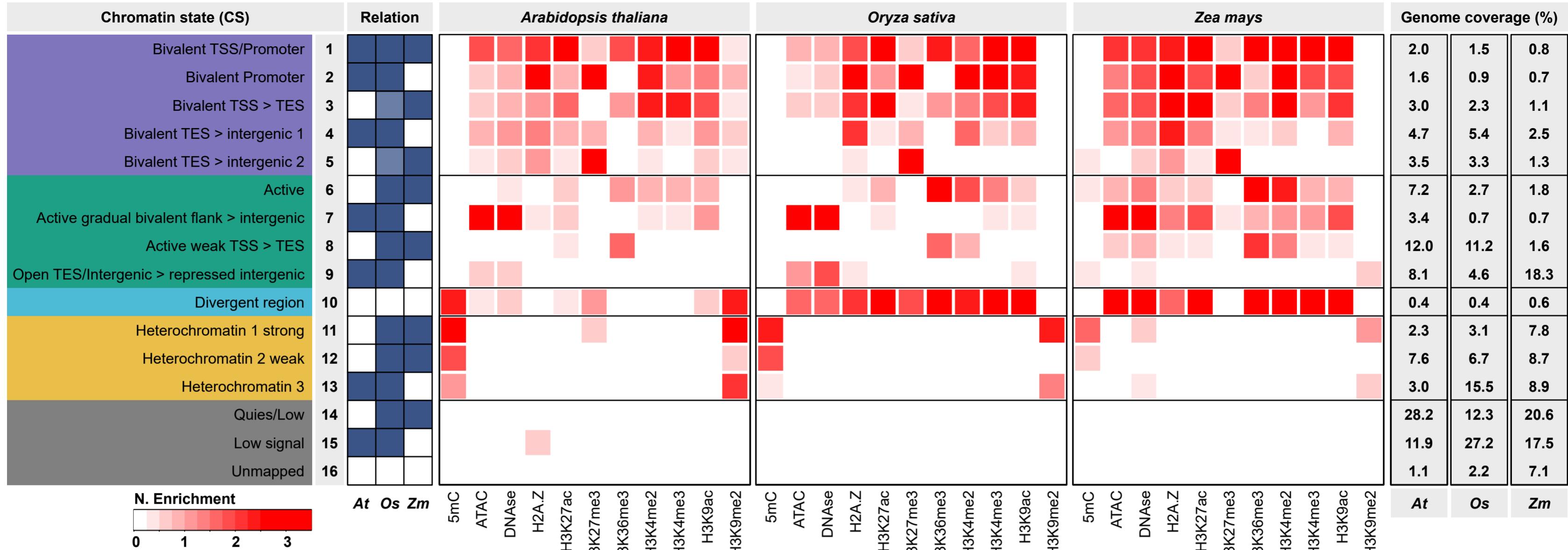
865 **Fig. 5. Predictive models of paralogs degree of functional divergence including  
866 chromatin states metrics. A)** Chromatin states metrics were obtained dividing promoter  
867 and genes in a fixed number of windows, calculating frequency and presence vectors  
868 and computing several distance and similarity coefficients between genes from the same  
869 paralog pair comparing equivalent vector types (see **Methods**). **B-E)** Results  
870 reproducing Ezoe, Shirai, and Hanada, 2021 models including CS metrics. **B)** Custom  
871 chromatin state metric (CCSM; see **Methods**) distribution of high and low diversified  
872 gene pairs. P-value, two-tailed Wilcoxon rank sum test. Numbers in parenthesis  
873 represent the number of duplicate pairs. **C)** Relative importance in explanatory variables.  
874 The relative importance was inferred based on the logistic regression algorithm. **D)**  
875 Receiver Operating Characteristic (ROC) and Precision-Recall (PR) curves in our  
876 prediction models. Colored lines indicate different generated models in six types of  
877 formula based on logistic regression algorithms using different sets of features. The area  
878 under the curve (AUC) values were calculated by the best prediction model in each  
879 formula. A perfect classification model would have AUC-ROC and AU-PRC score of 1.0;  
880 black dotted lines represent performance of random classification model, in which AUC-

