

1 **Main Manuscript for**  
2 **Repeated global adaptation across plant species**

3

4 Gabriele Nocchi<sup>1</sup>, James R. Whiting<sup>1</sup>, Samuel Yeaman<sup>1</sup>

5 1. Department of Biological Sciences, University of Calgary, Calgary, Canada T2N 1N4  
6

7 Corresponding author: Gabriele Nocchi

8 **Email:** [gabriele.nocchi@ucalgary.ca](mailto:gabriele.nocchi@ucalgary.ca)

9 **Author Contributions:** GN analysed the data and wrote the manuscript. JW performed SNP calling,  
10 developed the methods to assess pleiotropy and provided help and suggestions for the analyses. SY  
11 conceived the project, provided overall guidance, and helped with writing the manuscript.

12 **Competing Interest Statement:** NA

13 **Classification:** Biological Sciences - Evolution

14 **Keywords:** global adaptation, selective sweep, pleiotropy, repeated adaptation

15

16 **This PDF file includes:**

17 Main Text  
18 Figures 1 to 4

19

20

21

22

23

24

25

26

27

28

29

30

31

## 32 **Significance**

33 Global adaptation occurs when a species undergoes selection toward a common optimum throughout  
34 its natural range. While instances of global adaptation are widespread in the literature, there is a  
35 shortage of comparative studies aimed at understanding its genetic architecture and how it contrasts  
36 with that of local adaptation. This research compares global selective sweeps across 17 plant species  
37 to uncover the attributes of the genetic loci repeatedly involved in adaptation. We show that global  
38 adaptation tends to rely on genes with reduced pleiotropy and is characterized by increased levels of  
39 gene duplication. This finding contrasts with recent observations of increased pleiotropy in genes  
40 driving local adaptation, reflecting the opposing dynamics underlying these two evolutionary  
41 processes.

42

## 43 **Abstract**

44 Global adaptation occurs when all populations of a species undergo selection toward a common  
45 optimum. This can occur by a hard selective sweep with the emergence of a new globally  
46 advantageous allele that spreads throughout a species' natural range until reaching fixation. This  
47 evolutionary process leaves a temporary trace in the region affected, which is detectable using  
48 population genomic methods. While selective sweeps have been identified in many species, there  
49 have been few comparative and systematic studies of the genes involved in global adaptation.  
50 Building upon recent findings showing repeated genetic basis of local adaptation across independent  
51 populations and species, we asked whether certain genes play a more significant role in driving global  
52 adaptation across plant species. To address this question, we scanned the genomes of 17 plant and  
53 forest tree species to identify signals of repeated global selective sweeps. Despite the substantial  
54 evolutionary distance between the species analysed, we identified several gene families with strong  
55 evidence of repeated positive selection. These gene families tend to be enriched for reduced  
56 pleiotropy, consistent with predictions from Fisher's evolutionary model and the cost of complexity  
57 hypothesis. We also found that genes with repeated sweeps exhibit elevated levels of gene  
58 duplication. Our findings contrast with recent observations of increased pleiotropy in genes driving  
59 local adaptation, consistent with predictions based on the theory of migration-selection balance.

60

## 61 **Introduction**

62 Plant species occupy a wide range of niches, adopt very different life history strategies, and inhabit  
63 environments with drastically different biotic pressures (1). Due to their sessile nature, plants must  
64 contend with the biotic and abiotic stresses they encounter directly in their environment, therefore  
65 phenotypic plasticity and genomic adaptation are of critical importance in the plant kingdom (1). To no  
66 surprise, local adaptation has been more commonly detected in plants than in animals (2,3,4,5).

67 Local adaptation occurs when a species inhabits a heterogeneous environment with spatial variability  
68 in the optimal phenotype. This can lead to the evolution of spatially differentiated genotypes along  
69 environmental gradients that exhibit fitness trade-offs when transplanted, with local genotypes  
70 providing higher fitness than foreign (4,5,6). By contrast, global adaptation occurs when all  
71 populations of a species experience selection towards the same optimum, resulting in the gradual  
72 refinement of a trait that is advantageous throughout all the environments inhabited by a species,  
73 such as the evolution of opposable thumbs in ancestral humans (6). While both global and local  
74 adaptation involve positive selection, the tension between migration and divergent natural selection  
75 inherent in local adaptation can result in different predictions about the evolution of genetic  
76 architecture. Local adaptation will tend to involve alleles that are larger and more tightly linked than  
77 global adaptation, as such architectures are better able to resist the homogenizing effect of migration  
78 (6). However, while this is a clear theoretical prediction, it has not been extensively tested using  
79 empirical data.

80 At the molecular level, global adaptation can occur via hard or soft selective sweeps or through more  
81 subtle shifts in allele frequency at many loci (7). With a hard selective sweep (8,9,10) a new globally  
82 advantageous mutation rapidly spreads throughout a species' natural range until it reaches fixation.  
83 During this process, the affected genomic region displays a distinctive signature marked by  
84 diminished genetic diversity and increased homozygosity, a shift of the site frequency spectrum (SFS)  
85 toward low and high-frequency variants (11,12), and particular patterns of linkage disequilibrium (LD)  
86 characterized by elevated LD on each side of the selected site and decreased LD between sites  
87 located on different sides (13). This trace can be detected by scanning intraspecific genetic data, and  
88 various population genomics methods have been developed to detect sweeps (10,14,15). After  
89 fixation of the beneficial variant, new mutations and recombination in the region slowly decay the  
90 genomic signature typical of a hard selective sweep, therefore genomic scans can only detect hard  
91 sweeps within a restricted time frame (14). Global adaptation can also arise through positive selection  
92 on standing genetic variation, which tends to result in soft selective sweeps (16) or more subtle allele  
93 frequency shifts at many loci (7). When a beneficial variant is already segregating in a population  
94 before being subjected to positive selection, the selective process does not affect linked neutral  
95 polymorphisms as much as in hard-selective sweeps, making the detection of soft-sweeps much  
96 more challenging (10,16).

97 While selective sweeps have been reported in many species across both plants and animal kingdoms  
98 (17,18,19,20,21,22,23,24,25,26,27,28,29,30,31), there has been limited comparative and systematic  
99 genome-wide study of the repeatability of gene involvement in global adaptation across multiple  
100 species. On the other hand, several studies have demonstrated that the same genetic loci are  
101 observed repeatedly driving local adaptation in different populations or species  
102 (32,33,34,35,36,37,38,39,40,41). Whether adaptation is local or global, one important factor affecting  
103 the extent of repeatability is genetic redundancy. If a given adaptive trait is characterized by limited  
104 genetic redundancy (6), few loci are available to produce the mutations needed to reach the  
105 phenotypic optimum, therefore more repeatability is expected. In such cases, adaptive loci usually  
106 have large effects, a pattern that has recently emerged for several cases of local adaptation (42). On  
107 the other hand, if a trait is highly polygenic and driven by numerous alleles of small and  
108 interchangeable effect (6,43), a multitude of genotypes could potentially yield the same optimum  
109 phenotype. In such cases, fewer loci are expected to exhibit repeated contributions to adaptation in  
110 independent bouts of evolution across species or lineages.

111 As repeated genetic basis of adaptation has been identified among different populations and in  
112 closely related species, it is intriguing to assess whether repeatability is observed at greater  
113 phylogenetic distances (41, 44). It is expected that more recently diverged lineages will demonstrate a  
114 higher degree of shared gene utilization in adaptation, owing to limited functional differentiation,  
115 increased allele sharing, comparable genomic architecture, and similar life histories/adaptive  
116 strategies (41). Also, while factors such as pleiotropy, mutability, and average mutation effect size  
117 likely vary among genes and would be theoretically predicted to affect the repeatability of adaptation,  
118 there has been little systematic study of the importance of such factors in empirical datasets.

119 To address these questions, we retrieved publicly available high-quality whole-genome sequencing  
120 (WGS) data from thousands of individuals from 17 natural plant and forest trees species distributed  
121 across four continents, including woody and herbaceous angiosperms species (Fig. 1 and SI  
122 Appendix, Table S1). We employed linkage disequilibrium-based genomic scans for selective sweeps  
123 (45) within each species to identify putative genes under positive selection, and then used *PicMin* (46)  
124 to identify genes that are enriched for strong sweep signatures across multiple species. We explored  
125 the attributes of genes with repeated global sweeps by examining the pleiotropic potential and gene  
126 duplication levels of the repeatedly swept gene families identified. Our assessment of pleiotropy  
127 leveraged available gene expression data from *Arabidopsis thaliana* and *Medicago truncatula* to  
128 generate different pleiotropy measures based on gene tissue specificity and position within gene co-  
129 expression networks. We also assessed levels of gene duplication based on our orthology  
130 assignment. We contextualized these findings in view of Fisher's model of evolution and the cost of  
131 complexity hypothesis (47, 48), as well as recent theoretical work on migration-selection balance (6).  
132 To contrast our results for the architecture of global adaptation with that of local adaptation, we

133 compare our observations to results from a related study on local adaptation (49), which employed  
134 the same bioinformatic methods and incorporated 14 of the 17 datasets used in our analysis.

135

## 136 **Results**

### 137 **Datasets assembly and filtering**

138 Whole-genome sequencing data (WGS) from 17 angiosperm plants and forest trees derived from  
139 published studies was retrieved from SRA/ENA (Fig. 1 and SI Appendix, Table S1). These datasets  
140 included exclusively natural populations and ranged from closely related sister species known to  
141 hybridize (50) to species separated by up to 160 million years of evolution (Fig. 1B).

142 Single nucleotide polymorphisms (SNPs) were identified by applying a uniform pipeline to all datasets  
143 for consistency, with filtering steps to exclude singleton SNPs and individuals with high relatedness  
144 (SI Appendix, Fig. S1; see *Methods* for details). The SNP calling and filtering process resulted in  
145 between 20 and 565 individuals per dataset, and between 667,641 and 50,268,965 SNPs (Fig. 1D  
146 and SI Appendix, Table S1). We assessed population structure and found a wide range of patterns,  
147 with some species showing nearly panmictic structure across wide spatial scales (e.g. *Eucalyptus*  
148 *albens*; SI Appendix, Fig. S2) and others exhibiting substantial sub structuring, even across small  
149 spatial scales (e.g. *Helianthus argophyllus*) (Fig. 1C and SI Appendix, Fig. S2).

150

### 151 **Orthology inference**

152 The phylogenetic relationship among the genes of the 12 species with reference genome assemblies  
153 was inferred with *Orthofinder2* (51), which was used to group genes into orthology groups, or  
154 orthogroups. These included both orthologous genes, which are homologous genes separated by a  
155 speciation event, and paralogous genes, which are homologous genes diverged from a within-species  
156 duplication event. In total, *Orthofinder* assigned 376,881 genes out of 415,552 total genes to 31,582  
157 orthogroups. These included 9,521 orthogroups specific to individual species, 7,919 orthogroups with  
158 representation from all species, and 633 single-copy orthogroups. Subsequently, this assignment was  
159 refined by excluding cases where a given species had more than 10 paralogues within a single  
160 orthogroup (but retaining the orthogroup in other species with fewer than 10 paralogues). We also  
161 excluded any gene with insufficient sequencing coverage within a given species, and excluded any  
162 orthogroups with representation in fewer than 7 species. These filtering steps resulted in a final set of  
163 13,268 orthogroups for subsequent analysis, which exhibited low levels of paralogy and high species  
164 inclusion (Fig. 2). The mean and median number of paralogues per orthogroup was 1.72 and 1.67 per  
165 species respectively, while mean and median species number per orthogroup was 15.37 and 15.42  
166 (per species). Notably, *B. stricta* yielded results for significantly fewer orthogroups compared to the  
167 other species, owing to insufficient sequencing coverage of the data (52) across many genic regions  
168 (Fig. 2 and SI Appendix, Table S1). Overall, the topology of the tree inferred with *Orthofinder2* largely  
169 matched the species tree derived from *TimeTree* (<https://timetree.org/>), with the exception of *M.*  
170 *truncatula* and *E. grandis*. *Orthofinder2* placed *M. truncatula* closer to the Brassicaceae family, relative  
171 to *E. grandis* (Fig. 2D and SI Appendix, Fig. S3). It is worth noting that *TimeTree* estimates divergence  
172 between species using a simple average across published time estimates from scientific literature  
173 (53), therefore a minor mismatch with the protein sequences-based tree inferred with *Orthofinder2* is  
174 not a reason of concern.

175

### 176 **Repeated selective sweeps**

177 We used *OmegaPlus* (45) to scan for global selective sweeps within species followed by *PicMin* (46)  
178 to identify orthogroups with repeated sweep signatures across several species. After performing a  
179 False Discovery Rate (FDR) correction to the p-values from *PicMin* based on the number of

180 orthogroups tested, we detected 33 orthogroups with significant signatures of repeated sweeps, at  $q <$   
181 0.1 (Fig. 3A).

