

## RESEARCH ARTICLE

# Adaptation and Phenotypic Diversification in Arabidopsis through Loss-of-function Mutations in Protein-coding Genes

Yong-Chao Xu<sup>1,2</sup>, Xiao-Min Niu<sup>1,2</sup>, Xin-Xin Li<sup>1,2</sup>, Wenrong He<sup>3</sup>, Jia-Fu Chen<sup>1,2</sup>, Yu-Pan Zou<sup>1,2</sup>, Qiong Wu<sup>1</sup>, Yong E. Zhang<sup>2,4</sup>, Wolfgang Busch<sup>3</sup>, Ya-Long Guo<sup>1,2\*</sup>

<sup>1</sup> State Key Laboratory of Systematic and Evolutionary Botany, Institute of Botany, Chinese Academy of Sciences, Beijing 100093, China

<sup>2</sup> University of Chinese Academy of Sciences, Beijing 100049, China

<sup>3</sup> Salk Institute for Biological Studies, Plant Molecular and Cellular Biology Laboratory, La Jolla, CA 92037, USA

<sup>4</sup> State Key Laboratory of Integrated Management of Pest Insects and Rodents & Key Laboratory of the Zoological Systematics and Evolution, Institute of Zoology, Chinese Academy of Sciences, Beijing, China

\* Corresponding Author: [yalong.guo@ibcas.ac.cn](mailto:yalong.guo@ibcas.ac.cn)

**Short title:** Adaptation via loss-of-function mutations

**One-sentence summary:** A large-scale study exploring how gene loss contributes to evolutionary change reveals that loss-of-function mutations contribute to the adaptation and phenotypic diversification of plants.

The author responsible for distribution of materials integral to the findings presented in this article in accordance with the policy described in the Instructions for Authors ([www.plantcell.org](http://www.plantcell.org)) is: [yalong.guo@ibcas.ac.cn](mailto:yalong.guo@ibcas.ac.cn).

## ABSTRACT

According to the “less-is-more” hypothesis, gene loss is an engine for evolutionary change. Loss-of-function (LoF) mutations resulting in the natural knockout of protein-

coding genes not only provide information about gene function, but they also play important roles in adaptation and phenotypic diversification. Although the “less-is-more” hypothesis was proposed two decades ago, it remains to be explored on a large scale. In this study, we identified 60,819 LoF variants in 1,071 *Arabidopsis thaliana* genomes and found that 34% of *A. thaliana* protein-coding genes annotated in the Col-0 genome do not have any LoF variants. We found that nucleotide diversity, transposable element density, and gene family size are strongly correlated with the presence of LoF variants. Intriguingly, 0.9% of LoF variants with minor allele frequency larger than 0.5% are associated with climate change. In addition, in the Yangtze River basin population, 1% of genes with LoF mutations were under positive selection, providing important insights into the contribution of LoF mutations to adaptation. In particular, our results demonstrate that LoF mutations shape diverse phenotypic traits. Overall, our results highlight the importance of the LoF variants for the adaptation and phenotypic diversification of plants.

## 1 INTRODUCTION

2 Variation in gene copy number plays important roles in adaptation and  
3 diversification. The two major processes that result in the gain of genes or gene  
4 copies are *de novo* gene formation and gene duplication (Chen et al., 2010;  
5 Carvunis et al., 2012; Guo, 2013; Palmieri et al., 2014; Zhao et al., 2014;  
6 McLysaght and Guerzoni, 2015; Li et al., 2016). However, the “less-is-more”  
7 hypothesis proposes that gene loss also contributes to evolutionary change in  
8 the context of gene copy number variation (Olson, 1999). The dosage balance  
9 hypothesis suggests that altering the stoichiometry of members of  
10 macromolecular complexes can affect the function of the whole complex and  
11 could ultimately affect phenotype and evolutionary fitness (Edger and Pires,  
12 2009; Birchler and Veitia, 2012; Hao et al., 2018). Therefore, the dosage  
13 balance hypothesis provides a more compelling explanation of the advantages  
14 and disadvantages of removing gene copies from a genome and explains the  
15 “less-is-more” hypothesis in a mechanistic way. There are two major  
16 mechanisms underlying gene loss: the physical removal of a gene by  
17 recombination, and gene inactivation by loss-of-function (LoF) mutations (Albalat  
18 and Canestro, 2016). The latter mechanism can occur through the gain of a  
19 premature stop codon, splice site disruption, or disruption of a transcript reading

20 frame. LoF mutations are frequent and have recently gained attention owing to  
21 the improved detection power of advanced sequencing techniques (MacArthur et  
22 al., 2012; Narasimhan et al., 2016).

23 Consistent with the “less-is-more” hypothesis, many cases of gene loss have  
24 been reported in diverse organisms, including unicellular organisms such as  
25 bacteria (Will et al., 2010; Hottes et al., 2013), and multicellular organisms, such  
26 as plants (Zufall and Rausher, 2004; Hoballah et al., 2007; Song et al., 2007;  
27 Tang et al., 2007; Gujas et al., 2012; Amrad et al., 2016; Sas et al., 2016; Wu et  
28 al., 2017) and animals (Greenberg et al., 2003; Hodgson et al., 2014; Goldman-  
29 Huertas et al., 2015), including humans (MacArthur et al., 2012; Narasimhan et  
30 al., 2016). LoF mutations have been shown to affect important biological  
31 processes, such as development and stress resistance in plants (Gujas et al.,  
32 2012; Wu et al., 2017), and intellectual disabilities in humans (Green et al.,  
33 2017). However, the extent to which LoF mutations are correlated with  
34 adaptation and phenotypic diversification is largely unknown (Albalat and  
35 Canestro, 2016). Even though loss of function mutations will produce large and  
36 detrimental effects on fitness will be depleted in inbreeding species compared to  
37 outcrossing plant and animal species, inbreeders are an excellent system to  
38 study the effect of natural homozygous knockouts, as inbreeding leads to a high  
39 frequency of homozygous LoF mutations (Saleheen et al., 2017). The best-  
40 characterized selfing species is *Arabidopsis thaliana*, which is naturally occurring  
41 in almost all parts of the world.