ROC and AU-PRC values would be 0.5. **E**) Histogram of the inferred degree of functional divergence (DFD) in high and low duplicates of the training data. The inferred DFD was calculated for 463/111 high/low diversified pairs, respectively. The bottom 5% of the inferred high diversified DFD values were < 0.46 (i.e low DFD at 5% FDR). The top 5% of the inferred low diversified DFD values were > 0.93 (i.e high DFD at 5% FDR). Ka/Ks = protein divergence sequence rate, Re/Ks = gene expression similarity rate, FD = number of shared functional domains, GO = number of shared gene ontologies, PPI = protein-protein interactions. **F-I**) Results reproducing Cusack et al., 2021 models including CS metrics. **F**) Top 200 final selected features distribution across groups of variables for extreme-inclusive redundancy definitions without (RD4-RD9, respectively) and with (RD4C-RD9C, respetively) CS information. Numbers in parenthesis denote the median importance ranks for all the features in that group. Feature importance was determined using SVM with a linear kernel and normalized features values. Colors represent distinct redundancy definitions and features sets. RD4 (light green): extreme redundancy definition without CS information; RD4C (dark green): extreme redundancy definition with CS information; RD9 (light purple): inclusive redundancy definition without CS information; RD9C (dark purple): inclusive redundancy definition with CS information. All gene pairs in RD4/RD4C are contained in RD9/RD9C. **G**) ROC and PR curves of final SVM models for each redundancy definition/feature set. AUC values were calculated by the best prediction model in each formula. **H**) AUC-ROC and AU-PRC for the heldout tests for models built with each redundancy definition/feature set. **I**) Matrix layout for all intersections between top 200 variables in redundancy definition/feature sets, sorted by decreasing order. Dark circles in the matrix indicate sets that are part of the intersection.