182 Our application of *PicMin* tests for an enrichment of strong sweep signatures across multiple species  
183 but does not perform a test of which species contribute to the signal. To explore any patterns in the  
184 species driving these results, we classify any gene with an empirical p-value of  $< 0.1$  as a “driving  
185 gene”. Contribution towards the signatures of repeatability in the 33 OGs varied between species,  
186 ranging from approximately 14% (*B. stricta*) driving genes to about 68% (*E. sideroxylon*) (Fig. 3B).  
187 The repeatability signal is scattered throughout the phylogeny, with low *OmegaPlus* empirical p-values  
188 present in every species tested (Fig. 3C). However, two clusters of closely related species  
189 (*Eucalyptus* spp. and *Helianthus* spp.) were top contributors based on the count of driving genes and  
190 displayed a denser heat signature indicating an enrichment of driving genes (Fig. 3B, Fig. 3C and SI  
191 Appendix, Fig. S4). Despite this, results for only a single orthogroup appeared to be driven solely by  
192 species from the *Eucalyptus* genus, OG0010509, and none was solely driven by the *Helianthus*  
193 genus (Fig. 3C and SI Appendix, Fig. S4). However, in four orthogroups (OG0017633, OG0019910,  
194 OG0022308, OG0017585), driving genes were found exclusively in multiple *Eucalyptus* species in  
195 conjunction with multiple *Helianthus* species (Fig. 3C and SI Appendix, Fig. S4).

196 It is biologically expected that closely related species might evolve similarly and therefore might be  
197 more likely to experience sweeps in the same genes or genes families, particularly if they occupy an  
198 overlapping habitat and experience the same, or similar, selective pressures (41, 54). However, it is  
199 also possible that such signatures are bioinformatic artefacts, as a single reference genome was used  
200 to call variants in the species within each of the *Helianthus* and *Eucalyptus* genera, which means that  
201 any assembly errors could yield anomalous LD signatures that might confound *OmegaPlus* in each of  
202 the species. Consequently, the same errors could systematically propagate across closely related  
203 species sharing the same reference and potentially drive some of the observed *PicMin* repeatability.

204 To test whether this was a source of error, we assessed whether phylogenetic distance correlated with  
205 the amount of overlap in the driving genes between pairs of species, which we calculated as the ratio  
206 of observed versus expected overlap, according to the expectation of a hypergeometric distribution,  
207 and found no significant correlation (*Pearson's r* = 0.05, *p-value* = 0.57; *Spearman's rho* = -0.05, *p-*  
208 *value* = 0.6) (SI Appendix, Fig. S5A). Overall, the observed overlap in driving genes between closely  
209 related species did not appear to deviate from the expectation differently than that between more  
210 distantly related species. Further assessment of the average phylogenetic distance between driving  
211 genes species showed that all the candidates, except one (OG0010509), are characterized by  
212 considerable mean phylogenetic distances between the species with driving genes (SI Appendix, Fig.  
213 S5B). This confirms that despite an apparent enrichment of low empirical p-values in *Helianthus* and  
214 *Eucalyptus* species, the overall repeatability signal is driven by genes spread throughout the analysed  
215 phylogeny.

216 To preclude spurious contributions from closely related species arising from bioinformatic artefacts,  
217 we conducted follow-up repeatability testing by running *PicMin* with only one species from each of  
218 these two genera and taking the union of significant hits (FDR  $< 0.5$ ) across runs, evaluating all nine  
219 of the 13-species combinations [1 *Helianthus* + 1 *Eucalyptus* + all other species]. While this restricted  
220 analysis cannot detect true signals of repeated sweeps within a genus, it does not suffer from the risk  
221 of being driven by bioinformatic errors propagating due to closely related species sharing the same  
222 reference genome. Despite reduced power due the overall lower number of species included, this  
223 follow-up analysis successfully retrieved 15 of the 33 originally identified candidate orthogroups (FDR  
224  $< 0.1$ ) (SI Appendix, Table S2).

225 We explored the functional characteristics of this highest confidence set of 15 orthogroups with  
226 repeated sweep signatures, using the annotations of *A. thaliana* (55) and *M. truncatula* (56) genes (SI  
227 Appendix, Table S2). These orthogroups included several genes with pivotal roles in biotic and abiotic  
228 stress responses. For instance, TPS14 (TERPENE SYNTHASE 14) is a key protein involved in  
229 terpene metabolism. Terpenes are metabolites involved in plant defence against both pathogens and  
230 herbivores, and are responsible for the attraction of beneficial organisms such as pollinators and seed  
231 dispersers (57, 58). Similarly, TGSL12 (CALLOSE SYNTHASE 3) is a protein that participates in  
232 callose metabolism, a polysaccharide synthesized in the cell wall of a variety of higher plants in

233 response to pathogens infections and abiotic stress (59). The extensively studied cytochrome P450,  
234 belongs to a large family of proteins involved in NADPH- and O<sub>2</sub>-dependent hydroxylation reactions  
235 across all domains of life and is closely linked to environmental stress response in plants (60).  
236 Notably ATERF019 (ERF019, ERF19), is a critical transcription factor involved in the response to a  
237 range of stressors, such as drought, osmotic, and oxidative stress, as well as bacterial and fungal  
238 infection (61, 62). Intriguingly, this gene has been found to negatively regulate plant resistance to  
239 *Phytophthora parasitica*, underscoring its significance in mitigating the impact of destructive plant  
240 pathogens causing significant crop losses worldwide (61). Another notable protein group identified  
241 was the Plant U-box type E3 ubiquitin ligase (PUB62, PUBs), which encompasses proteins with  
242 diverse functions in abiotic and biotic stress responses and is also associated with pollen self-  
243 incompatibility (63).

244 In addition to stress responses, several orthogroups with signatures of repeated sweeps play roles in  
245 plant growth and development. For instance, the clathrin adaptor complexes medium subunit family  
246 protein is involved in clathrin-mediated endocytosis (CME). This process regulates various aspects of  
247 plant development, such as hormone signaling, and has also been linked to response against  
248 environmental stresses (64). Interestingly, the machinery of CME has been shown to be evolutionarily  
249 conserved in plants (64). Furthermore, we identified the p300/CBP acetyltransferase-related protein,  
250 which is linked to cell differentiation, growth, and homeostasis. Notably, this protein has a 600 amino  
251 acid C-terminal segment which appears highly conserved across plants and animals, suggesting an  
252 essential role across metazoan organisms (65). Other key players in growth and development  
253 included ATSPLA2-ALPHA (PHOSPHOLIPASE A2-ALPHA), which is involved in various growth-  
254 related processes (66), and ROXY1/ROXY2 of the CC-type glutaredoxin (ROXY) family, which is a  
255 group of proteins crucial for organogenesis, particularly anther development (67). Lastly, among our  
256 top candidates was ATIPT4 (ARABIDOPSIS THALIANA ISOPENTENYLTRANSFERASE 4), a protein  
257 that plays a significant role in cytokinin biosynthesis. Cytokinins are hormones essential for regulating  
258 various aspects of plant growth and development (68).

259

## 260 **Spatial scale of adaptation**

261 Selective sweeps are typically interpreted as evidence of a mutation spreading throughout the range  
262 of an entire species (i.e. global adaptation), but they can also occur within a restricted portion of the  
263 range due to local adaptation. It is possible that some of the signatures of repeated sweeps we  
264 observed could be driven by strong local adaptation within a subsection of the species range, rather  
265 than global adaptation. To explore this possibility, we estimated  $F_{ST}$  for each SNP within each species  
266 and took the average across all SNPs within each gene. Driving genes in orthogroups with repeated  
267 sweep signatures displayed low to average  $F_{ST}$  within species, with few exceptions, suggesting that  
268 the signatures of repeated selective sweeps were driven by global, rather than local adaptation (SI  
269 Appendix, Fig. S6 and Fig. S7). Even though we observed few driving genes with high  $F_{ST}$  (8 genes  
270 out of 78 fell within the top 10% of  $F_{ST}$  values genome wide; SI Appendix, Fig. S7), it is worth noting  
271 that a mutation making a global selective sweep can generate transient genetic differentiation as  $F_{ST}$   
272 outliers, so even the relatively rare cases of high  $F_{ST}$  that we observed within species could still have  
273 been driven by global adaptation (69).

274

## 275 **Genes with repeated sweeps have low pleiotropy**

276 We explored the pleiotropic characteristics of orthogroups with repeated sweep signatures to assess  
277 the genetic architecture of global adaptation in comparison to the classical theoretical expectation  
278 based on Fisher's geometric model of universal pleiotropy (47,48), other empirical studies  
279 (32,33,34,35,36,37,38,39,70,71,72,73,74), as well as more recent findings on the architecture of local  
280 adaptation derived from a large comparative study across plant species sharing the same methods  
281 and majority of datasets, including the same SNPs sets (49).

282 To evaluate the amount of pleiotropy in genes with repeated selective sweeps, we utilized gene  
283 expression data for *A. thaliana* and *M. truncatula* genes sourced from Expression Atlas (75) and

284 ATTED-II (76). We estimated two kinds of gene property related to pleiotropy: A) tissue specificity of  
285 gene expression (77) and B) statistics describing a gene's position/importance within a gene co-  
286 expression network (78). For tissue specificity (A), we employed the  $\tau$  metric (79) of *A. thaliana* genes,  
287 which should be inversely proportional to pleiotropy; highly tissue-specific genes are likely to influence  
288 fewer processes compared to genes expressed across several tissues. Tissue specificity has been  
289 found to be positively associated with rates of molecular evolution (80). Using gene co-expression  
290 networks (B), we estimated pleiotropy based on four centrality measures: degree, strength,  
291 closeness, and betweenness (81). Node degree signifies the number of nodes connected to a gene  
292 and thus the number of co-expressed genes; node strength represents the combined weights of these  
293 connections; closeness reflects a node's capacity to interact with all other nodes, even those not  
294 directly connected to it, and hence its co-expression capability across the entire network;  
295 betweenness indicates a node's ability to act as a bridging node in the network, linking co-expression  
296 subnetworks together. Each of these measures should be positively correlated to pleiotropy  
297 (82,83,84). Genes with high centrality measures in co-expression networks are hubs; hence, their  
298 expression affects many other genes and processes, potentially imposing an evolutionary constraint  
299 as changes in them can more likely be lethal or have strong effects. Conversely, genes with low  
300 centrality measures can be considered less pleiotropic, with changes in them affecting fewer  
301 processes (82,83,84).

302 To evaluate the amount of pleiotropy in the candidate orthogroups with repeated sweep signatures,  
303 each of the above measures was calculated for each orthogroup based on the metrics described  
304 above, using data from *A. thaliana* and *M. truncatula* (see *Methods* for details). We assessed  
305 pleiotropy in the candidate set including 33 Orthogroups (FDR < 0.1), as well as within a larger set of  
306 107 orthogroups with a more relaxed false discovery rate cutoff (FDR < 0.2). While using a more  
307 relaxed FDR results in more false positives, it also allows a larger number of true positives, which  
308 improves the power to test patterns in pleiotropy for genes involved in global adaptation. Consistent  
309 with the Fisher-Orr model of evolution (47,48), both sets of results showed the same pattern:  
310 pleiotropy appeared significantly lower in the candidate orthogroups compared to random expectation  
311 (Fig. 4A, Fig. 4B). All measures used as proxies for pleiotropy, except *Medicago* betweenness,  
312 showed either a significant ( $p < 0.05$ ) association with decreased pleiotropy or tended strongly  
313 towards this direction. It is noteworthy that another study focusing on local (rather than global)  
314 adaptation found the opposite relationship with increased pleiotropy for genes driving repeated local  
315 adaptation (49), using many of the same datasets and the same methods to identify repeated  
316 adaptation and pleiotropy.

317 We further assessed pleiotropy within the more conservative set of 15 orthogroups identified by  
318 intersecting the 33 candidate OGs identified in the main analysis (FDR < 0.1) with results derived from  
319 nine additional *PicMin* omitting closely related *Eucalyptus* and *Helianthus* species. The correlation  
320 between global adaptation and decreased pleiotropy largely persisted ( $p < 0.05$ ) when the  
321 assessment was restricted to this small set, indicating a robust association between genes driving  
322 global adaptation and decreased pleiotropy (Fig. 4C).

323

## 324 **Duplications are enriched in genes with repeated sweeps**

325 Newly duplicated genes may experience selective sweeps due to positive selection if the resulting  
326 copies undergo sub- or neo-functionalization. Additionally, sub-functionalization reduces the pleiotropy  
327 of the original gene by distributing the function of the parent gene among new copies (85). This  
328 process may aid adaptation by easing the evolutionary constraints associated with a more pleiotropic  
329 ancestral gene (85,86). Considering this, we utilized the results obtained from *Orthofinder2* to  
330 investigate the involvement of gene duplications in recurrent global adaptation. We counted the  
331 number of paralogs within orthogroups as a representation of the number of duplication events and  
332 compared our candidates with repeated sweeps against a set of randomly chosen orthogroups, using  
333 the same statistical approach as the pleiotropy assessment. The 33 orthogroups with repeated  
334 sweeps (FDR < 0.1) exhibited a significant enrichment in gene duplications ( $p < 0.05$ ), with a 3.3 fold-  
335 increase relative to the overall average for all orthogroups (Fig. 4D).

336

### 337 Possible biases and artefacts

338 As *PicMin* identifies orthogroups with consistently extreme signatures across multiple species, any  
339 source of bias arising from some characteristic of a gene can also bias *PicMin* if the characteristic  
340 tends to be conserved over deep time. Gene length is one such characteristic that might contribute to  
341 bias, based on how we ran *OmegaPlus* using one scan near the centre of each gene. This analysis  
342 choice might lead to reduced detection power in longer genes, as the signal of a selective sweep near  
343 a gene edge would be more attenuated at the middle of the gene due to the more extensive decay of  
344 linkage disequilibrium, compared to shorter genes. Consequently, *OmegaPlus* might exhibit reduced  
345 sensitivity for larger genes. Our investigation revealed that this concern was unfounded, as there were  
346 generally low correlations between empirical p-values for the sweep signatures and gene length  
347 within each species (SI Appendix, Fig. S8). No consistent pattern was observed in the relationship  
348 between gene length and empirical p-value across species, as 10 species had  $r > 0$ , while 7 had  $r < 0$ ;  
349 (SI Appendix, Fig. S9). Furthermore, when assessing the average gene length in the 33 orthogroups  
350 candidates for repeated global adaptation (FDR < 0.1), we found no significant difference in mean  
351 gene length compared to random expectations. If anything, our candidate orthogroups exhibit a  
352 tendency toward increased gene length (SI Appendix, Fig. S10), counter to the expectation favouring  
353 shorter genes under this type of bias.