42 In this study, to investigate the evolutionary patterns of LoF mutations in  
43 natural populations, we explored 1,071 genomes of *A. thaliana* from the 1001  
44 Genomes Project (The 1001 Genomes Consortium, 2016), the African genomes  
45 project (Durvasula et al., 2017), and our own sequencing project (Zou et al.,  
46 2017). Across these 1,071 *Arabidopsis* genomes, we identified 60,819 LoF  
47 variants, including 17,453 stop codon-introducing (stop-gain) variants, 37,935  
48 disruptions of a transcript reading frame (frameshift), and 5,431 splice site-  
49 disrupting single nucleotide variants (splice site). We found that the presence of  
50 LoF mutations is correlated with the level of nucleotide diversity, the density of  
51 transposable elements (TEs), and gene family size. Intriguingly, 1% of genes  
52 with LoF mutations are under positive selection in the Yangtze River basin

53 population, suggesting that these genes with LoF mutations are crucial for  
54 adaptation. Overall, this study highlights the importance of LoF mutations for  
55 adaptation and phenotypic diversification.

56

## 57 **RESULTS**

### 58 **The Presence of LoF Variants is Correlated with Nucleotide Diversity, TE** 59 **density, and Gene Family Size**

60 We scanned the genomes of 1,071 accessions for LoF variants (Figure 1A and  
61 1B), including 893 genomes downloaded from the 1001 Genomes Project  
62 (Supplemental Data Set 1) (The 1001 Genomes Consortium, 2016), 61 from the  
63 African genomes project (Supplemental Data Set 2) (Durvasula et al., 2017), and  
64 117 from our own genome project (Supplemental Data Set 3) (Zou et al., 2017).  
65 Given that LoF mutations annotated by mapping sequences to a reference  
66 genome will have high false positive rates, a series of stringent filtering steps  
67 were implemented, which were similar to those employed in previous studies of  
68 human genomes (MacArthur et al., 2012; Narasimhan et al., 2016). First, single  
69 nucleotide variants (SNVs) and indels (insertions/deletions) were removed if they  
70 matched multiple times to the reference genome or were located in short tandem  
71 repeats or if the SNVs were in close proximity (3 bp or less) to an indel. Second,  
72 stop-gain variants were excluded if they were within a codon that was linked to  
73 other SNVs, and frameshift variants were removed if they were occurring  
74 together with other frameshifts that resulted in the restoration of the reading  
75 frame. Third, SNVs and indels were removed if the inferred LoF mutations  
76 occurred within the last 5% of the transcript or were observed in the ancestral  
77 states inferred from *Arabidopsis lyrata* and *Capsella rubella* (see Methods).  
78 Fourth, LoF variants that affect all known transcripts of the protein-coding gene  
79 in the reference genome (Col-0) were retained for subsequent analysis. After  
80 filtering according to these rules, we detected 17,453 (214.7 per accession)  
81 premature stop codons that were caused by SNVs (stop-gain), 37,935 (512.2  
82 per accession) insertion/deletion-induced frameshift variants (frameshift) leading  
83 to the disruption of a transcript reading frame, and 5,431 (103.0) splice site-  
84 disrupting SNVs (splice site) (Table 1). Taking all of these stop-gain, frameshift,

85 and splice site variants together, we identified 60,819 (829.9 per accession) LoF  
86 variants across the 1,071 accessions (Table 1).

87 To assess the power of our detection method and to estimate the false  
88 discovery rate (FDR) for LoF variants, we used RNA sequencing data for eight  
89 *de-novo* assembled genomes (Gan et al., 2011), the Landsberg erecta genome  
90 assembled with PacBio Sequel data (Zapata et al., 2016), and the KBS-Mac-74  
91 genome assembled with nanopore data (Michael et al., 2018). The results  
92 suggest that our method has a very low FDR (2.6% in each accession on  
93 average) (Supplemental Data Set 4). Given that the accumulated FDR across  
94 1,071 accessions might be higher than that in a single genome, we calculated  
95 the accumulated false positive rates based on validation data in these 10  
96 accessions. While with the addition of genomes to the analysis, the accumulated  
97 false positive LoF as well as the total number of LoF increased, the accumulated  
98 FDR was largely stable (Supplemental Figure 1). Therefore, we estimated that in  
99 each *A. thaliana* accession, an average of 808.2 genes have LoF mutations  
100 (3.0% of all protein-coding genes in the reference Col-0) with a FDR of 2.6%. In  
101 particular, 80 false positive LoF mutations (22 stop-gains, 56 frameshifts, and 2  
102 splice sites) identified in the validation process were excluded from the  
103 subsequent analysis.

104 The density of genes with stop-gain, frameshift, splice site, and LoF variants  
105 (summed up for all the stop-gain, frameshift, and splice site variants), as  
106 calculated by dividing the number of genes with LoF variants (genes with LoF  
107 variants detected in any of the 1,071 accessions) by the gene number in each  
108 200 kb window in the Col-0 reference genome, varied across the genome  
109 (Figure 1C, Supplemental Figure 2). By contrast, the density of stop-gain and  
110 frameshift variants within gene bodies was roughly evenly distributed, with a  
111 slightly higher frequency at the beginning and at the end of genes (Supplemental  
112 Figure 3). This pattern is similar to the distribution in human and great ape  
113 genomes (MacArthur et al., 2012; de Valles-Ibanez et al., 2016).