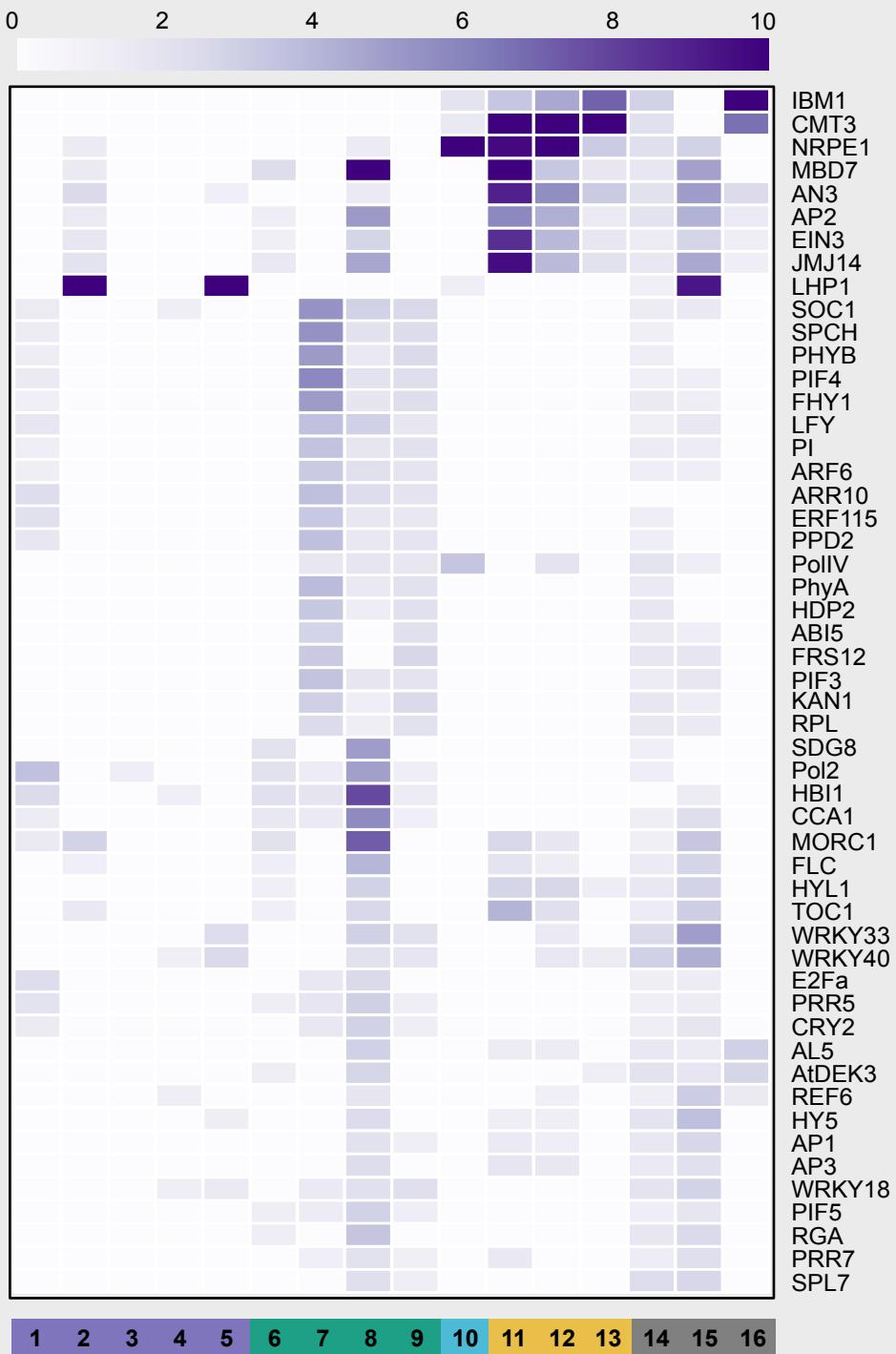
**Fig. 6. Functional genomics conservation (LECIF) score overview and downstream analyses.** This figure is constituted by 4 panels (**overview (A)**, *Arabidopsis thaliana* (**B, E, H**), *Oryza sativa* (**C, F, I**), and *Zea mays* (**D, G, J**)). **A**) Overview of the LECIF-score. Very briefly, LECIF algorithm was applied integrating epigenomic, chromatin states, whole genome alignments and transcriptomic information to obtain functional genomics conservation scores for all pairwise comparisons. These scores, together with previosuly generated resources, are stored in PlantFUNCO database to allow future applications and further hypothesis testing such as paralog functional evolution. **Arabidopsis thaliana (B, E, H), Oryza sativa (C, F, I), and Zea mays (D, G, J) panels** illustrate LECIF-score downstream analyses for *A. thaliana* (*At*), *O. sativa* (*Os*) and *Z. mays* (*Zm*), respectively. Each of this panels are divided into two sides according to the two remaining target species and three description modules: **B,C and D**) Genetic variability as genomic overlap-enrichment of GWAS significant SNPs over regions

917 divided into five bins based on LECIF scores. Black bars indicate significance ( $p < 0.05$ ).  
918 **E, F and G)** Chromatin states module with genome-wide (histogram) and state-specific  
919 (violinplot) LECIF scores distribution. Additionally, this module is covered by chromatin  
920 state similarity between high/low (percentile rank  $> 60 / < 40$ ; dark colors) and low/high  
921 (light colors) functional (LECIF) /comparative (PhyloP) genomics score regions,  
922 respectively (horizontal grouped barplot); and between regions with low, medium and  
923 high LECIF score (lineplot). Chromatin state similarity was computed using the Dice  
924 coefficient. **H, I and J)** Comparative genomics represented by boxplots showing the  
925 distribution of LECIF scores against PhatCons elements/CNEs and correlation values  
926 for LECIF versus PhyloP scores (PCC = Pearson correlation coefficient; SCC =  
927 Spearman correlation coefficient). Gray lines in boxplots denote genome-wide median  
928 and mean. Coverage (%) refers to the aligning regions overlap. PlantFUNCO DB is  
929 available at <https://rocesv.github.io/PlantFUNCO>.