354 To strengthen our confidence in the results, we next assessed the candidates derived from each of  
355 three additional approaches to implementing sweep scans using *OmegaPlus*, with variations in  
356 settings (refer to *Methods* for details). All three approaches showed high correlations between the  
357 resulting *OmegaPlus* empirical p-values and the empirical p-values obtained with the main approach  
358 implemented above (*OmegaPlus* with minimum window of 500 bp) (SI Appendix, Table S3).  
359 Remarkably, all approaches yielded consistent results in the association with pleiotropy, highlighting a  
360 robust relationship between global adaptation and low pleiotropy (SI Appendix, Fig. S11). While our  
361 analysis of pleiotropy assumes that estimates of tissue specificity from *A. thaliana* are representative  
362 of patterns in other species, specificity in expression has been shown to decline slowly among  
363 orthologues (87). In addition, our findings are robust to different measures of pleiotropy (Figure 4).

364 Another potential source of bias is in the recombination landscape around a gene, as regions with  
365 lower recombination rates may harbour higher LD and lower diversity (88), which could lead to an  
366 enrichment of low *OmegaPlus* empirical p-values in these regions. It is also possible that sweep  
367 signatures may take longer to degrade in regions of low recombination. However, recombination  
368 events play a crucial role in generating the typical selective sweep linkage disequilibrium patterns  
369 recognized by *OmegaPlus*, therefore recombination must be substantial to obtain higher  $\omega$  scores  
370 (13,15). This was largely reflected in our assessment (see *Methods*), which showed higher density of  
371 lower *OmegaPlus* empirical p-values in regions with moderate to high recombination, with the  
372 exception of *H. argophyllus* (SI Appendix, Fig. S12). Overall, *OmegaPlus* demonstrated robustness to  
373 misidentifying regions with lower recombination rates as sweeps (SI Appendix, Fig. S12).

374

## 375 Discussion

376 Here, we have uncovered evidence of 33 orthogroups (FDR < 0.1) with repeated global selective  
377 sweeps in multiple species (Fig. 3C). This discovery is noteworthy given the considerable evolutionary  
378 divergence among the analysed plant species and their diverse life history strategies and growth  
379 forms. These range from the herbaceous perennial monocot *P. hallii*, to eudicots encompassing both  
380 herbaceous annuals (*A. thaliana*, *C. rubella*, *M. truncatula*, *A. tuberculatus*, *Helianthus spp.*) and  
381 herbaceous perennials (*A. halleri*), to long-lived trees such as *Betula spp.*, *Eucalyptus spp.*, and  
382 *Populus spp.* (Fig. 1). While our study stands out for its focus on global adaptation and the exploration  
383 of more extensive phylogenetic distances, other studies have identified repeated signatures of local  
384 adaptation across independent populations and species residing in similar habitats or experiencing  
385 analogous selective pressures (32,33,34,35,36,37,38,39,40,41,49). Taken together, these studies  
386 suggest that while adaptation may commonly be polygenic and involve only subtle changes in allele



387 frequency across many loci (89), it can also involve the contribution of key genes that appear to play  
388 particularly important roles, given their repeated involvement. In attempting to understand what  
389 determines the propensity of a gene to contribute to adaptation, it is important to consider both the  
390 biological characteristics of the gene and the influence of the methods used to detect selection.

391 Our set of genes with repeated sweeps includes many with demonstrated roles in abiotic and biotic  
392 stress responses in plants, as well as a few candidates with important functions related to growth and  
393 developmental processes (SI Appendix, Table S2). This echoes prior research based on rates of  
394 sequence substitution showing that protein function significantly influences these signatures of  
395 selection (90,91,92,93), as found in *Drosophila* (94,95), *Arabidopsis* (96), hominids (97), and other  
396 mammals (98,99). There is strong support for the hypothesis that host-pathogen interactions drive  
397 particularly rapid protein evolution, as immune and stress/defence response genes are often identified  
398 as targets of natural selection (90,91,92,93). Similarly, genes linked to defence responses, such as  
399 cytochrome P450 proteins (OG0001991), have been previously reported for their higher rates of  
400 sequence evolution (90,91). These findings are derived from methods based on high rates of  
401 nonsynonymous substitution that have very limited power to detect evolution that involves only a  
402 small number of impactful mutations and will only tend to capture particularly strong and persistent  
403 selection pressure that drives many recurrent sweeps. By contrast, our sweep-based methods only  
404 detect more recent adaptation, as genomic signatures of selective sweeps degrade with time (14) but  
405 are able to detect the signatures arising from individual selective sweeps that may be missed by  
406 substitution-based methods. It is interesting that methods suited to detecting selection operating at  
407 very different timescales both tend to find genes involved in biotic interactions and defence as the  
408 targets of repeated natural selection.

409 One characteristic of a gene that appears particularly important for adaptation is pleiotropy, where a  
410 single genetic locus influences multiple phenotypic traits (100). Pleiotropy is thought to constrain  
411 adaptation due to the possible detrimental effects of a mutational change affecting multiple biological  
412 processes (42,101). Fisher's model of evolution provides a mathematical representation of this effect:  
413 pleiotropy is equated with "organismal complexity" and the model proposes that the greater the  
414 number of trait dimensions, the lower the chance that a random mutation is beneficial, thus posing a  
415 reduced adaptive potential for genes with high pleiotropy (47). If adaptation involves optimizing trait  
416 values on many dimensions, there is a higher chance that genes with high pleiotropy will overshoot  
417 the optimum or cause a change in trait value that is maladaptive on at least one dimension. Kimura  
418 (102) re-assessed Fisher's model by incorporating the probability of fixation and concluded that  
419 mutations of intermediate effect would be most likely to contribute to adaptation. Later, Orr (103,104),  
420 considered the distribution of mutations fixed over an adaptive walk and found this would be  
421 approximately exponential, with alleles of larger effect fixing only earlier in the process, and later  
422 stages of evolution dominated by mutations of smaller effect, a prediction which has found empirical  
423 support in stickleback (105). All else being equal, this family of models predicts that the rate of  
424 adaptation will be slower for organisms with more traits due to the greater amount of pleiotropy (48).

425 This "cost of complexity" led to the widely held view that evolution occurs via mutations of small effect,  
426 and genes and mutations with low pleiotropy were often found as the target of parallel genetic and  
427 phenotypic evolution (42,101). Consistent with this hypothesis, our research has pinpointed genes  
428 exhibiting low pleiotropy as recurrent targets of global adaptation across multiple plant species (Fig.  
429 4). Furthermore, this evolutionary model, despite its simplicity, has also found support in previous  
430 empirical studies, such as in yeast, where mutations affecting more phenotypic traits showed higher  
431 fitness costs hence implying a negative relationship between pleiotropic and fitness effects of  
432 mutations (70). Similarly, a study on Norwegian graylings populations showed that gene pleiotropy  
433 constrains both plastic and adaptive gene expression responses and highlighted the importance of  
434 genes with low pleiotropy in evolution (71). Taken together, these findings align with Fisher's view of  
435 evolution, suggesting that functional changes that favour one process are more likely to have  
436 deleterious trade-offs on others if pleiotropy is high.

437 Further theoretical work on the importance of pleiotropy was motivated by empirical observations of  
438 substantial modularity (i.e. most genes affect few traits) along with larger per-trait mutation effects in  
439 more pleiotropic genes (106,107). These observations were incorporated into models showing that

440 organisms with intermediate levels of pleiotropy would have the fastest adaptive rates, mitigating the  
441 cost of complexity (107). These species-level models have subsequently been interpreted as  
442 predicting the greatest contribution to adaptation within a species by genes with intermediate  
443 pleiotropy (73), which is reasonable but worthy of further directed theoretical study. Consistent with  
444 this prediction, empirical studies in *A. thaliana* (73), stickleback (32), ragweed (72), and *Heliconius*  
445 butterflies (74) have found evidence of intermediate pleiotropy driving repeated adaptation.

446 Another factor that can affect the impact of pleiotropy on adaptation is the spatial pattern of natural  
447 selection. If natural selection favours different trait optima in different locations of a species range (i.e.,  
448 local adaptation), then migration will tend to counteract divergence, resulting in a general advantage  
449 for alleles of larger effect, as they can withstand the homogenizing effect of migration (6). To the  
450 extent that pleiotropic effects of a mutation are aligned with the direction of divergence in phenotypic  
451 optima, which seems common for modular traits (106,107), we would therefore expect more  
452 pleiotropic genes to contribute to local adaptation more readily, due to their larger effects overcoming  
453 the homogenizing effect of migration. By contrast, when a species adapts to a similar phenotypic  
454 optimum across its range (i.e., global adaptation), there is no tension with migration and therefore no  
455 additional advantage for alleles of larger effect (6). In accordance with this prediction (6), we found a  
456 strong association between decreased pleiotropy and genes involved in global adaptation (Fig. 4),  
457 which is also in line with Fisher's and Orr's models of evolution and with evidence of reduced  
458 pleiotropy in rapidly evolving adaptive genes reported in other studies (70,80,82,83). In further  
459 agreement with migration-selection theory predictions (6), the opposite pattern was found for genes  
460 driving local adaptation (i.e. increased pleiotropy) by another study using the same bioinformatic  
461 pipeline and statistical tests, and conducted on many of the same species, but studying signatures of  
462 local, rather than global adaptation (49). Taken together, these two studies provide the first controlled  
463 comparison of local vs. global adaptation and the importance of pleiotropy, suggesting that this  
464 interaction also play an important role in determining the propensity of genes to contribute to  
465 adaptation.

466 Finally, we explored the role of gene duplications in global sweeps, as these have been shown to  
467 possibly be a source of genetic flexibility and facilitators of adaptation, through the means of sub and  
468 neo functionalization (85,86). Indeed, candidates for convergent local adaptation to temperature in  
469 two distantly related conifers, lodgepole pine (*Pinus contorta*) and interior spruce (*Picea glauca*, *Picea*  
470 *engelmannii*, *Picea glauca* x *Picea engelmannii*), were enriched for duplicated genes (34), but this  
471 was not found in the recent study of repeated local adaptation across plants (49). In this study we  
472 found a strong association between the orthogroups with signatures of repeated global adaptation  
473 and increased duplication number, with a 3.3x-fold increase relative to the average, despite our  
474 method penalizing duplicated genes with a conservative paralogue correction per orthogroup (Fig.  
475 4D). However, we note that pleiotropy and number of duplications tend to be inversely correlated,  
476 therefore it can be tricky to separate the causative effects (SI Appendix, Table S4). Gene duplication  
477 may facilitate global adaptation by decreasing pleiotropy of a parent gene via sub-functionalization  
478 (85). Alternatively, if sub/neo-functionalization result in a new selective landscape for the novel copy of  
479 a gene, repeated sweeps would be expected as it evolves to improve this novel function. In either  
480 case, if the propensity for gene duplication is conserved over deep time, these mechanisms  
481 associated with duplication could partly explain our findings.

482

## 483 **Materials & Methods**

### 484 **Dataset selection**

485 We downloaded raw sequencing data of 17 angiosperm plants and forest trees whole-genome  
486 sequencing (WGS) datasets from SRA and ENA (Fig. 1A and SI Appendix, Table S1).

487 The choice of WGS as the sole sequencing method for inclusion was imposed by the need for high  
488 quality dense SNP data by the software used in this study, *OmegaPlus* (45). Additionally, the datasets  
489 were chosen based on the following specific criteria: they encompassed natural populations in their  
490 native habitats, comprised non-invasive and non-domesticated species, included a minimum of 20

491 unrelated individuals sampled from five or more locations, and featured a high-quality reference  
492 genome, or one of a closely related species.

493 We explored the phylogenetic relationship between reference genome species using *TimeTree*  
494 (<https://timetree.org/>). Phylogenetic distances between species were calculated based on the  
495 *TimeTree* phylogeny using the R package *ape* (function: *cophenetic*).

496

## 497 **SNP calling**

498 We applied a uniform SNP calling pipeline to all datasets for consistency. This pipeline was derived  
499 from a previous study and was selected as it optimizes the trade off between SNP quality and  
500 processing times (108). Raw *fastq* files were trimmed using *fastp* (109) with default settings. Clean  
501 reads were then aligned to reference genomes with *bwa-mem* (v0.7.17-r1188) (110), using 12 distinct  
502 reference genomes to map 17 datasets. If a species reference genome was not available, we used  
503 that of a closely related species. Three clusters of closely related species were mapped to the same  
504 reference genome (*B. pendula* and *B. platyphylla* mapped to *B. pendula*; *E. magnificata*, *E.*  
505 *sideroxylon* and *E. albens* to *E. grandis*; *H. annuus*, *H. petiolaris* and *H. argophyllus* to *H. annuus*).  
506 Following mapping, *samtools* was used to convert the alignment files from sequence alignment map  
507 (SAM) format to sorted, indexed binary alignment map (BAM) files, while discarding any alignment  
508 with mapping quality below 10 (-q 10) (111). The *MarkDuplicates* tool (112) from *Picard tools* was  
509 used to remove potential PCR duplicates and to set read groups. Indels were realigned using *GATK*  
510 *RealignerTargetCreator* followed by *GATK IndelRealigner* (113). After indel realignment, SNP calling  
511 was performed using *BCFtools mpileup*, computing genotype likelihoods based on alignments with a  
512 minimum mapping quality of 5 (-q 5), followed by *BCFtools call* to identify single nucleotide  
513 polymorphisms (SNPs) from the pileup output and generate VCFs (111). Finally, we filtered raw VCF  
514 files with *VCFtools* (114) to retain only biallelic SNPs genotyped in at least 70% of the individuals,  
515 SNPs with quality value above 30 (--minQ 30), genotype quality above 20 (--minGQ 20) and minimum  
516 read depth above 5 (--minDP 5).