114 To infer mechanisms that favor LoF, we analyzed the correlations between  
115 LoF and various gene features. The density of genes with LoF mutations (stop-  
116 gain, frameshift, splice site and LoF variants) across the Col-0 genome showed  
117 a strong correlation with nucleotide diversity ( $\pi$ ) across the 1,071 genomes and

118 the TE density in the Col-0 genome (Figure 1D). Large gene families tend to  
119 have high gene redundancy (Wagner, 2005), and gene loss in large gene  
120 families may be an advantage according to the dosage balance hypothesis  
121 (Edger and Pires, 2009; Birchler and Veitia, 2012; Hao et al., 2018). We  
122 therefore evaluated whether these genes were prone to acquiring LoF  
123 mutations. We classified the 27,206 protein-coding genes of the reference Col-0  
124 genome into 7,430 gene families based on amino acid sequence similarity. We  
125 then estimated the correlation between gene family size and the proportion of  
126 genes with common LoF variants (variants displaying a minor allele frequency  
127 [MAF] larger than 5% across the 1,071 accessions). For all gene families, there  
128 was a positive correlation between gene family size and the LoF ratio (the  
129 percentage of genes with LoF variants in each gene family) (Spearman  $r = 0.21$ ,  
130  $p = 2.2 \times 10^{-16}$ ). This suggests that genes belonging to larger gene families are  
131 more likely to have LoF variants. Similarly, the median sizes of gene families for  
132 genes with and without LoF variants were 11 and 6, respectively (Figure 2A,  
133 Wilcoxon sum test,  $p < 2.2 \times 10^{-16}$ ). This indicates that genes belonging to larger  
134 gene families are more likely to acquire LoF mutations. Furthermore, analogous  
135 to previous studies (Prachumwat and Li, 2008; Guo, 2013), we divided gene  
136 families into three classes: one gene or singleton (class 1: 3,825 families), two  
137 genes (class 2: 1,353 families), and three genes or more (class 3: 2,252  
138 families). The genes in class 3 were more likely to contain LoF variants than the  
139 others (Supplemental Figure 4, class 3 vs. class 1 or 2, Wilcoxon sum test, all  $p$   
140  $< 0.01$ ). Therefore, the presence of LoF variants is correlated with nucleotide  
141 diversity, TE density, and gene family size.

## 142 **A Large Fraction of *A. thaliana* Genes could Affect Fitness in Natural** 143 **Environments**

144 A genome-wide study of LoF variants can provide insight into gene lethality by  
145 analyzing the frequency of all null alleles for a given gene in natural  
146 environments (population gene lethality) (Albalat and Canestro, 2016). We found  
147 that 9,249 protein-coding genes (34.0% of all protein-coding genes in the Col-0  
148 reference genome) do not have any LoF variant within our panel of 1,071  
149 accessions. To ask whether known essential genes are depleted of LoF variants,  
150 we used a database of 358 lethal genes identified in laboratory studies (Meinke

151 et al., 2008). 88.3% of the lethal genes do not have any natural LoF mutation,  
152 and 7.3% of the lethal genes display a LoF allele frequency less than 0.5%  
153 across the 1,071 accessions (Figure 2B). This suggests that lethal genes in *A.*  
154 *thaliana* are significantly depleted of natural LoF mutations (Pearson's chi-  
155 squared test,  $p < 2.2 \times 10^{-16}$ ). We then performed gene ontology (GO)  
156 enrichment analysis of the four different gene classes based on the LoF allele  
157 frequency in the 1,071 accessions (genes without LoF, with LoF, with a LoF  
158 allele frequency of  $<0.05$  and  $\geq 0.05$ ). Although both the gene sets (with and  
159 without LoF variants) were enriched for the GO term 'unknown function', the only  
160 gene set highly enriched for many other GO terms was the set without LoF  
161 variants (Figure 2C). These results suggest a model in which natural knockouts  
162 of 9,249 *A. thaliana* genes without LoF variants would reduce fitness in natural  
163 environments, and individuals with LoF mutations in these non-essential genes  
164 can likely be removed from populations by purifying selection.

#### 165 **LoF Variants are Correlated with Climate Variables**

166 LoF variants could be functionally important and associated with climate  
167 variables. To investigate this notion, we performed a genome-wide association  
168 study (GWAS) of 14,270 LoF variants with MAF larger than 0.5% and 20  
169 environmental variables, including latitude and 19 climate variables  
170 (Supplemental Data Set 5). To validate that the associations were not random,  
171 we compared the total number of significantly associated LoF variants with the  
172 number obtained from 1,000 GWAS permutations based on 14,270 randomly  
173 resampled SNPs (single-nucleotide polymorphisms; MAF  $> 0.5\%$ ) in the coding  
174 regions of protein-coding genes across the whole genome for each of the 20  
175 environmental variables. If LoF variants were randomly associated with  
176 environmental variables, then we would expect the number of associated  
177 variants obtained in the permutation analysis and the observed value for each  
178 environmental variable to be similar. For 10 of the 20 environmental variables,  
179 the observed number was significantly higher than the number identified by  
180 permutation analysis (Supplemental Figure 5, one-sample Wilcoxon signed rank  
181 test,  $p < 0.05$ ), suggesting that LoF variants are more prone to associate with  
182 climate variables. 125 of the 14,270 LoF variants (0.9%) were significantly  
183 associated with at least one of the 10 environmental variables, involving 124

184 genes (Supplemental Figures 6-8 and Supplemental Data Set 6). Of these  
185 genes, some are known to be functionally important, such as *CRK36*  
186 (Supplemental Data Set 6). *CRK36* encodes an abiotic stress-inducible receptor-  
187 like protein kinases that negatively regulates abscisic acid signaling (Tanaka et  
188 al., 2012).