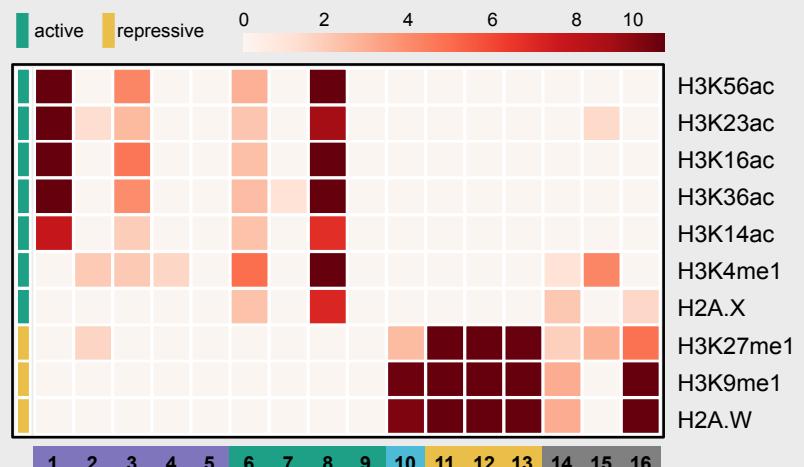
930 **Fig. 7. Experimental validation of potential high diversified AOX.** From left to right  
931 degree of functional divergence (DFD) values, genic models, chromatin states and  
932 LECIF scores, when applicable, for each of the AOX paralogs evaluated. Rows represent  
933 genotypes and columns indicate distinct conditions. For each column representative  
934 images of 5 days seedlings and cotyledons after 3,3-Diaminobenzidine (DAB) staining  
935 are displayed. The white bar represents 1 cm. Furthermore, root phenotype boxplots of  
936 root length, hypocotyl length and root:hypocotyl length ratio are presented in the bottom  
937 panel projection of the column. After two paired conditions (Control vs PEG x Heat; Mock  
938 vs Antimycin A) an additional column is added to illustrate DAB quantification intra-  
939 genotype results. The staining intensity was quantified after 32-bit gray scale  
940 transformation as: integrated density – (area selected \* mean intensity of background  
941 readings). Phenotypic differences were determined based on at least twelve biological  
942 replicates for root phenotypes and at least three biological replicates for DAB staining. A  
943 difference is considered significant with  $p < 0.05$ . “ns”:  $p > 0.05$ ; “\*”:  $p < 0.05$ ; “\*\*”:  $p <$   
944 0.01; “\*\*\*”:  $p < 0.001$ ; “\*\*\*\*”:  $p < 0.0001$ . KW = Kruskall-Wallis.