517 For downstream analyses, we retained SNP present at all minor allele frequencies except singletons.  
518 This retention of frequencies, often overlooked, was critical since the method employed to detect  
519 selective sweeps relies upon evaluating the patterns of linkage disequilibrium across genic regions  
520 (45). Therefore, filtering by allele frequency would introduce a bias by distorting LD patterns and could  
521 significantly decrease the detective power of the analysis, as excess of low frequency variants  
522 constitutes the main signature of a selective sweep (115). However, singletons are often bioinformatic  
523 artifacts and distinguishing them from real mutations can be challenging, hence their exclusion (116).

524 Finally, each dataset was filtered based on genomic relatedness (Fig. 1D). This step was taken to  
525 remove closely related individuals, as cryptic relatedness can potentially distort the estimation of  
526 regions under selection similarly to how it can confound GWAS (117). We used *Plink* (--genome  
527 function) to calculate relatedness (118), and the R package *PlinkQC* for filtering (119). *PlinkQC* aims  
528 to find the minimum number of individuals to remove to keep relatedness between any pairs below a  
529 chosen threshold (119). Individuals were systematically excluded from each species dataset to ensure  
530 that no relatedness scores between pairs exceeded 0.2, effectively removing any first and second-  
531 degree relatedness.

532

## 533 **Orthology inference**

534 To assign genes to orthogroups, we first retrieved the amino acid sequences (proteomes) for each of  
535 the 12 reference species (Table S1). For each gene, we selected a primary transcript according to the  
536 longest isoform using custom scripts. Sequences were sorted by amino acid length, and each protein  
537 was given a name corresponding to the genomic coordinates of its gene. A Perl script was used to  
538 scan the proteins FASTA files. Upon encountering a duplicate sequence (i.e., a sequence with the  
539 same header), the script retained the first occurrence (the longest) and discarded subsequent  
540 occurrences. The decision is made based on whether the header has been encountered before. If it's

541 the first time encountering a particular header, the script writes both the header and sequence to the  
542 output file. Finally, *Orthofinder2* (51) was run with default settings using as input the 12 filtered  
543 proteomes, including only a single transcript per gene.

544

#### 545 **Detection of selective sweeps using OmegaPlus**

546 We used *OmegaPlus* (45) to scan each species dataset for global selective sweeps. *OmegaPlus*  
547 searches for specific LD patterns characteristic of recent hard selective sweeps and outputs the  $\omega$ -  
548 statistic (13). LD-based methods have been shown to outperform SFS-based methods in the search  
549 for hard selective sweeps resulting in higher true positive rate (TPR), and *OmegaPlus* (45) has been  
550 consistently reported as the most sensitive tool to detect potential hard selective sweeps (10,15, 115).  
551 For our analysis of repeated sweeps across multiple species using *PicMin*, it is preferable to use the  
552 most sensitive method at the cost of a higher false positive rate, in order to detect the highest number  
553 of true positive sweeps. This is not a reason of concern in comparative studies using *PicMin*, as the  
554 same false positives are unlikely to arise from independent analyses in different species (46). We  
555 extracted genes from each species' VCF and added 1000 bp flanking regions on either side to include  
556 potential promoter and regulatory regions. Subsequently, *OmegaPlus* was run on each gene using a  
557 grid size of 3, resulting in 3 measurements: one at the first SNP, one at the last SNP and one  
558 equidistant between those 2 measurements. The minimum and maximum sizes of the subregion  
559 around a position that was included in the calculation of the  $\omega$ -statistic were set to 500 and 100,000  
560 base pairs, respectively. For each gene we retained the second scan, which was expected to lie  
561 approximately near the center of genes.

562 Following *OmegaPlus*, we ranked genes within each species by converting  $\omega$ -scores to empirical p-  
563 values. An empirical p-value corresponds to the rank of a given gene's selection score relative to all  
564 other genes from that species, therefore it reflects the strength of evidence against a null hypothesis  
565 of no selection (46). Gene empirical p-values were further summarized within each orthogroup by  
566 taking the lowest empirical p-value (i.e. the strongest selection evidence for a sweep) and applying a  
567 Dunn-Sidak correction to account for the number of paralogues within the orthogroup. This method of  
568 correcting for multiple comparisons within orthogroups only represents the contribution of the member  
569 gene with the strongest sweep signature within a given species (i.e. if two genes experience sweeps  
570 within an orthogroup in a given species, only the gene with the stronger signal contributes to the test).

571 We excluded orthogroups with more than 10 paralogues within a species, as they would be heavily  
572 penalized by the Dunn-Sidak correction and suffer from low power. Furthermore, we removed  
573 orthogroups with genes from less than 7 species from the analysis, as *PicMin* sensitivity would be  
574 reduced under such conditions. Finally, after these exclusions, we re-ranked the empirical p-values  
575 within each species to ensure uniform distributions of orthogroups empirical p-values for *PicMin*.

576

#### 577 **Detection of selective sweeps using OmegaPlus: additional approaches**

578 *OmegaPlus* was tested using the methodology outlined in the preceding section, with variations in the  
579 minimum window settings to explore robustness of our results. We experimented with both 200 and  
580 1000 as minimum window sizes before ultimately opting for the intermediate value of 500. This was  
581 crucial, as research indicates that this setting can potentially significantly impact results and introduce  
582 increased stochasticity, particularly when employing very small minimum windows (115).

583 Additionally, we explored a genome-wide scan strategy, covering the entire genomes at intervals of  
584 1000 bases. Subsequently, we extracted the scans falling within genes and aggregated these  
585 measurements on a per-gene basis by calculating the average for each gene, before applying the  
586 same methodology as before. However, we deemed this latter approach unsuitable due to its bias  
587 towards shorter genes, which tended to report more extreme average  $\omega$ -scores.

588 Nevertheless, we evaluated the correlation between all different *OmegaPlus* runs per dataset and  
589 systematically examined the pleiotropy of the candidate orthogroups derived from each approach to  
590 enhance our confidence in the results.

591

## 592 **Testing for repeated global sweeps - PicMin**

593 We used *PicMin* (46) to test for repetitive selective sweeps in 13,268 orthogroups across 17 species.  
594 *PicMin* uses order statistics to perform hypothesis tests on a set of ranked values, in this case  
595 empirical p-values derived from *OmegaPlus*  $\omega$  scores, and identifies orthogroups enriched for large  
596 numbers of genes with low empirical p-values (46). Lower *OmegaPlus* empirical p-values correspond  
597 to higher  $\omega$  scores, and indicate genes with stronger evidence of selection. For each orthogroup,  
598 *PicMin* provides a p-value representing the probability of generating ranks as extreme or more  
599 extreme than the observation under the null hypothesis of random genetic drift driving the *OmegaPlus*  
600 scores within each species. The method works as follows: for an orthogroup with genes in  $n$  species  
601 or lineages (17 in this case), under the null hypothesis the  $n$  empirical p-values representing the  
602 strength of evidence for selective sweeps in each species should follow a uniform distribution. If we  
603 order the empirical p-values within the orthogroup, the theory of order statistics shows they have  
604 marginal distributions that belong to the beta distribution, which *PicMin* uses to compute one-sided p-  
605 values. If the  $x^{\text{th}}$  rank in the list of ordered empirical p-values is low relative to the beta distribution, this  
606 indicates that the  $x$  species with the lowest empirical p-values all have stronger signatures of selective  
607 sweeps than would be expected by chance. Because some genes will always have low a rank within  
608 one species, we ignore the  $x = 1^{\text{st}}$  (lowest) ranked empirical p-value and consider all remaining higher  
609 ranks, effectively conducting tests of repeatability across two or more species. *PicMin* applies a  
610 multiple comparisons correction to the minimum p-value across all  $x = \{2 \dots n\}$  contrasts, based on the  
611 methods by Tippett (120), Dunn and Sidak (121,122). This results in a final p-value that reflects the  
612 evidence that a particular orthogroup exhibits repeated adaptation. Finally, a multiple testing  
613 correction to account for the number of tested orthogroups was applied to the final per orthogroup p-  
614 values according to the Benjamini & Hochberg (1995) (123) formula, implemented in the R function  
615 *p.adjust* (method = "fdr").

616

## 617 **Population structure assessment**

618 For population structure assessment, SNPs with minor allele frequencies  $< 0.05$  were discarded  
619 (*VCFtools*) (114) and each dataset was pruned by linkage disequilibrium ( $r^2 > 0.4$ ) using the *indep-*  
620 *pairphase* function of *Plink* in windows of 50 and step of 5 (118). Population structure was explored  
621 using principal component analysis (PCA) with *Plink* and ancestry inference with *fastSTRUCTURE*  
622 (124), testing  $K$ s from 1 to 10. The representative admixture model for each dataset was determined  
623 using *fastSTRUCTURE* built-in *chooseK* function, which selects the model that maximises the log-  
624 marginal likelihood lower bound (LLBO) of the data and best explain strong population structure (124).

625 Weir & Cockerham (1984)  $F_{ST}$  (125) was calculated between the populations identified with  
626 *fastSTRUCTURE* on a per SNP site basis with the *vcfTools* function *weir-fst-pop* (114). Within each  
627 species, individuals were assigned to populations according to the best  $K$  model  $Q$  coefficient, using  
628 as threshold of inclusion  $Q > 0.9$ . Average  $F_{ST}$  per gene was calculated by taking the mean across  $F_{ST}$   
629 values of SNPs within each gene.

630

## 631 **Gene length**

632 We assessed the association between gene length and *OmegaPlus* empirical p-values by calculating  
633 their Pearson correlation in each species, and examining the consistency of any patterning across  
634 species. We tested whether the mean gene length of the orthogroups with repeated sweep signatures  
635 identified with *PicMin* differed significantly from the expectation for randomly chosen genes. This was  
636 done by taking 10,000 random orthogroup draws of the same size as the significant set (33

637 orthogroups) and calculating the mean gene length of each draw. We then compared the mean of our  
638 candidate set against this null distribution.

639

## 640 **Recombination landscape**

641 We tested the correlation between gene recombination rate and evidence for selective sweeps  
642 derived from the *OmegaPlus* analysis. To achieve this, we downloaded recombination rate data for  
643 *Arabidopsis thaliana* (126) and *Helianthus annuus* (127), corresponding to four distinct datasets in our  
644 analysis (*A. thaliana*, *H. annuus*, *H. petiolaris*, *H. argophyllus*). We then constructed density plots (in  
645 *R*) comparing *OmegaPlus* per-gene empirical p-values against the ranked average gene  
646 recombination rate.

647

## 648 **Estimation of pleiotropy**

649 We estimated the amount of pleiotropy for each gene using two characteristics based on gene  
650 expression data: A) specificity of expression across tissues (77) and B) centrality within co-expression  
651 networks (78). To measure tissue specificity (A), we obtained *Arabidopsis thaliana* tissue expression  
652 data from Expression Atlas, accession *E-MTAB-7978* (75). This dataset includes tissue expression  
653 (transcripts per million - TPM) across developmental stages, tissue types and sub-tissue types. We  
654 computed the mean TPM across all developmental stages and sub-tissue types within each tissue  
655 type, to result in the mean TPM for each of the 23 tissue types. The tissue specificity metric  $\tau$  was  
656 determined following the method by Yanai et al. (2005) (79) as:

$$657 \quad \tau = \frac{\sum_{i=1}^N (1 - x_i)}{N - 1}$$

658 where, for a given gene,  $x_i$  corresponds to the mean TPM for a given tissue type normalised by the  
659 maximum mean TPM across  $N$  tissue types.

660 We converted  $\tau$  scores into per-gene empirical p-values, treating higher  $\tau$  estimates (corresponding to  
661 higher tissue specificity) as higher empirical p-values. We then applied the same methodology used to  
662 summarize *OmegaPlus* per-gene empirical p-values into orthogroups values, retaining the minimum  
663 empirical p-value (corresponding to minimum  $\tau$  and therefore maximum pleiotropy) per orthogroup  
664 and applying a Dunn-Šidák correction to correct for the number of paralogues. Finally, we transformed  
665 per-orthogroup empirical p-values into Z-scores with a mean of 0 and standard deviation of 1 across  
666 all orthogroups.

667 This approach avoids assuming that specificity/pleiotropy is maintained across paralogues, which  
668 should be more representative than taking the mean  $\tau$  per-orthogroup. In fact, calculating the mean  $\tau$   
669 per orthogroup significantly decreased the prevalence of high  $\tau$  values in the genome-wide  
670 distribution. This suggests that tissue specificity within orthogroups varies significantly among  
671 paralogues, possibly due to neofunctionalization or sub functionalization. Additionally, we explored the  
672 same methodology in the reverse direction, considering lower  $\tau$  estimates as lower empirical p-values.  
673 However, this alternative approach did not yield any substantial differences in the results.