189 We used RNA-seq data for 541 accessions to compare the expression levels  
190 of alleles with and without LoF mutations for each of these environmentally  
191 associated genes (Kawakatsu et al., 2016) (Figure 2D). Among the 124  
192 associated genes, 49 genes have both LoF and non-LoF alleles in the 541  
193 accessions for which RNA-seq data were available. There were no LoF variants  
194 of the other 75 genes in this set of accessions. We therefore focused on the 49  
195 genes with both LoF and non-LoF alleles (different LoF variants of the same  
196 gene were combined). The transcript levels of 16.4% of these genes were  
197 significantly altered (8.2% upregulated and 8.2% downregulated) in accessions  
198 harboring the LoF variants (Wilcoxon sum test, FDR corrected  $p < 0.05$ ) (Figure  
199 2D). Furthermore, genes with LoF variants associated with environmental  
200 variables exhibited greater rates of upregulation (8.2% vs. 3.0%) or  
201 downregulation (8.2% vs. 6.5%) than genes with LoF variants at the genome  
202 level (Figure 2E). To determine whether the expression changes were caused by  
203 the different genetic backgrounds of various natural accessions, we calculated  
204 the up- and downregulated gene numbers based on 30,000 randomly selected  
205 alleles with non-LoF variants in coding regions across the 541 accessions. A  
206 higher proportion of genes with LoF variants associated with environmental  
207 variables are up- (8.2% vs. 1.7%) or downregulated (8.2% vs. 2.4%) compared  
208 to those detected when only the genetic background was taken into account  
209 (Supplemental Figure 9). The higher proportion of upregulated genes with LoF  
210 variants associated with environmental variables could result from diverse  
211 mechanisms, such as a dominant gain-of-function effect, in which LoF mutations  
212 are effectively gain-of-function mutations, thereby mimicking increased gene  
213 function (Schild et al., 1995; Xie et al., 1998).

#### 214 **LoF Variants of *KUK*, *PRR5*, and *LAZY1* Shape Phenotypic Variation**

215 Common LoF alleles are nearly always less deleterious than rare LoF alleles,  
216 but some could have an impact on phenotypes (MacArthur et al., 2012).

217 Considering that stop-gain mutations can exert severe effects by directly  
218 truncating genes, we performed an additional GWAS analysis for stop-gain  
219 mutations and the 10 environmental variables that are significantly associated  
220 with LoF variants (Supplemental Figure 10). We found that many genes with  
221 stop-gain mutations are associated with environmental variables (Supplemental  
222 Data Set 7). One gene with a significant GWAS signal for a precipitation-related  
223 trait (precipitation in the wettest month) is the F-box protein gene *KUK* (Figure  
224 2F, Supplemental Figures 10 and 11), which is involved in regulating root  
225 development in Arabidopsis (Meijon et al., 2014). *KUK* alleles with and without  
226 LoF variants are widely distributed across Eurasia, and most accessions from  
227 the Yangtze River basin carry stop-gain variants (Supplemental Figure 12A).  
228 Based on the phenotypic data for 129 accessions available from a previous  
229 study (Figure 3A and Supplemental Data Set 8) (Meijon et al., 2014), the  
230 meristematic zone and mature cortical cells of accessions with LoF variants  
231 (stop-gain, frameshift, or total LoF; no putative splice variation was observed)  
232 were significantly longer than those of accessions without LoF variants (Figure  
233 3B and Supplemental File 1; population structure-corrected ANOVA,  $p < 0.05$ ).

234 To independently test the hypothesis that truncated versions of *KUK* are  
235 associated with increased meristem and mature cell sizes, we phenotyped  
236 accessions containing specific LoF alleles that were not evaluated in the  
237 previous study (with one exception: Bur-0) (Meijon et al., 2014). We chose three  
238 accessions that did not contain a stop codon (control), three accessions with the  
239 most common stop codon located in the first half of *KUK* (111 G/A: TGG-TGA,  
240 SG111), the only two accessions in the 1001 Genomes panel that had a distal  
241 stop codon (652 C/T: CGA-TGA, SG652), including one of the accessions with  
242 the largest meristem and mature cell traits (Bur-0) according to a previous study  
243 (Meijon et al., 2014), and two accessions with a distal frameshift (738 T/-,  
244 FS738) (Supplemental Figure 13 and Supplemental Data Set 9). Compared with  
245 the control, all three classes showed significant increases (pairwise Wilcoxon  
246 rank sum test,  $p < 0.05$ ) in meristem length, with the most common stop codon  
247 associated with the smallest increase (Supplemental Figure 13 and  
248 Supplemental Data Set 9). Overall, these data suggest that truncated alleles of  
249 *KUK* influence cellular traits, most likely via a dominant mechanism. Further

250 transgenic studies should be performed to confirm the effects of these genetic  
251 variants in the same genetic background.

252 In addition to the LoF alleles of genes that were associated with environmental  
253 variables based on the GWAS, we studied *PSEUDO-RESPONSE*  
254 *REGULATOR5* (*PRR5*), which is involved in regulating various circadian-  
255 associated biological events, such as flowering time (Nakamichi et al., 2016).  
256 *PRR5* alleles with and without LoF variants are widely distributed in our set of  
257 1,071 accessions (Supplemental Figure 12B). Of the 812 accessions with  
258 flowering time data (The 1001 Genomes Consortium, 2016) (Figure 3C and  
259 Supplemental Data Set 10), the flowering time of accessions with *PRR5* LoF  
260 variants (frameshift, or total LoF; no putative splice variation was observed) was  
261 significantly delayed compared to accessions without LoF variants at either 10°C  
262 or 16°C (Figure 3D and Supplemental File 1; population structure-corrected  
263 ANOVA,  $p < 0.05$ ). We note that while flowering time for the stop-gain allele was  
264 delayed relative to the non-LoF allele, the difference was not significant,  
265 probably owing to the small number of accessions.