## Chromatin-binding proteins



### Histone marks

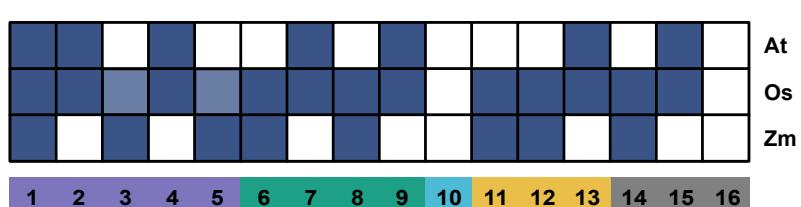


*Arabidopsis thaliana*

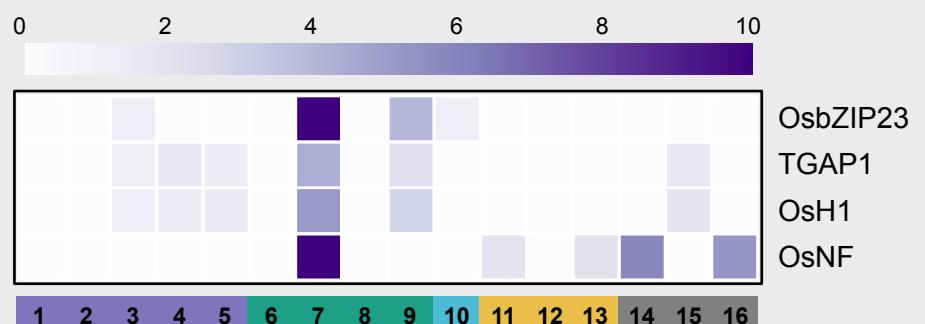
## Chromatin States

- |                               |  |
|-------------------------------|--|
| 1 Bivalent TSS/Promoter       | 4 Bivalent TSS > intergenic 2                |
| 2 Bivalent Promoter           | 5 Active                                     |
| 3 Bivalent TSS > TES          | 6 Active gradual bivalent flank > intergenic |
| 4 Bivalent TES > intergenic 1 | 7 Active weak TSS > TES                      |