674 To calculate the four centrality measures of co-expression networks (B), we built two networks using  
675 co-expression data from ATTED-II for *A. thaliana* and *M. truncatula*. The ATTED-II database  
676 summarises gene co-expression data derived from RNA-seq and microarray sources in a condition-  
677 independent manner, which is given as a standardised z-score between a given pair of genes (76). To  
678 construct co-expression networks, genes were treated as nodes and z-scores as edges, where  
679 positive z-scores denote positive co-expression and vice-versa for negative z-scores. Co-expression  
680 gene tables were downloaded for both species: *A. thaliana* = Ath-u.c3-0 and *M. truncatula* = Mtr-u.c3-  
681 0. We discarded all edges with  $-5 < z < 2.33$  following the recommendations for significant  
682 negative/positive co-expression.

683 The *A. thaliana* network included 18,570 genes, while *M. truncatula* network included 17,786 genes.  
684 Networks were generated using the *igraph* package in R. Node betweenness and closeness were  
685 calculated using the *estimate\_betweenness()* and *closeness()* functions respectively. Node degree  
686 and strength were calculated as the number of edges and the sum of all edge absolute z-scores  
687 respectively. We then condensed the resulting gene centrality measures into per-orthogroup Z-scores  
688 with mean 0 and standard deviation of 1 across all orthogroups, with the same approach used for  
689 tissue specificity  $\tau$  scores and testing both directions for the initial conversion of centrality scores into  
690 empirical p-values.

691 To test whether genes with repeated sweep signatures had high/low values of pleiotropy, for both the  
692 tissue specificity (A) and centrality measures (B) we used a bootstrapping method to compare their  
693 values to those of randomly chosen genes. For each measure, we calculated the mean z-score of the  
694 candidate set for repeatability, including 33 OGs. We then performed 10,000 random draws, each  
695 comprising the same number of orthogroups as the candidate set and calculated the mean of each  
696 random draw. Finally, we assessed whether the mean of the candidate set fell within the 95%  
697 confidence interval of the means from the 10,000 random draws. For the smaller sets of candidates,  
698 we assessed pleiotropy by calculating Stouffer's Z score (128) for tissue specificity and centrality  
699 measures using the following formula:

$$Z = \frac{\sum_{i=1}^n Z_i}{\sqrt{n}}$$

700

701

## 702 Gene duplication

703 We used the output from *Orthofinder2* to count duplication events within orthogroups, including both  
704 terminal and nonterminal nodes. We then used this data to test whether the candidate set with  
705 repeated sweep signatures was significantly enriched for duplications, using the same bootstrapping  
706 method used for pleiotropy (described in the previous section).

707

## 708 Data Availability

709 The scripts for SNP calling are available at:

710 [https://github.com/GabrieleNocchi/snp\\_calling\\_bcftools\\_slurm](https://github.com/GabrieleNocchi/snp_calling_bcftools_slurm)

711 The scripts for the population structure analysis are available at:

712 [https://github.com/GabrieleNocchi/population\\_structure\\_analysis](https://github.com/GabrieleNocchi/population_structure_analysis)

713 The scripts for the main analyses are available at: <https://github.com/GabrieleNocchi/RepSweeps>

714 References and links to the genomic resources for each dataset are available in SI Appendix, Table  
715 S1.

716

## 717 Acknowledgements

718 Funding was provided by NSERC Discovery and Alberta Innovates, with computational resources and  
719 support provided by the Digital Research Alliance of Canada.

720

## 721 References

722 1. P.J. Flood, A.M. Hancock, The genomic basis of adaptation in plants. *Curr Opin Plant Biol*  
723 36, 88-94 (2017).

- 724 2. J. Hereford, A quantitative survey of local adaptation and fitness trade-offs. *Am Nat* 173,  
725 579-588 (2009).
- 726 3. E. Sanford, M.W. Kelly, Local adaptation in marine invertebrates. *Ann Rev Mar Sci* 3, 509-  
727 535 (2011).
- 728 4. O. Savolainen, M. Lascoux, J. Merilä, Ecological genomics of local adaptation. *Nat Rev*  
729 *Genet* 14, 807-820 (2013).
- 730 5. O. Savolainen, T. Pyhäjärvi, T. Knürr, Gene flow and local adaptation in trees. *Annual*  
731 *Review of Ecology, Evolution, and Systematics* 38, 595–619 (2007).
- 732 6. S. Yeaman, Evolution of polygenic traits under global vs local adaptation. *Genetics* 220,  
733 iyab134 (2022).
- 734 7. I. Höllinger, P.S. Pennings, J. Hermisson, Polygenic adaptation: From sweeps to subtle  
735 frequency shifts. *PLoS Genet* 15, e1008035 (2019).
- 736 8. J.M. Smith, J. Haigh, The hitch-hiking effect of a favourable gene. *Genet Res* 89, 391-403  
737 (2007).
- 738 9. R. Nielsen et al., Genomic scans for selective sweeps using SNP data. *Genome Res* 15,  
739 1566-1575 (2005).
- 740 10. P. Pavlidis, N. Alachiotis, A survey of methods and tools to detect recent and strong  
741 positive selection. *J Biol Res (Thessalon)* 24, 7 (2017).
- 742 11. J.M. Braverman, R.R. Hudson, N.L. Kaplan, C.H. Langley, W. Stephan, The hitchhiking  
743 effect on the site frequency spectrum of DNA polymorphisms. *Genetics* 140, 783-796  
744 (1995).
- 745 12. J.C. Fay, C.I. Wu, Hitchhiking under positive Darwinian selection. *Genetics* 155, 1405-  
746 1413 (2000).
- 747 13. Y. Kim, R. Nielsen, Linkage disequilibrium as a signature of selective sweeps. *Genetics*  
748 167, 1513-1524 (2004).
- 749 14. H. Weigand, F. Leese, Detecting signatures of positive selection in non-model species  
750 using genomic data. *Zool. J. Linn. Soc.* 184, 528-583 (2018)
- 751 15. A. Koropoulis, N. Alachiotis, P. Pavlidis, Detecting Positive Selection in Populations Using  
752 Genetic Data. *Methods Mol Biol* 2090, 87-123 (2020).
- 753 16. W. Stephan, Selective Sweeps. *Genetics* 211, 5-13 (2019).
- 754 17. H.M.T. Vy, Y.J. Won, Y. Kim, Multiple Modes of Positive Selection Shaping the Patterns of  
755 Incomplete Selective Sweeps over African Populations of *Drosophila melanogaster*. *Mol*  
756 *Biol Evol* 34, 2792-2807 (2017).
- 757 18. S. Ihle, I. Ravaoarimanana, M. Thomas, D. Tautz, An analysis of signatures of selective  
758 sweeps in natural populations of the house mouse. *Mol Biol Evol* 23, 790-797 (2006).
- 759 19. C.D. Huber, M. DeGiorgio, I. Hellmann, R. Nielsen, Detecting recent selective sweeps  
760 while controlling for mutation rate and background selection. *Mol Ecol* 25, 142-156  
761 (2016).
- 762 20. T.R. Booker, B.C. Jackson, P.D. Keightley, Detecting positive selection in the genome.  
763 *BMC Biol* 15, 98 (2017).
- 764 21. J.M. Akey et al., Population history and natural selection shape patterns of genetic  
765 variation in 132 genes. *PLoS Biol* 2, e286 (2004).
- 766 22. E.C. Andersen et al., Chromosome-scale selective sweeps shape *Caenorhabditis*  
767 *elegans* genomic diversity. *Nat Genet* 44, 285-290 (2012).
- 768 23. S. Glinka, L. Ometto, S. Mousset, W. Stephan, D. De Lorenzo, Demography and natural  
769 selection have shaped genetic variation in *Drosophila melanogaster*: a multi-locus  
770 approach. *Genetics* 165, 1269-1278 (2003).
- 771 24. Y. Kim, W. Stephan, Detecting a local signature of genetic hitchhiking along a  
772 recombining chromosome. *Genetics* 160, 765-777 (2002).
- 773 25. D.J. Orengo, M. Aguadé, Detecting the footprint of positive selection in a european  
774 population of *Drosophila melanogaster*: multilocus pattern of variation and distance to  
775 coding regions. *Genetics* 167, 1759-1766 (2004).
- 776 26. C.J. Rubin et al., Whole-genome resequencing reveals loci under selection during  
777 chicken domestication. *Nature* 464, 587-591 (2010).
- 778 27. M.S. Wang et al., Positive selection rather than relaxation of functional constraint drives  
779 the evolution of vision during chicken domestication. *Cell Res* 26, 556-573 (2016).



- 780 28. Y. Yuan et al., Selective sweep with significant positive selection serves as the driving  
781 force for the differentiation of japonica and indica rice cultivars. *BMC Genomics* 18, 307  
782 (2017).
- 783 29. Z. Zhang et al., Whole-genome resequencing reveals signatures of selection and timing  
784 of duck domestication. *Gigascience* 7, giy027 (2018).
- 785 30. J.Y. Choi et al., Natural variation in plant telomere length is associated with flowering time.  
786 *Plant Cell* 33, 1118-1134 (2021).
- 787 31. K. Wei, G.A. Silva-Arias, A. Tellier, Selective sweeps linked to the colonization of novel  
788 habitats and climatic changes in a wild tomato species. *New Phytol* 237, 1908-1921  
789 (2023).
- 790 32. D.J. Rennison, C.L. Peichel, Pleiotropy facilitates parallel adaptation in sticklebacks. *Mol*  
791 *Ecol* 31, 1476-1486 (2022).
- 792 33. L.R. Moreira, B.T. Smith, Convergent genomic signatures of local adaptation across a  
793 continental-scale environmental gradient. *Sci Adv* 9, eadd0560 (2023).
- 794 34. S. Yeaman et al., Convergent local adaptation to climate in distantly related conifers.  
795 *Science* 353, 1431-1433 (2016).
- 796 35. S. Soudi et al., Repeatability of adaptation in sunflowers: genomic regions harbouring  
797 inversions also drive adaptation in species lacking an inversion. *eLife*, 12:RP88604  
798 (2023).
- 799 36. L. Wang et al., Molecular Parallelism Underlies Convergent Highland Adaptation of Maize  
800 Landraces. *Mol Biol Evol* 38, 3567-3580 (2021).
- 801 37. H.A. Poore et al., Repeated genetic divergence plays a minor role in repeated phenotypic  
802 divergence of lake-stream stickleback. *Evolution* 77, 110-122 (2023).
- 803 38. I.S. Magalhaes et al., Intercontinental genomic parallelism in multiple three-spined  
804 stickleback adaptive radiations. *Nat Ecol Evol* 5, 251-261 (2021).
- 805 39. M. Bohutínská et al., Genomic basis of parallel adaptation varies with divergence in  
806 *Arabidopsis* and its relatives. *Proc Natl Acad Sci U.S.A* 118, e2022713118 (2021).
- 807 40. G. Montejo-Kovacevich et al., Repeated genetic adaptation to altitude in two tropical  
808 butterflies. *Nat Commun* 13, 4676 (2022).
- 809 41. M. Bohutinska, C. L. Peichel, Divergence time shapes gene reuse during repeated  
810 adaptation. *Trends Ecol. Evol.* (2023).
- 811 42. A. Martin, V. Orgogozo, The Loci of repeated evolution: a catalog of genetic hotspots of  
812 phenotypic variation. *Evolution* 67, 1235-1250 (2013).
- 813 43. H.A. Orr, J.A. Coyne, The genetics of adaptation: a reassessment. *Am Nat* 140, 725-742  
814 (1992).
- 815 44. G.L. Conte, M.E. Arnegard, C.L. Peichel, D. Schluter, The probability of genetic  
816 parallelism and convergence in natural populations. *Proc Biol Sci* 279, 5039-5047 (2012).
- 817 45. N. Alachiotis, A. Stamatakis, P. Pavlidis, OmegaPlus: a scalable tool for rapid detection of  
818 selective sweeps in whole-genome datasets. *Bioinformatics* 28, 2274-2275 (2012).
- 819 46. T.R. Booker, S. Yeaman, M.C. Whitlock, Using genome scans to identify genes used  
820 repeatedly for adaptation. *Evolution* 77, 801-811 (2023).
- 821 47. R.A. Fisher, *The Genetical Theory of Natural Selection* (Oxford, Clarendon Press, 1930).
- 822 48. H.A. Orr, Adaptation and the cost of complexity. *Evolution* 54, 13-20 (2000).
- 823 49. J. Whiting et al., Core genes driving climate adaptation in plants. Research Square  
824 [Preprint] (2023). <https://doi.org/10.21203/rs.3.rs-3434061/v1>
- 825 50. G. Nocchi et al., Genomic signals of local adaptation and hybridization in Asian white  
826 birch. *Mol Ecol* 32, 595-612 (2023).
- 827 51. D.M. Emms, S. Kelly, OrthoFinder: phylogenetic orthology inference for comparative  
828 genomics. *Genome Biol* 20, 238 (2019).
- 829 52. B. Wang et al., Ancient polymorphisms contribute to genome-wide variation by long-term  
830 balancing selection and divergent sorting in *Boechera stricta*. *Genome Biol* 20, 126  
831 (2019).
- 832 53. S. Kumar et al., TimeTree 5: An Expanded Resource for Species Divergence Times. *Mol*  
833 *Biol Evol* 39, msac174 (2022).
- 834 54. D.L. Stern, The genetic causes of convergent evolution. *Nat Rev Genet* 14, 751-764  
835 (2013).