266 Given the above two cases of genes having common LoF variants, we further  
267 studied *LAZY1* as an example for a rare LoF variant. Only one accession from  
268 the Yangtze River basin population (29-8) showed a frameshift mutation in  
269 *LAZY1*, which that was caused by a 20 bp deletion in exon 3 and resulted in a  
270 premature stop codon (Figure 4D and Supplemental Figure 14). *LAZY1* mutant  
271 plants display a much larger branch angle (81°) than wild type (42°) (Yoshihara  
272 et al., 2013). As expected, we found that accession 29-8 has a large branch  
273 angle (74°) (Figure 4A). To validate that this frameshift mutation was the causal  
274 mutation, we crossed accession 29-8 with accessions displaying smaller branch  
275 angles: Col-0 (30°) and 3-2 (30°). In the Col-0 × 29-8 F<sub>2</sub> population, the  
276 segregation ratio for branch angles was roughly 3:1, suggesting that a major,  
277 recessively acting gene was responsible for the large branch angle of 29-8  
278 (Figure 4B). Using a mapping-by-sequencing approach based on 57 F<sub>2</sub> plants  
279 from the Col-0 × 29-8 F<sub>2</sub> population with large branch angles (>85°), we  
280 identified a causal region on chromosome 5 (Figure 4C). To refine the target  
281 interval, we used F<sub>2</sub> plants in both Col-0 × 29-8 and 3-2 × 29-8 populations, and  
282 finally narrowed down the causal locus to a 240 kb interval containing *LAZY1*

283 (Figure 4D, Supplemental Data Set 11 and 12). The truncated transcript of  
284 *LAZY1* in the 29-8 accession caused the loss of the only functional domain of  
285 the nuclear localization signal (NLS), which is located between conserved  
286 regions 3 and 4. This NLS is a key component during the nuclear import process  
287 (Supplemental Figure 14) (Yoshihara et al., 2013). Overall, these results provide  
288 evidence that the frameshift variant of *LAZY1* is likely responsible for the larger  
289 branch angle in accession 29-8.

### 290 **A subset of LoF variants exhibits a signature of natural selection**

291 Compared with other SNP variants across the whole genome in the 1,071  
292 accessions, LoF variants are biased towards low frequencies (Figure 5A). This  
293 suggests that many natural LoF variants might be deleterious in *Arabidopsis* and  
294 are under purifying selection. Similarly, compared with 1,000 permutation results  
295 of non-LoF variants associated with environmental variables (GWAS non-LoF  
296 variants), LoF variants associated with environmental variables (GWAS LoF  
297 variants) were biased towards low allele frequencies (Figure 5A). This suggests  
298 that a large fraction of LoF variants that are associated with environmental  
299 variables are likely under purifying selection as well.

300 Although LoF mutations are usually deleterious, adaptive LoF mutations can  
301 occur and spread rapidly in small populations (Olson, 1999). To determine  
302 whether genes with LoF mutations were influenced by natural selection, we  
303 focused on the Yangtze River basin population (86 accessions, Supplemental  
304 Data Set 3), a population that recently adapted to this region (Zou et al., 2017).  
305 In total, 2,709 genes had LoF variants (3,349 LoF variants) in this population;  
306 among them, 54 genes (64 LoF variants) were under positive selection based on  
307 a selection sweep analysis in our recent study (Zou et al., 2017) (Figure 5B, see  
308 Supplemental Data Set 13 and 14 for details). Of the 54 genes in these selected  
309 regions, 27 genes (1.0% of the genes with LoF variants) had a LoF allele  
310 frequency exceeding 95% and were more likely under positive selection,  
311 including *KUK*, in which one LoF allele (111 G/A: TGG-TGA type, Supplemental  
312 Figure 11) was fixed in the Yangtze River basin population (Supplemental Data  
313 Set 14).

314 Finally, we performed ecological niche modeling of non-LoF (571) and LoF  
315 (500) alleles of *KUK*. In niche identity tests, the simulated values were  
316 significantly higher than the observed value ( $Io = 0.877$ ,  $p < 0.001$ ), indicating  
317 that there is significant ecological differentiation between the two groups (Figure  
318 5C, Supplemental Figure 15). This result suggests that the LoF allele of *KUK* is  
319 associated with adaptation to the Yangtze River basin. Overall, our results  
320 indicate that LoF mutations can be associated with adaptation and phenotypic  
321 diversification; more importantly, at least 1.0% of the genes with LoF variants are  
322 under positive selection in the Yangtze River basin.