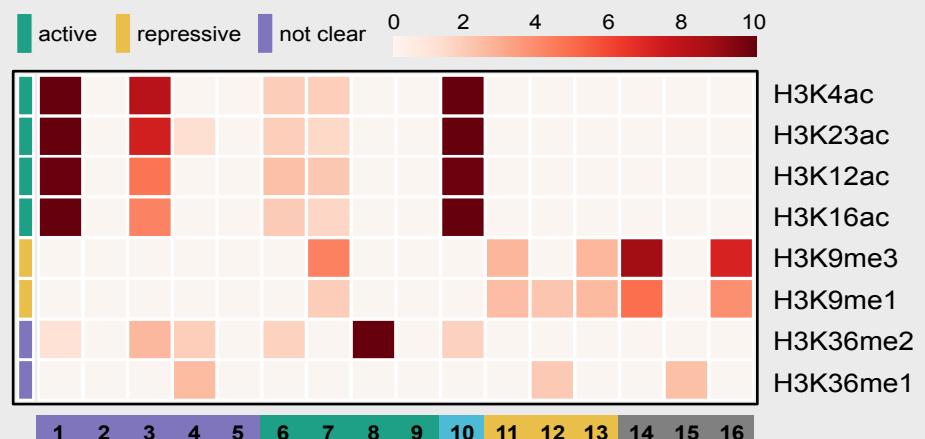
## Relation



## Chromatin-binding proteins



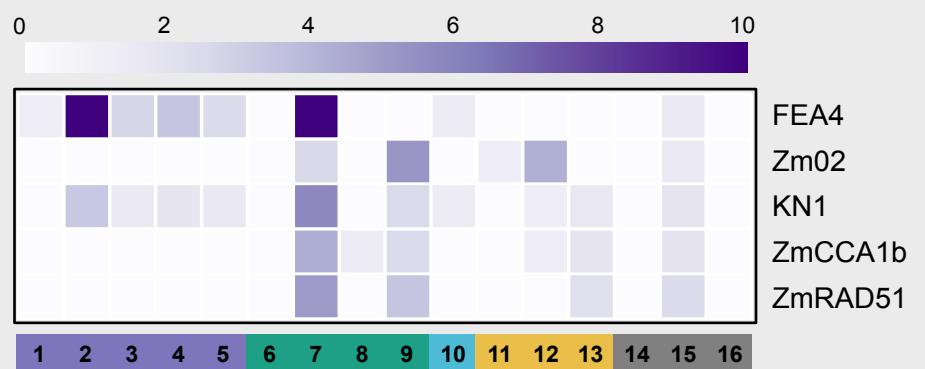
### Histone marks



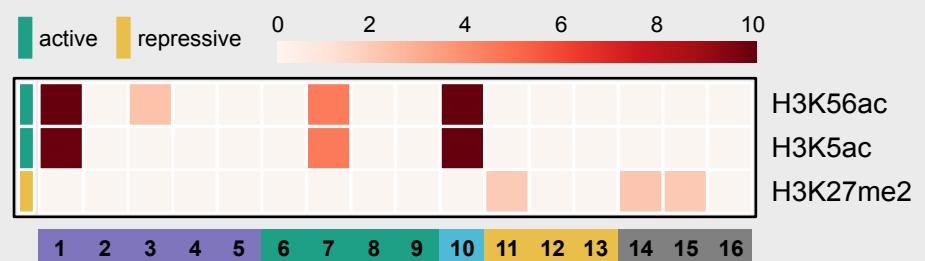
*Oryza sativa*

*Zea mays*

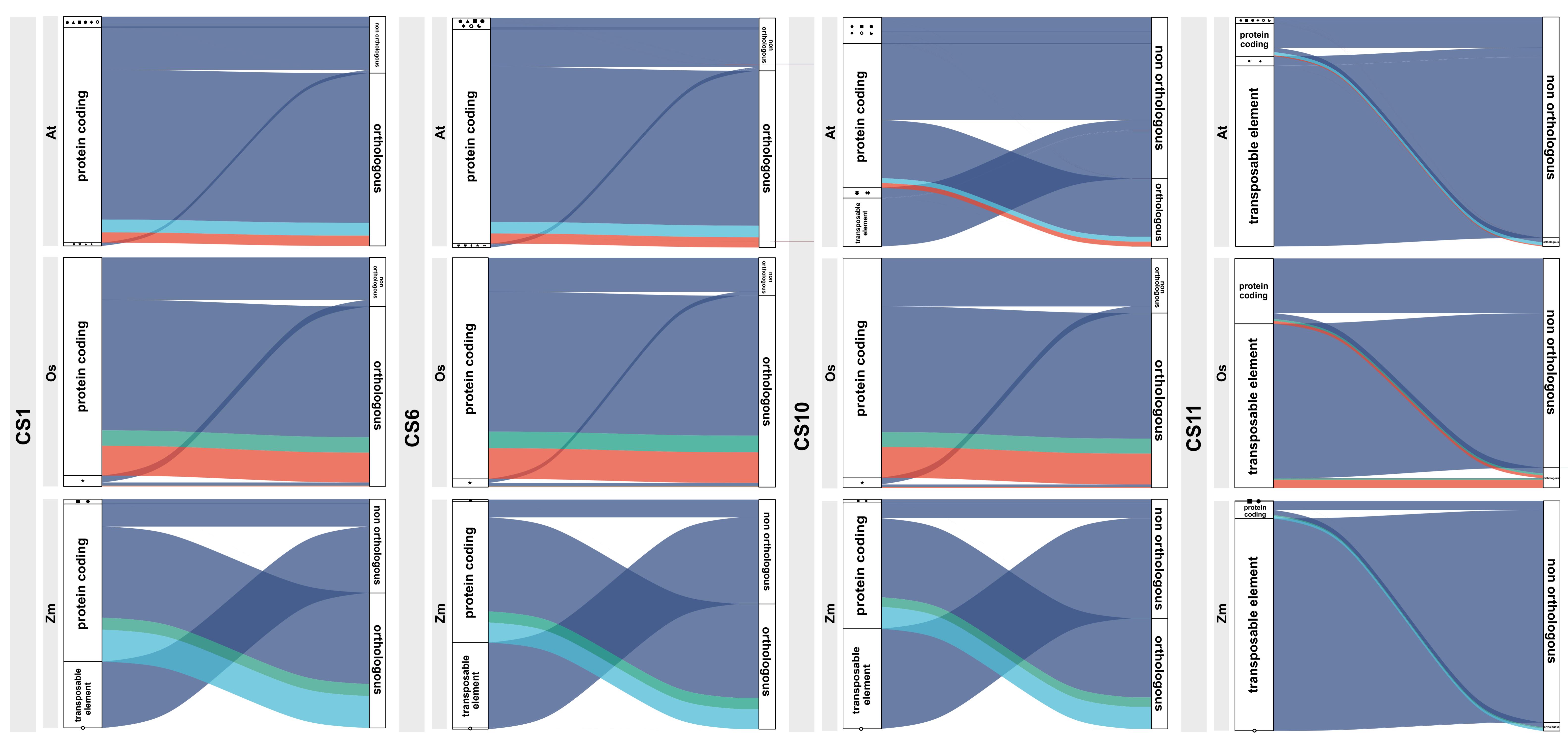
## Chromatin-binding proteins



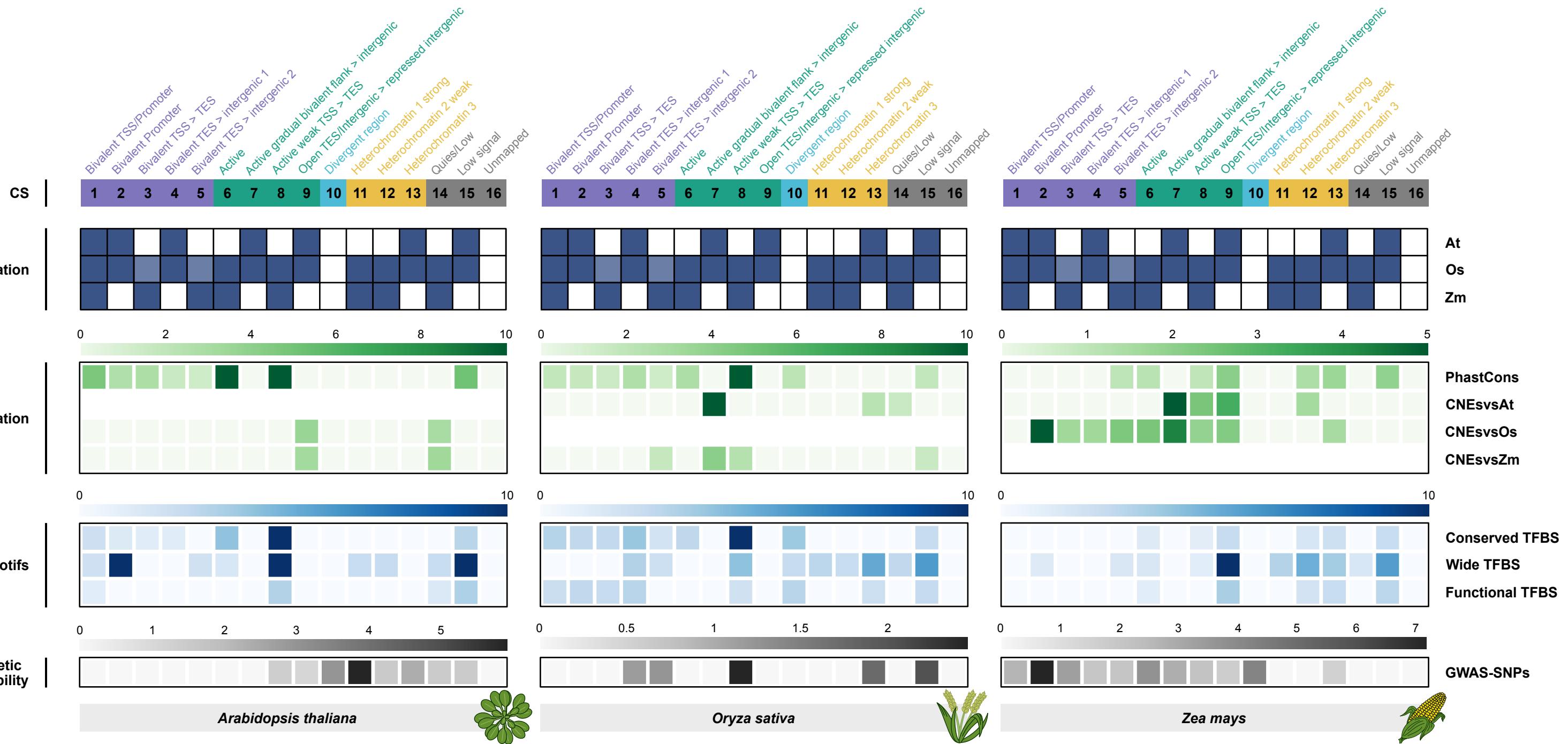
### Histone marks

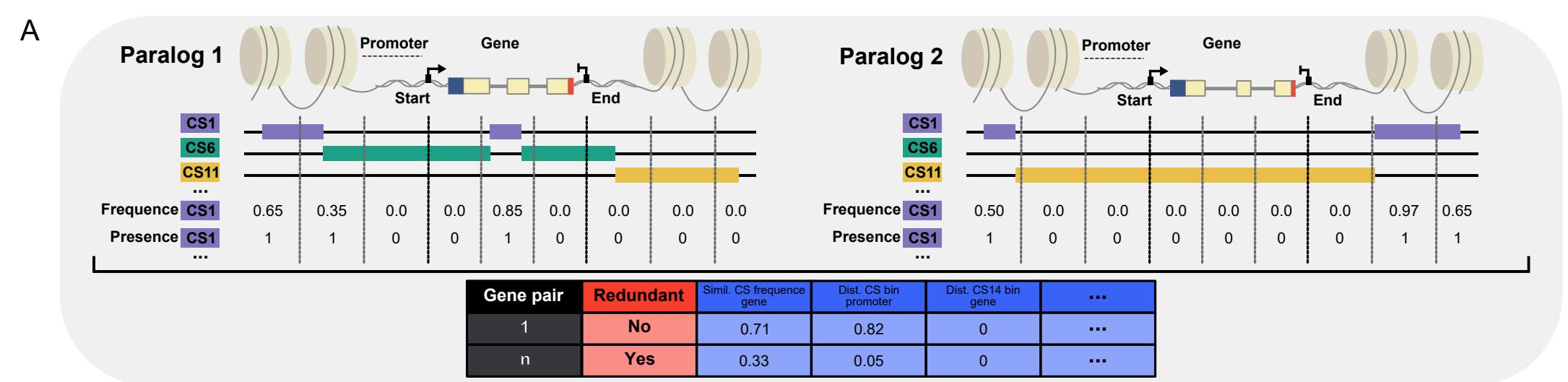


- |  |                      |
|--|----------------------|
| 8 Open TES/Intergenic > repressed intergenic | 12 Heterochromatin 3 |
| 9 Divergent region                           | 13 Quies/Low         |
| 10 Heterochromatin 1 strong                  | 14 Low signal        |
| 11 Heterochromatin 2 weak                    | 15 Unmapped          |



♦ antisense long non coding RNA   ♦ novel transcribed region   ♦ pseudogene   ▲ antisense RNA   ○ pre tRNA / tRNA   ♥ small nuclear RNA   ■ long non coding RNA / lincRNA  
 ★ transposable element   ♦ small nucleolar RNA   ● miRNA primary transcript / miRNA   ● other RNA   ↗ ribosomal RNA





Degree of Functional Divergence in Duplicates Is Associated with Distinct Roles in Plant Evolution (Ezoe et al., 2021)