- 836 55. T. Z. Berardini et al., The Arabidopsis information resource: Making and mining the "gold  
837 standard" annotated reference plant genome. *Genesis* 53, 474-485 (2015).
- 838 56. F. Cunningham et al., Ensembl 2022. *Nucleic Acids Res* 50, D988-D995 (2022).
- 839 57. D. Tholl, S. Lee, Terpene Specialized Metabolism in Arabidopsis thaliana. *Arabidopsis*  
840 *Book* 9, e0143 (2011).
- 841 58. F. Zhou, E. Pichersky, The complete functional characterisation of the terpene synthase  
842 family in tomato. *New Phytol* 226, 1341-1360 (2020).
- 843 59. N. Li et al., The multifarious role of callose and callose synthase in plant development and  
844 environment interactions. *Front Plant Sci* 14, 1183402 (2023).
- 845 60. B. A. Pandian, R. Sathishraj, M. Djanaguiraman, P. V. V. Prasad, M. Jugulam, Role of  
846 Cytochrome P450 Enzymes in Plant Stress Response. *Antioxidants* (Basel) 9, 454  
847 (2020).
- 848 61. W. Lu et al., The Arabidopsis thaliana gene AtERF019 negatively regulates plant  
849 resistance to Phytophthora parasitica by suppressing PAMP-triggered immunity. *Mol Plant*  
850 *Pathol* 21, 1179-1193 (2020).
- 851 62. P. Y. Huang et al., NINJA-associated ERF19 negatively regulates Arabidopsis pattern-  
852 triggered immunity. *J Exp Bot* 70, 1033-1047 (2019).
- 853 63. M. Trujillo, News from the PUB: plant U-box type E3 ubiquitin ligases. *J Exp Bot* 69, 371-  
854 384 (2018).
- 855 64. S. Di Rubbo et al., The clathrin adaptor complex AP-2 mediates endocytosis of  
856 brassinosteroid insensitive1 in Arabidopsis. *Plant Cell* 25, 2986-2997 (2013).
- 857 65. L. Bordoli, M. Netsch, U. Lüthi, W. Lutz, R. Eckner, Plant orthologs of p300/CBP:  
858 conservation of a core domain in metazoan p300/CBP acetyltransferase-related proteins.  
859 *Nucleic Acids Res* 29, 589-597 (2001).
- 860 66. H.J. Kim et al., Endoplasmic reticulum- and Golgi-localized phospholipase A2 plays  
861 critical roles in Arabidopsis pollen development and germination. *Plant Cell* 23, 94-110  
862 (2011).
- 863 67. S. Xing, S. Zachgo, ROXY1 and ROXY2, two Arabidopsis glutaredoxin genes, are  
864 required for anther development. *Plant J* 53, 790-801 (2008).
- 865 68. J. Sun et al., The Arabidopsis AtIPT8/PGA22 gene encodes an isopentenyl transferase  
866 that is involved in de novo cytokinin biosynthesis. *Plant Physiol* 131, 167-176 (2003).
- 867 69. T.R. Booker, S. Yeaman, M. C. Whitlock, Global adaptation complicates the interpretation  
868 of genome scans for local adaptation. *Evol Lett* 5, 4-15 (2021).
- 869 70. T.F. Cooper, E. A. Ostrowski, M. Travisano, A negative relationship between mutation  
870 pleiotropy and fitness effect in yeast. *Evolution* 61, 1495-1499 (2007).
- 871 71. S. Papakostas et al., Gene pleiotropy constrains gene expression changes in fish  
872 adapted to different thermal conditions. *Nat. Commun* 5, 4071 (2014).
- 873 72. T. Hämälä, A. J. Gorton, D. A. Moeller, P. Tiffin, Pleiotropy facilitates local adaptation to  
874 distant optima in common ragweed (*Ambrosia artemisiifolia*). *PLoS Genet* 16, e1008707  
875 (2020).
- 876 73. L. Frachon et al., Intermediate degrees of synergistic pleiotropy drive adaptive evolution  
877 in ecological time. *Nat Ecol Evol* 1, 1551-1561 (2017).
- 878 74. J.J. Lewis et al., Parallel evolution of ancient, pleiotropic enhancers underlies butterfly  
879 wing pattern mimicry. *Proc Natl Acad Sci U.S.A.* 116, 24174-24183 (2019).
- 880 75. I. Papatheodorou et al., Expression Atlas: gene and protein expression across multiple  
881 studies and organisms. *Nucleic Acids Res* 46, D246-D251 (2018).
- 882 76. T. Obayashi, H. Hibara, Y. Kagaya, Y. Aoki, K. Kinoshita, ATTED-II v11: A Plant Gene  
883 Coexpression Database Using a Sample Balancing Technique by Subagging of Principal  
884 Components. *Plant Cell Physiol* 63, 869-881 (2022).
- 885 77. K. Watanabe et al., A global overview of pleiotropy and genetic architecture in complex  
886 traits. *Nat Genet* 51, 1339-1348 (2019).
- 887 78. S.R. Proulx, D.E. Promislow, P.C. Phillips, Network thinking in ecology and evolution.  
888 *Trends Ecol Evol* 20, 345-353 (2005).
- 889 79. I. Yanai et al., Genome-wide midrange transcription profiles reveal expression level  
890 relationships in human tissue specification. *Bioinformatics* 21, 650-659 (2005).

- 891 80. K.L. Mack, M. Phifer-Rixey, B. Harr, M. W. Nachman, Gene Expression Networks Across  
892 Multiple Tissues Are Associated with Rates of Molecular Evolution in Wild House Mice.  
893 *Genes* (Basel) 10, 225 (2019).
- 894 81. F.J. Azuaje, Selecting biologically informative genes in co-expression networks with a  
895 centrality score. *Biol Direct* 9, 12 (2014).
- 896 82. M. W. Hahn, A.D. Kern, Comparative genomics of centrality and essentiality in three  
897 eukaryotic protein-interaction networks. *Mol Biol Evol* 22, 803-806 (2005).
- 898 83. N. Mähler et al., Gene co-expression network connectivity is an important determinant of  
899 selective constraint. *PLoS Genet* 13, e1006402 (2017).
- 900 84. K.C. Wollenberg Valero, Aligning functional network constraint to evolutionary outcomes.  
901 *BMC Evol Biol* 20, 58 (2020).
- 902 85. J. A. Birchler, H. Yang, The multiple fates of gene duplications: Deletion,  
903 hypofunctionalization, subfunctionalization, neofunctionalization, dosage balance  
904 constraints, and neutral variation. *Plant Cell* 34, 2466-2474 (2022).
- 905 86. S. Ohno, Evolution by Gene Duplication (Springer-Verlag, 1970).
- 906 87. N. Kryuchkova-Mostacci, M. Robinson-Rechavi, Tissue-Specificity of Gene Expression  
907 Diverges Slowly between Orthologs, and Rapidly between Paralogs. *PLoS Comput Biol*  
908 12, e1005274 (2016).
- 909 88. W. Stephan, Genetic hitchhiking versus background selection: the controversy and its  
910 implications. *Philos Trans R Soc Lond B Biol Sci* 365, 1245-1253 (2010).
- 911 89. N. Barghi, J. Hermisson, C. Schlötterer, Author Correction: Polygenic adaptation: a  
912 unifying framework to understand positive selection. *Nat Rev Genet* 21, 782 (2020).
- 913 90. A.F. Moutinho, F.F. Trancoso, J.Y. Duthel, The Impact of Protein Architecture on Adaptive  
914 Evolution. *Mol Biol Evol* 36, 2013-2028 (2019).
- 915 91. A. F. Moutinho, T. Bataillon, J.Y. Duthel, Variation of the adaptive substitution rate  
916 between species and within genomes. *Evolutionary Ecology* 34, 315–338 (2019).
- 917 92. A.F. Moutinho, A. Eyre-Walker, J.Y. Duthel, Strong evidence for the adaptive walk model  
918 of gene evolution in *Drosophila* and *Arabidopsis*. *PLoS Biol* 20, e3001775 (2022).
- 919 93. T.B. Sackton, Studying Natural Selection in the Era of Ubiquitous Genomes. *Trends*  
920 *Genet* 36, 792-803 (2020).
- 921 94. T.B. Sackton et al., Dynamic evolution of the innate immune system in *Drosophila*. *Nat*  
922 *Genet* 39, 1461-1468 (2007).
- 923 95. D.J. Obbard, J. J. Welch, K. W. Kim, F. M. Jiggins, Quantifying adaptive evolution in the  
924 *Drosophila* immune system. *PLoS Genet* 5, e1000698 (2009).
- 925 96. T. Slotte et al., Genomic determinants of protein evolution and polymorphism in  
926 *Arabidopsis*. *Genome Biol Evol* 3, 1210-1219 (2011).
- 927 97. R. Nielsen et al., A scan for positively selected genes in the genomes of humans and  
928 chimpanzees. *PLoS Biol* 3, e170 (2005).
- 929 98. C. Kosiol et al., Patterns of positive selection in six Mammalian genomes. *PLoS Genet* 4,  
930 e1000144 (2008).
- 931 99. F.J. Alberto et al., Convergent genomic signatures of domestication in sheep and goats.  
932 *Nat. Commun* 9, 813 (2018).
- 933 100. F.W. Stearns, One hundred years of pleiotropy: a retrospective. *Genetics* 186, 767-773  
934 (2010).
- 935 101. D.L. Stern, V. Orgogozo, The loci of evolution: how predictable is genetic evolution.  
936 *Evolution* 62, 2155-2177 (2008).
- 937 102. Kimura M., The neutral theory of molecular evolution. Cambridge: Cambridge University  
938 Press, (1983).
- 939 103. H.A. Orr, The population genetics of adaptation: the distribution of factors fixed during  
940 adaptive evolution. *Evolution* 52, 935-949 (1998).
- 941 104. H.A. Orr, The evolutionary genetics of adaptation: a simulation study. *Genet Res* 74, 207-  
942 214 (1999).
- 943 105. S.M. Rogers et al., Genetic signature of adaptive peak shift in threespine stickleback.  
944 *Evolution* 66, 2439-2450 (2012).
- 945 106. G.P. Wagner, J. Zhang, The pleiotropic structure of the genotype-phenotype map: the  
946 evolvability of complex organisms. *Nat Rev Genet* 12, 204-213 (2011).

- 947 107. Z. Wang, B.Y. Liao, J. Zhang, Genomic patterns of pleiotropy and the evolution of  
948 complexity. *Proc Natl Acad Sci U.S.A.* 107, 18034-18039 (2010).
- 949 108. R.J. Jasper et al., Evaluating the accuracy of variant calling methods using the frequency  
950 of parent-offspring genotype mismatch. *Mol Ecol Resour* 22, 2524-2533 (2022).
- 951 109. S. Chen, Y. Zhou, Y. Chen, J. Gu, fastp: an ultra-fast all-in-one FASTQ preprocessor.  
952 *Bioinformatics* 34, i884-i890 (2018).
- 953 110. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler transform.  
954 *Bioinformatics* 25, 1754-1760 (2009).
- 955 111. P. Danecek et al., Twelve years of SAMtools and BCFtools. *Gigascience* 10, giab008  
956 (2021).
- 957 112. "Picard Toolkit." Broad Institute, GitHub Repository. <http://broadinstitute.github.io/picard/>;  
958 Broad Institute (2018).
- 959 113. M.A. DePristo et al., A framework for variation discovery and genotyping using next-  
960 generation DNA sequencing data. *Nat Genet* 43, 491-498 (2011).
- 961 114. P. Danecek et al., The variant call format and VCFtools. *Bioinformatics* 27, 2156-2158  
962 (2011).
- 963 115. N. Alachiotis, P. Pavlidis, Scalable linkage-disequilibrium-based selective sweep  
964 detection: a performance guide. *Gigascience* 5, 7 (2016).
- 965 116. R.C. Edgar, H. Flyvbjerg, Error filtering, pair assembly and error correction for next-  
966 generation sequencing reads. *Bioinformatics* 31, 3476-3482 (2015).
- 967 117. E. Uffelmann et al., Genome-wide association studies. *Nat. Rev. Methods Primers* 1, 59  
968 (2021).
- 969 118. C.C. Chang et al., Second-generation PLINK: rising to the challenge of larger and richer  
970 datasets. *Gigascience* 4, 7 (2015).
- 971 119. H.V. Meyer, plinkQC: Genotype quality control in genetic association studies. Zenodo.  
972 <https://doi.org/10.5281/zenodo.3934294> (2020).
- 973 120. L. Tippett, The methods of statistics: An introduction mainly for workers in the biological  
974 sciences. (Williams & Norgate Ltd., 1931).
- 975 121. O.J. Dunn, Estimation of the means of dependent variables. *Ann. Math. Stat.* 29, 1095–  
976 1111 (1958).
- 977 122. Z. Sidak, Rectangular confidence regions for the means of multivariate normal  
978 distributions. *J. Am. Stat. Assoc.* 62, 626 (1967).
- 979 123. Y. Benjamini, Y. Hochberg, Controlling the false discovery rate: A practical and powerful  
980 approach to multiple testing. *J. R. Stat.* 57, 289–300 (1995).
- 981 124. A. Raj, M. Stephens, J. K. Pritchard, fastSTRUCTURE: variational inference of population  
982 structure in large SNP data sets. *Genetics* 197, 573-589 (2014).
- 983 125. B. S. Weir, C. C. Cockerham, Estimating f-statistics for the analysis of population  
984 structure. *Evolution* 38, 1358-1370 (1984).
- 985 126. B. A. Rowan et al., An Ultra High-Density Arabidopsis thaliana Crossover Map That  
986 Refines the Influences of Structural Variation and Epigenetic Features. *Genetics* 213,  
987 771-787 (2019).
- 988 127. K. Huang et al., Mutation Load in Sunflower Inversions Is Negatively Correlated with  
989 Inversion Heterozygosity. *Mol Biol Evol* 39, msac101 (2022).
- 990 128. S. Stouffer, L. DeVinney, E. Suchmen, The American Soldier: Adjustment During Army  
991 Life, Vol. 1. (Princeton University Press, Princeton, NJ, 1949).