323

## 324 **DISCUSSION**

325 The “less-is-more” hypothesis proposes an evolutionary process involving  
326 adaptive gene loss (Olson, 1999). LoF mutations have gained increasing  
327 attention (MacArthur et al., 2012; Yang et al., 2015; Narasimhan et al., 2016).  
328 For example, recent studies have demonstrated the functional importance of  
329 natural LoF mutations (Gujas et al., 2012; Lek et al., 2016; Saleheen et al.,  
330 2017; Wu et al., 2017). However, our understanding of the evolutionary pattern  
331 of LoF mutations at the genome level and the adaptive effects of LoF variants at  
332 the population level is highly limited. In this study, we investigated the  
333 evolutionary pattern of LoF mutations in *A. thaliana* and found that a high level of  
334 nucleotide diversity, high TE density, and high gene redundancy (large gene  
335 family size) are associated with LoF mutations. The latter two factors were also  
336 observed in a previous study of LoF mutations in human populations (MacArthur  
337 et al., 2012). Given that these three factors themselves are largely correlated  
338 with rates of evolution, rapid sequence evolution could be associated with a high  
339 frequency of LoF mutations. More importantly, we hypothesize that 34% of *A.*  
340 *thaliana* genes can affect fitness in natural environments, and furthermore, 1% of  
341 genes with LoF variants are under positive selection in the Yangtze River basin  
342 population.

343 The mechanisms by which LoF mutations can produce functional effects are  
344 complicated. Some of the profound effects of LoF can be explained by the  
345 dosage balance hypothesis (Birchler and Veitia, 2007; Hou et al., 2018; Kremling

346 et al., 2018). A simple mechanism for a beneficial effect of a null mutation is the  
347 removal of a protein that is detrimental in the current environment (Hottes et al.,  
348 2013). For example, *brx* (*brevis radix*) LoF alleles in Arabidopsis help roots  
349 adapt to acidic soil (Gujas et al., 2012). However, it is also possible that LoF  
350 variants can act as dominant mutations. For instance, stop-gain variants of the  
351 SOX9 transcription factor MiniSOX9 act in a dominant-negative manner, thereby  
352 counteracting the activity of the wild-type variant (Abdel-Samad et al., 2011).  
353 Other stop-gain and missense variants can act as dominant gain-of-function  
354 mutations, thereby mimicking increased gene function (Schild et al., 1995; Xie et  
355 al., 1998). In the case of *KUK*, a gene that acts as a positive regulator of  
356 meristem and cell size (Meijon et al., 2014), a hypermorphic gain-of-function is a  
357 more likely explanation than a LoF. Although most premature stop codons result  
358 in a LoF or dominant effects, more complex possibilities exist. For instance, in  
359 the fruit fly *Drosophila sechellia*, an olfactory receptor pseudogene encodes a  
360 functional receptor as a result of translation read-through of the premature  
361 termination codon (Prieto-Godino et al., 2016).

362 Of the many LoF variants (or pseudogenes) that we identified in over 1,000  
363 natural accessions, some are presumably functional, with a minor or altered  
364 function, rather than non-functional. More efforts are needed to understand their  
365 functional effects. For example, transgenic analyses are needed to clarify the  
366 functional differences between the predicted pseudogenes and their ancestral  
367 genes. While more function research is needed, LoF mutational variation also  
368 provides a way of investigating the phenotypic consequences of the loss of  
369 specific genes within diverse genomes while at the same time establishing their  
370 role in the evolutionary process. Overall, our study highlights the importance of  
371 natural knockouts for adaptation and phenotypic diversification and emphasizes  
372 the need for more in-depth studies of LoF variants.

373

## 374 **METHODS**

### 375 **Data analyzed in this study**

376 The genomes of 1,071 accessions were used, including 893 representative  
377 genomes of *A. thaliana* from the 1001 Genomes Project (The 1001 Genomes

378 Consortium, 2016) (Supplemental Data Set 1), 61 from the African sequenced  
379 genomes project (Durvasula et al., 2017) (Supplemental Data Set 2), and 117  
380 genomes from our own sequencing project (Zou et al., 2017) (Supplemental  
381 Data Set 3). The 541 transcriptomes used in this study were obtained from a  
382 previous study, as indicated in Supplemental Data Set 1 (Kawakatsu et al.,  
383 2016). All 1,071 accessions with clean data were mapped against the Col-0  
384 reference genome, and SNPs and indels were called using Genome Analysis  
385 Toolkit (GATK version 2.1.8) with default parameters (DePristo et al., 2011).

### 386 **Identification of candidate LoF variants**

387 High-quality homozygous SNPs and indels (quality (Q)  $\geq$  30, quality-by-depth  
388 ratio (QD)  $\geq$  10 (QD  $\geq$  5 for indels), ReadPosRankSum  $\geq$  -8.0, depth of coverage  
389 (DP)  $\geq$  3, probability of strand bias (FS)  $\leq$  10.0 (FS  $\leq$  200.0 for indels)) were used  
390 to identify LoF mutations using SnpEff software (SnpEff 3.3f) (McLaren et al.,  
391 2010). SNPs annotated as STOP GAINED, SPLICE SITE DONOR, and SPLICE  
392 SITE ACCEPTOR (SPLICE SITE DONOR and SPLICE SITE ACCEPTOR were  
393 combined as splice site), and indels annotated as FRAMESHIFT, were regarded  
394 as LoF variants and included in subsequent analyses.

### 395 **Filtering of candidate LoF variants**

396 To remove false positive LoF variants caused by sequencing and mapping  
397 errors, annotation errors, and the effect of nearby variants, a series of stringent  
398 filters were used. First, a sequence context filter was applied. Variants that could  
399 be mapped to multiple regions of the reference genome (Col-0) were removed  
400 using Genome Multitool (Derrien et al., 2012; Marco-Sola et al., 2012). SNVs  
401 (stop-gain and splice site) and indels (frameshift) overlapping with tandem  
402 repeats and SNVs in close proximity (3 bp or less) to a known indel were  
403 excluded. Second, multi-nucleotide polymorphism filters were applied. To filter  
404 incorrect premature stop codon annotation, stop-gain variants that were found  
405 within the same codon linked to other SNVs were removed. Frameshift variants  
406 were filtered out if two frameshift variants resulted in the restoration of the  
407 reading frame. Three or more frameshift variants occurring in each gene of the  
408 same accession were excluded; most resulted in the restoration of the reading  
409 frame. Third, annotation filters were applied. SNVs and indels were filtered out if

410 the inferred LoF allele was also the ancestral state, as this implies that such  
411 variants were “gain-of-function” or the sequencing errors of the gene model at  
412 this location in the reference.