992

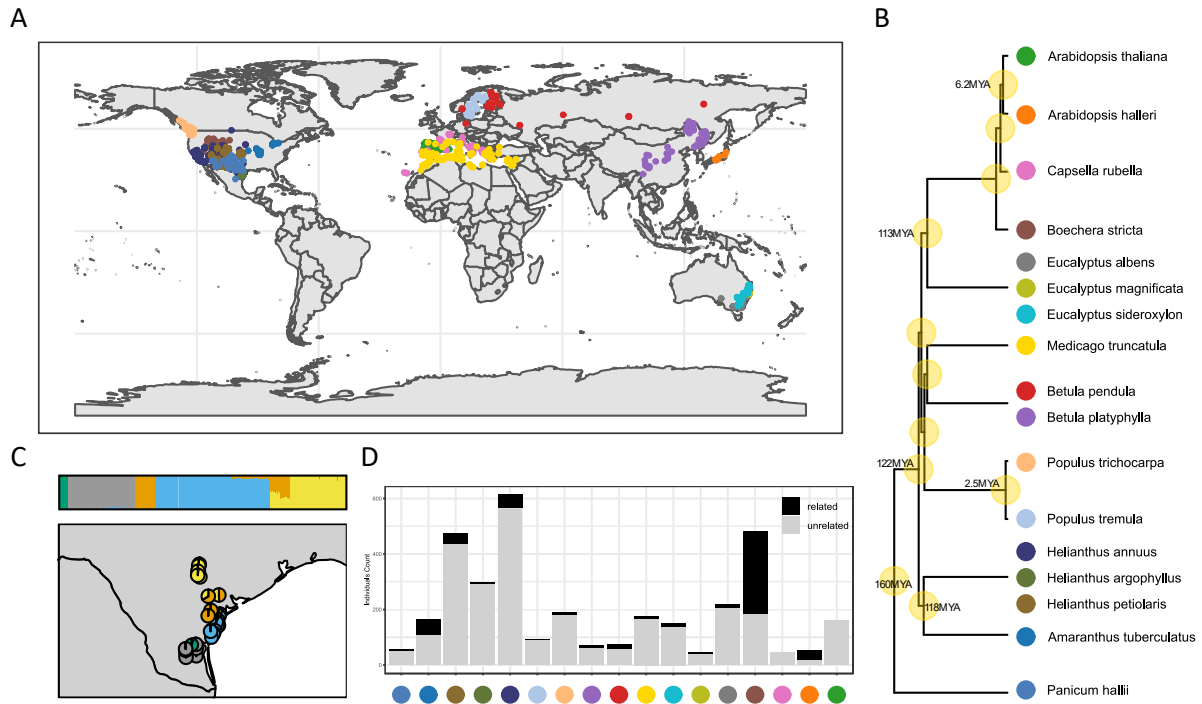
993

994

995

996

997 **Figures and Tables**



998

999 **Fig. 1.** Geographical distribution and relatedness within and among the study species. (A) Sampling  
 1000 locations of the 17 datasets included in the study. (B) Time-calibrated phylogenetic tree (retrieved  
 1001 from <https://timetree.org/>) of the 17 datasets, based on 12 reference species. (C) *fastSTRUCTURE*  
 1002 ancestry pie plot ( $K = 5$ ) of the *Helianthus argophyllus* dataset in Texas (USA), showing substantial  
 1003 sub structuring. (D) Relatedness filtering summary bar plot by dataset. Datasets labelled by  
 1004 corresponding colour from B.

1005

1006

1007

1008

1009

1010

1011

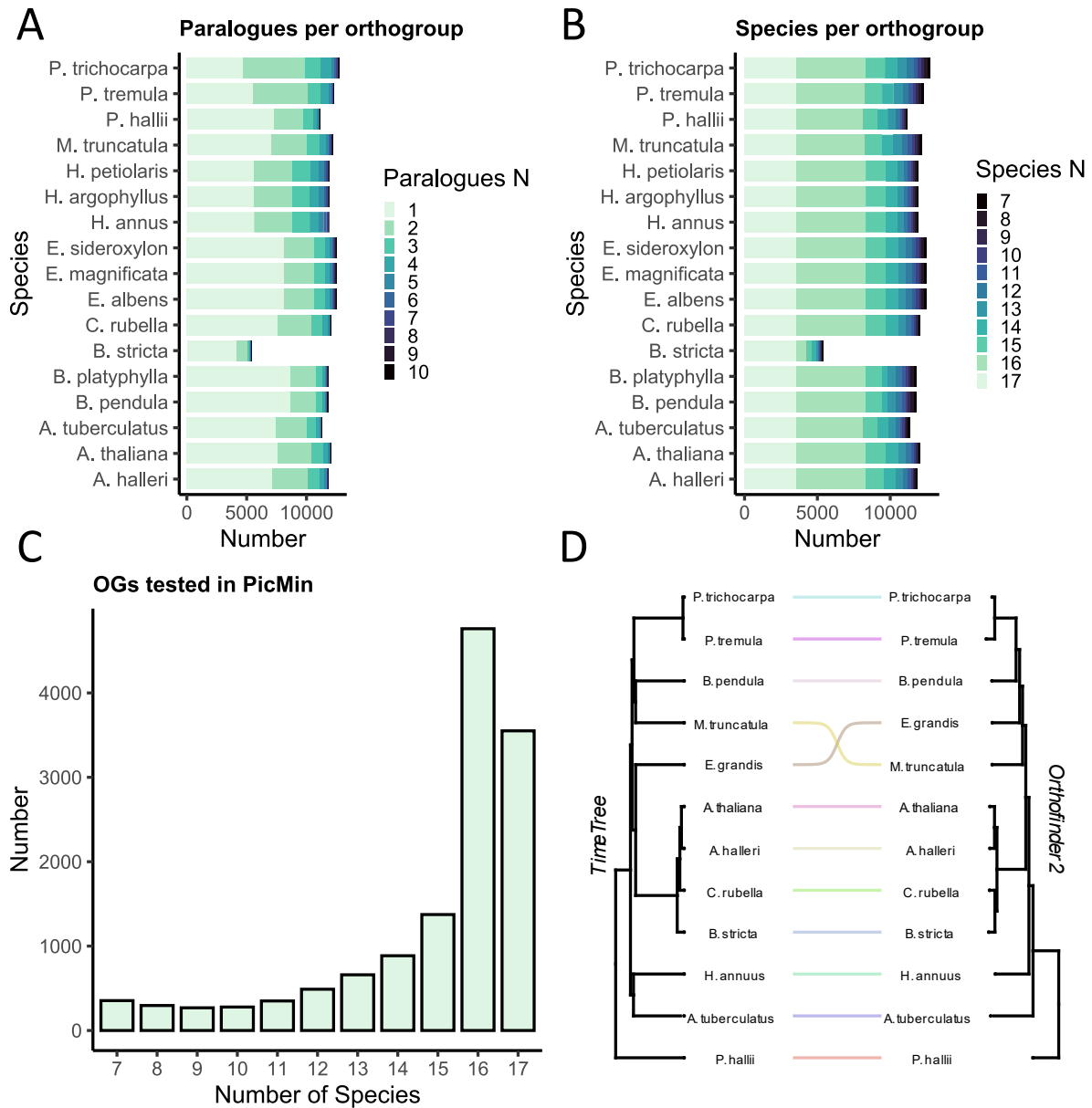
1012

1013

1014

1015

1016



1017

1018 **Fig. 2.** Orthology assignment summary of the final set of 13,268 orthogroups. (A) Bar plot of the  
 1019 number of paralogues per orthogroup for each species. (B) Bar plot of the number of species included  
 1020 in each orthogroup, for the orthogroups of each species. (C) Distribution of the 13,268 tested  
 1021 orthogroups across species number. (D) Comparison between the *TimeTree* and *Orthofinder2*  
 1022 phylogenies, each based on 12 reference species.

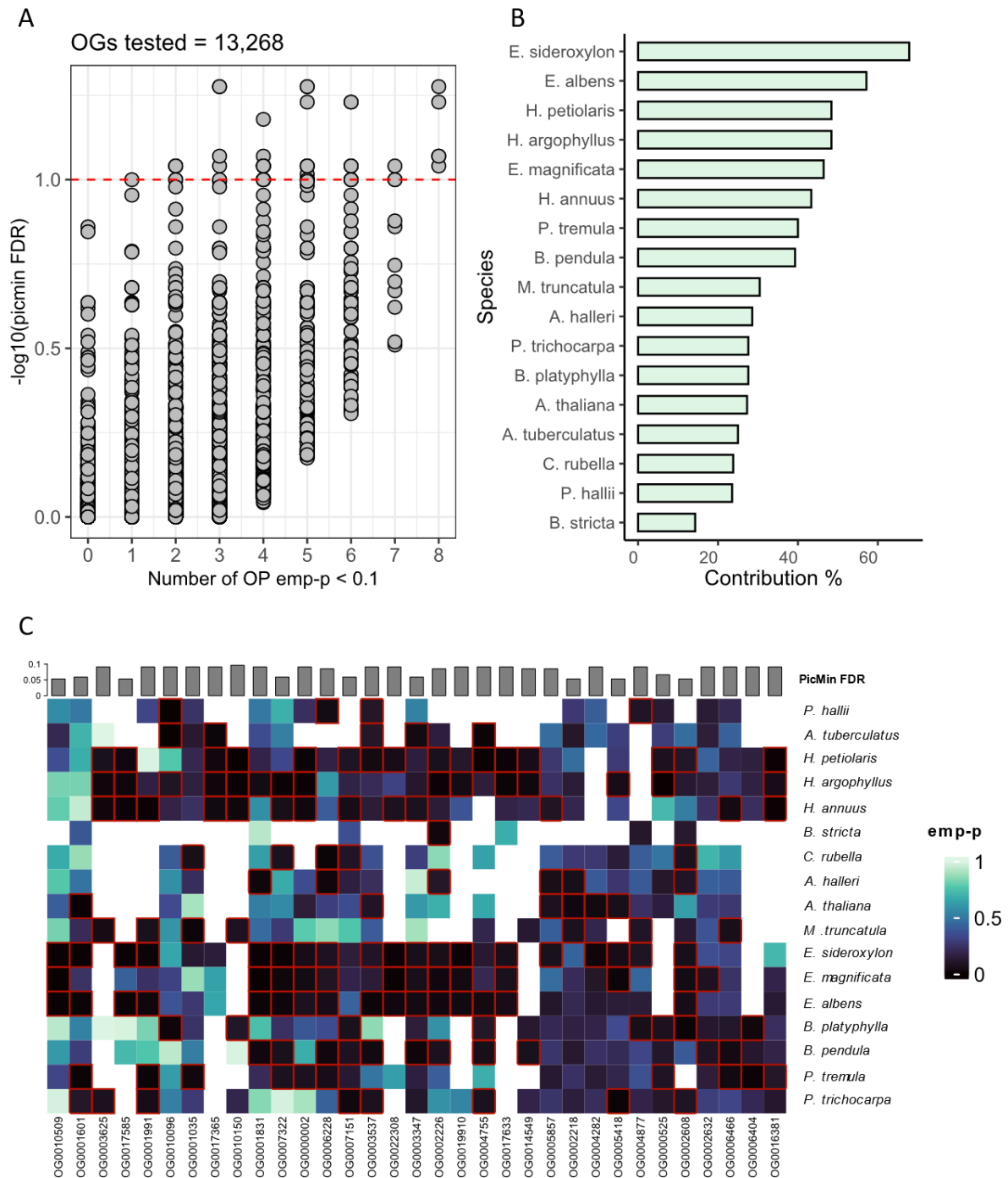
1023

1024

1025

1026

1027

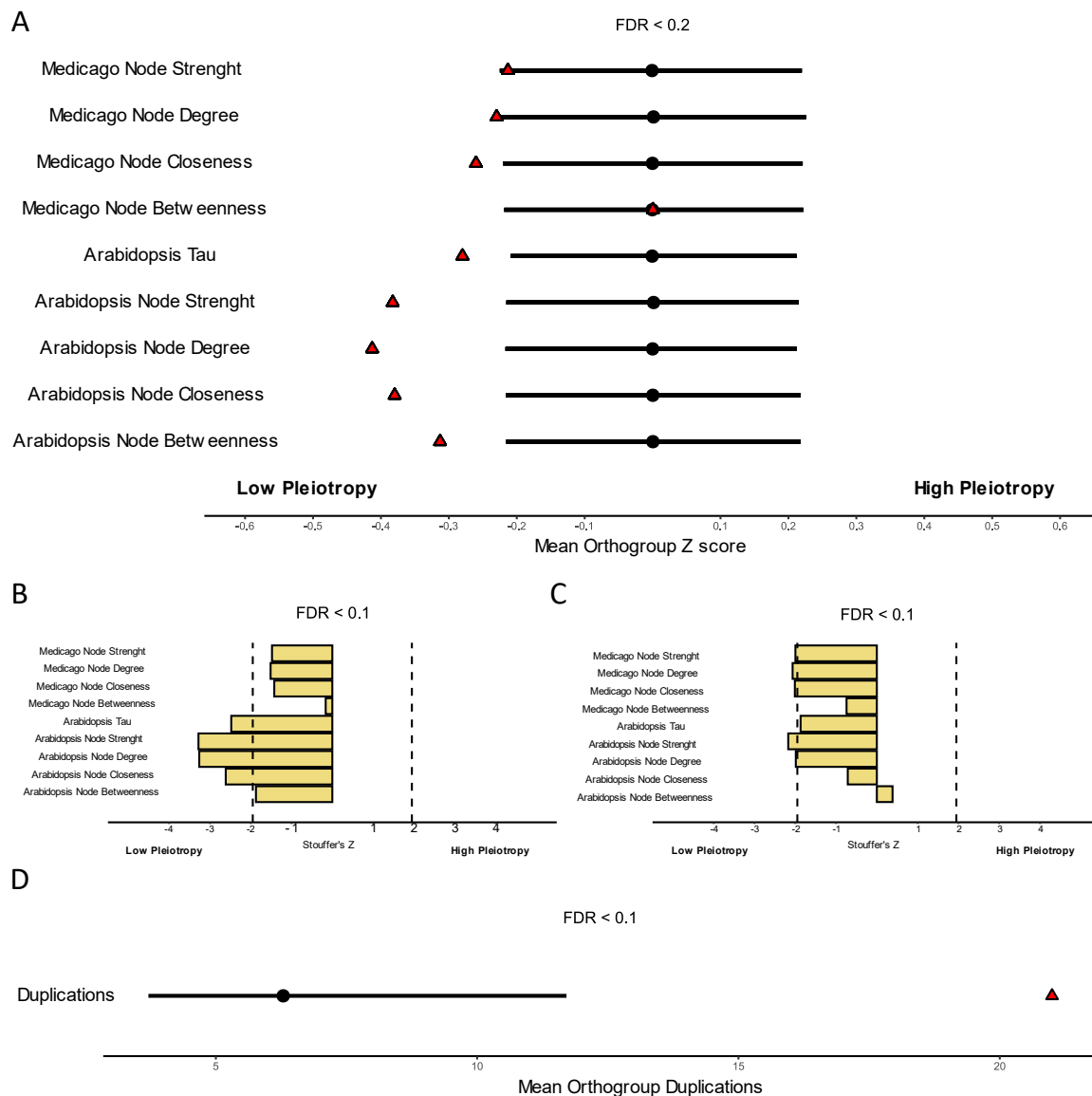


1028

1029 **Fig. 3.** Evidence for repeated selective sweeps across multiple species. (A) Distribution of *PicMin*  
 1030 FDR ( $-\log_{10}$  *PicMin* FDR on the Y axis) for the 13,268 tested orthogroups, ordered by number of  
 1031 putative driving genes on the X axis (number of *OmegaPlus* emp-p < 0.1). Points above the red line  
 1032 have FDR < 0.1. (B) Species contribution to *PicMin* top candidates (FDR < 0.1), calculated for each  
 1033 species as: [number of empirical p-values < 0.1 in the *PicMin* FDR < 0.1 orthogroups]/[total number of  
 1034 orthogroups with *PicMin* FDR < 0.1 tested]. (C) Heatmap of *OmegaPlus* empirical p-values for the 33  
 1035 candidate orthogroups (*PicMin* FDR < 0.1); driving genes cells (emp-p < 0.1 – black cells) are outlined  
 1036 in red, species are ordered by phylogenetic distance along the Y axis. A white cell indicates that an  
 1037 orthogroup was not tested in a species.

1038

1039

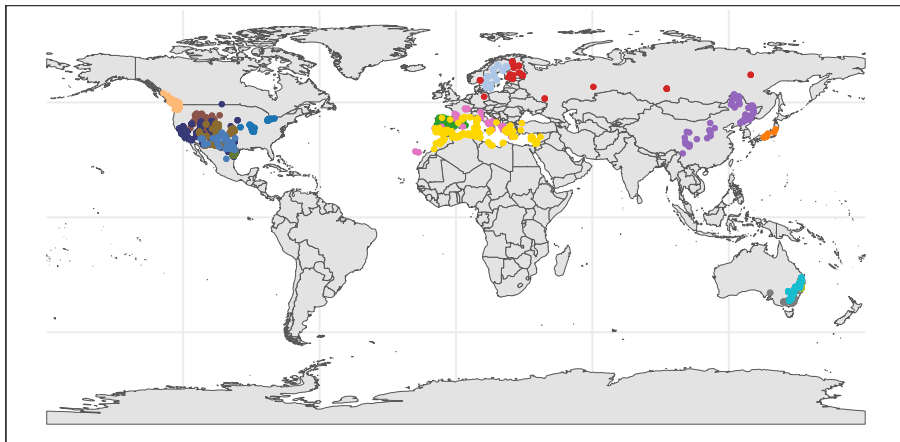


1040

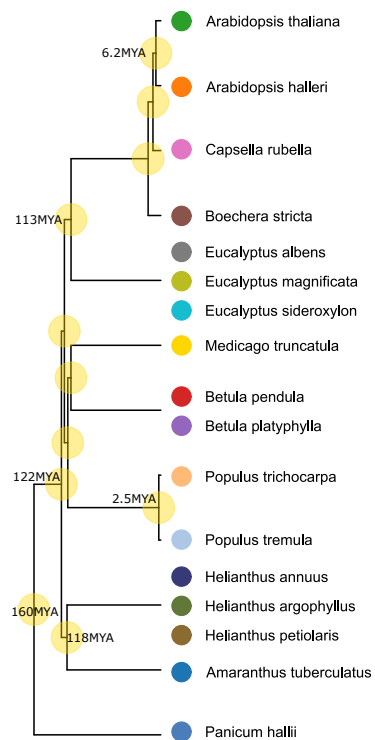
1041 **Fig. 4.** Orthogroups with repeated sweeps tend to be less pleiotropic and harbour more duplications.  
 1042 (A) Pleiotropy assessment of the relaxed set of 107 orthogroups with FDR < 0.2, against 10,000  
 1043 random draws. Red triangles represent the average pleiotropy of the candidate set, while black circles  
 1044 represent the mean of 10,000 random draws. Black lines represent the 95% interval of the random  
 1045 draws means. Each row represents a different pleiotropy measure, labelled on the left. (B) Pleiotropy  
 1046 Stouffer's Z of the top 33 orthogroups (FDR < 0.1), calculated using different pleiotropy measures  
 1047 labelled on the left. Dotted lines delimit the 95% confidence interval. (C) Pleiotropy Stouffer's Z for  
 1048 the more conservative set of 15 orthogroups identified by intersecting the 33 candidate OGs of the main  
 1049 analysis (FDR < 0.1) with results derived from nine additional *PicMin* omitting closely related  
 1050 *Eucalyptus* and *Helianthus* species. Dotted lines delimit the 95% confidence interval. (D) Duplication  
 1051 bootstrapping assessment of the top 33 orthogroups (FDR < 0.1) against 10,000 random draws. The  
 1052 red triangle represents the average duplication number in the candidate set of orthogroups, while the  
 1053 black circle and line represent mean and 95% interval of 10,000 random draws respectively.



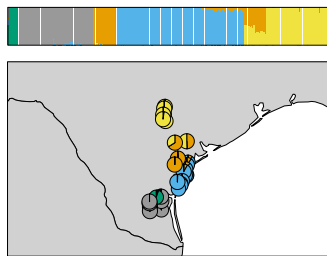
A



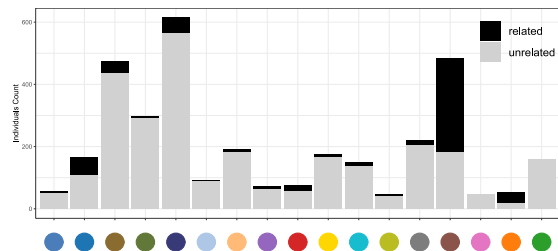
B

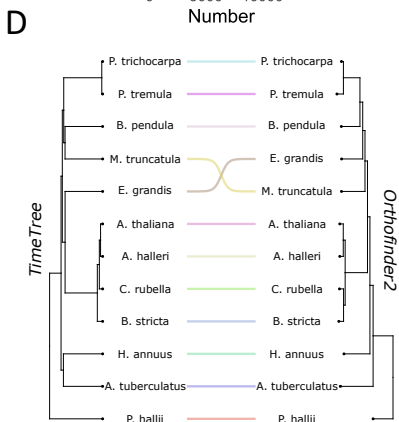
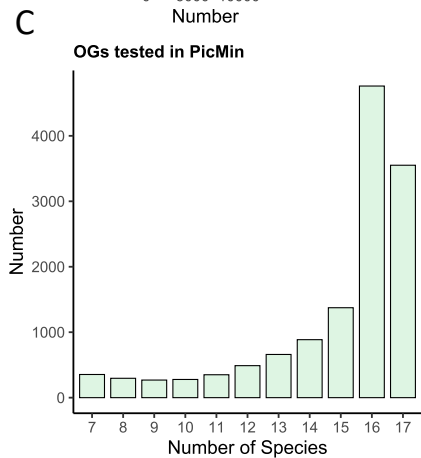
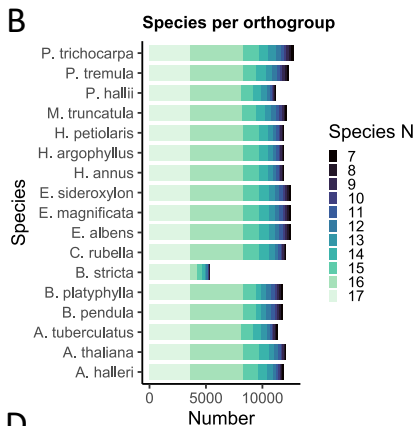
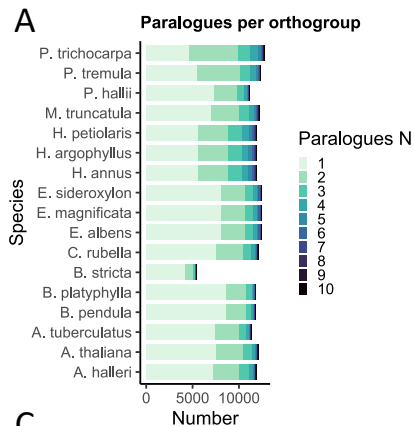


C



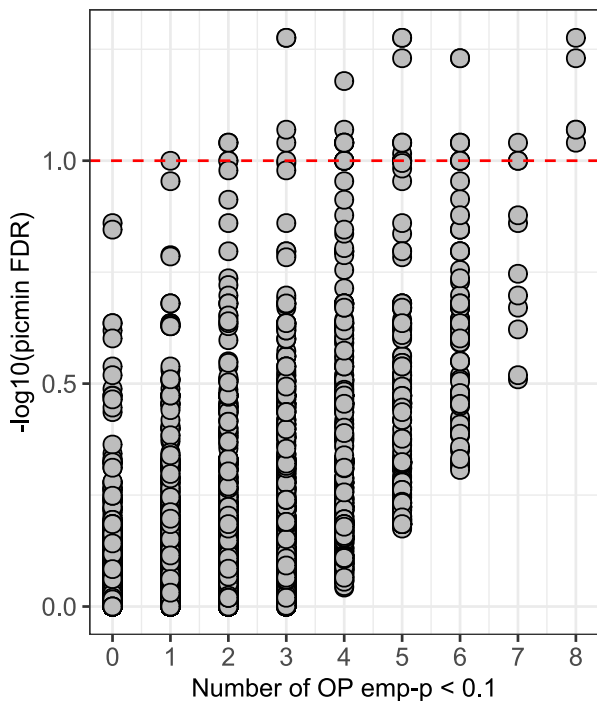
D



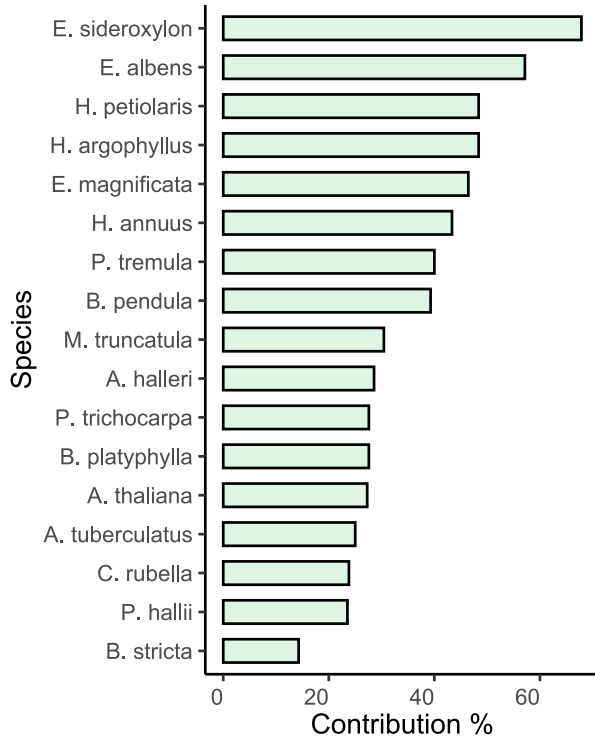


A

OGs tested = 13,268



B



C