413 Two related species, *Arabidopsis lyrata* (MN47) and *Capsella rubella* (MTE),  
414 were used to infer the ancestral state of SNVs and indels. For SNVs and indels,  
415 ancestral states were determined using alignments of Col-0 with the *A. lyrata*  
416 and *C. rubella* reference genomes, as described in our previous study (Li et al.,  
417 2016). SNVs and indels were filtered out if either of the two outgroups showed  
418 the same stop-gain or frameshift mutations observed in the 1,071 accessions.

419 The disrupted fraction of the coding sequence of the longest transcript caused  
420 by stop-gain and frameshift variants was calculated, and stop-gain or frameshift  
421 mutations occurring within the last 5% of the transcript were removed. Because  
422 we did not determine the position of splice site disruption on the final transcript  
423 (MacArthur et al., 2012), we were not able to perform this analysis for splice site  
424 variation. Finally, if there were two or more LoF variants in the same allele, only  
425 the LoF variant closest to the start codon (ATG) was considered.

#### 426 **Validation of LoF mutations**

427 We estimated the FDR using the assembled transcripts of eight accessions from  
428 a previous study (Gan et al., 2011), as well as gene sequences that were  
429 annotated using two published long reads assembled genomes: *Landsberg*  
430 *erecta* assembled from PacBio Sequel data (Zapata et al., 2016) and KBS-Mac-  
431 74 assembled from nanopore data (Michael et al., 2018). The genes of KBS-  
432 Mac-74 and *Landsberg erecta* were annotated using exonerate v2.2.0 (Slater  
433 and Birney, 2005) with CDS (coding sequences) of Col-0. All assembled  
434 transcript sequences for genes with LoF mutations in the eight accessions, and  
435 annotated gene sequences of KBS-Mac-74, *Landsberg erecta*, and Col-0 gene  
436 sequences were aligned using MUSCLE v3.8.31 (Edgar, 2004). All LoF  
437 mutations were manually checked based on the aligned sequences. Stop-gain  
438 and frameshift variants of the 10 accessions were used to assess the curve of  
439 accumulated false positive rate (all false positive LoF variants found in each  
440 accession combined). Splice site variants were not used in the analysis of  
441 accumulated false positive rate, because assembled transcripts of eight

442 accessions were assembled based on short reads, which are not robust enough  
443 to validate splice sites. Nevertheless, the FDRs for splice sites in the two long  
444 reads assembled genomes are similar or even lower than either stop-gain or  
445 frameshift. The accumulated false positive LoF number, accumulated total LoF  
446 number (the number of LoF variants with transcriptional or long reads data in  
447 each accession), and accumulated false positive LoF ratio were calculated with  
448 all probable accumulated groups.

#### 449 **LoF polymorphism variation**

450 All parameters were calculated using a 200 kb window size and 10 kb step size  
451 along the Col-0 reference genome. The density of genes with stop-gain variants  
452 was calculated by dividing the number of genes with stop-gains by the number of  
453 genes in each 200 kb window in the Col-0 reference genome. The same  
454 rationale was used to calculate the density of frameshift, splice site, and total  
455 LoF variants (summed over all the three LoF variants). Nucleotide diversity ( $\pi$ )  
456 was calculated based on the SNP matrix of 1,071 accessions in each 200 kb  
457 window. The GC content was calculated based on the Col-0 genome in each  
458 200 kb window. Exon number, CDS length, and gene expression level were  
459 calculated for each gene based on the Col-0 genome in each 200 kb window.  
460 Gene expression levels in Col-0 were obtained from a previous study  
461 (Kawakatsu et al., 2016). TE density was calculated as the number of TEs in  
462 each 200 kb window in the Col-0 reference genome with TE annotations  
463 downloaded from TAIR10.

#### 464 **Gene family, sequence, and expression analysis**

465 Gene family analysis was performed according to our previous study (Guo,  
466 2013). GO term enrichment analysis was performed using agriGO (Tian et al.,  
467 2017). All genes (5,050) with candidate LoF variants were removed from the  
468 non-LoF category in the GO analysis. PCR was performed to obtain the *LAZY1*  
469 upstream region (~1.5 kb), gene body, and downstream regions (~0.5 kb) from  
470 Col-0, 3-2, and 29-8 (marked with an asterisk in Supplemental Data Set 3) using  
471 Q5 polymerase (New England Biolabs, Ipswich, MA, USA). Sequences were  
472 aligned using Lasergene Seqman (DNASTAR, Madison, WI, USA). Markers  
473 used in the analysis of plants with large branch angles are listed in Supplemental

474 Data Set 11 and 12. Primers used for the sequencing and amplification are listed  
475 in Supplemental Data Set 15.

476 A total of 541 transcriptomes were available from the 1001 Genomes Project  
477 (Supplemental Data Set 1) (Kawakatsu et al., 2016), and 6,675 genes with LoF  
478 variants whose mean expression levels (FPKM) were larger than 3 across 541  
479 accessions were kept for subsequent analysis (Meng et al., 2016). For each of  
480 the 6,675 genes with LoF variants, 541 alleles were divided into two groups:  
481 alleles with and without LoF variants. Then, the Wilcoxon sum test was used to  
482 evaluate the expression data. All  $p$  values were adjusted for multiple testing by  
483 computing their FDR (Benjamini and Hochberg, 1995). Up- or downregulation  
484 was calculated using the mean expression data in the two groups (alleles with  
485 and without LoF variants).

#### 486 **Genome-wide association study**

487 GWAS was performed by using the compressed mixed linear model (Zhang et  
488 al., 2010) from the GAPIT R package (Lipka et al., 2012) to associate data for 20  
489 environmental variables (latitude and 19 climate variables from the  
490 WORLDCLIM database) (Supplemental Data Set 5) with the 14,270 LoF variants  
491 (MAF > 0.5%) of the 1,071 accessions. Principal components (“Q” matrix) and a  
492 kinship matrix (“K” matrix) were included to account for population structure.  
493 Bayesian information criterion-based model selection as implemented in the  
494 GAPIT R package was used to find the optimal number of principal components  
495 for each trait. The whole-genome significance cutoff for associations was set to  
496  $0.01/\text{total LoF variants}$  ( $-\log_{10}(p) = 6.12$ ). For the GWAS permutation analysis,  
497 the SNP matrix that contained 14,270 SNPs (MAF > 0.5%) was randomly  
498 resampled 1,000 times for those SNPs in the coding region of protein-coding  
499 genes across the whole genome. Whether the actual value is larger than the  
500 permutation results was tested using the one-sample Wilcoxon signed rank test.  
501 GWAS analysis was also performed using 3,822 stop-gain variants (MAF > 0.5%)  
502 with these 20 environmental variables within the 1,071 accessions. The  
503 significance cutoff in this analysis was set to  $0.05/\text{total LoF variants}$  ( $-\log_{10}(p) =$   
504  $4.88$ ).

#### 505 **Ecological niche modeling analysis**

506 MaxEnt 3.3.3 was used to perform ecological niche modeling via maximum  
507 entropy with default settings (Phillips et al., 2006) to examine ecological  
508 divergence between the two *KUK* groups (excluding redundant and unreliable  
509 records) combined with data for 19 ecological variables downloaded from the  
510 WORLDCLIM database (Phillips and Dudik, 2008). For the analysis, default  
511 settings of MaxEnt 3.3.3 were used. Ten of the environmental variables with  
512 pairwise Pearson correlation coefficients of  $-0.7 < r < 0.7$  were selected for final  
513 analysis (marked red in Supplemental Data Set 5) with outputs set at 25% for  
514 testing and 75% for training the model. The model accuracy was assessed  
515 based on the area under the curve (AUC), with scores between 0.7 and 0.9  
516 indicating good discrimination, and  $AUC > 0.9$  indicating reliable discrimination  
517 (Swets, 1988). The range of suitable distributions was drawn using DIVA-GIS  
518 v7.5 (Hijmans et al., 2005).

519 ENMTools 1.3 was used to perform niche identity test by calculating Warren's  
520 *I*, which ranged from 0 (no niche overlap) to 1 (identical niches), with 100  
521 pseudo-replicates, between the LoF and non-LoF *KUK* groups using the 10  
522 environmental variables that were used in MaxEnt (Warren et al., 2010). The  
523 observed *I* (*I<sub>o</sub>*) was calculated using the "ecological niche overlap" function  
524 implemented in ENMTools 1.3.

### 525 **Measurements of root traits**

526 Measurements were performed using 3.5-day-old seedlings grown on 1× MS  
527 medium with 1% sucrose in 16 h light/8 h darkness at 22°C ( $120 \mu\text{mol m}^{-2} \text{s}^{-1}$   
528 light intensity). The photon flux density was  $120 \mu\text{mol m}^{-2} \text{s}^{-1}$  (five cold white  
529 fluorescent Philips T5 28W/840 light bulbs and one warm white and yellow  
530 Philips T5 28W/830 light bulb). For root length measurement, roots were imaged  
531 using CCD scanners and measured using Fiji software (Schindelin et al., 2012).  
532 More than 15 seedlings were measured per line. To measure cellular traits, each  
533 seedling was treated with 10  $\mu\text{g/mL}$  propidium iodide (PI) solution for 2 min and  
534 washed with water twice. Each root tip was subsequently imaged under a Zeiss  
535 710 confocal microscope with a 20× objective. Both meristematic zone and  
536 mature cell sizes were measured using Fiji software (Schindelin et al., 2012).  
537 The length of the meristematic zone of each seedling was determined by  
538 measuring the distance from the quiescent center to the root cortex cell below

539 the first cortex cell that was twice as long as the cell below. The values for the  
540 left and right side of root longitudinal sections were averaged to obtain the length  
541 of the meristematic zone. The length of mature cells of each seedling was  
542 determined by averaging cell lengths of three mature cortex cells in the zone in  
543 which the xylem strands became visible. Five seedlings were measured per  
544 accession.

#### 545 **Mapping-by-sequencing of branch angle**

546 Plants of three natural accessions (Col-0, 3-2, and 29-8), and F<sub>2</sub> populations,  
547 were grown under LD conditions (16 h light/8 h dark, 120  $\mu\text{mol m}^{-2} \text{s}^{-1}$  light  
548 intensity) at 20°C. When the lateral shoot was longer than 5 cm, a protractor was  
549 used to measure the angle between the lateral branch and the main stem. The  
550 branch angles were measured in 2016 and 2017.

551 DNA was extracted from pooled leaves of 57 plants with a large branch angle  
552 (>85°) selected from a group of 684 F<sub>2</sub> plants of Col-0 × 29-8 with a large branch  
553 angle (>85°). The pooled DNA sample was sequenced using Illumina HiSeq  
554 Xten (150 bp pair-end reads, insert size of 450 bp). In total, 41,757,716 total  
555 reads (approximately 34.5-fold coverage) were mapped to the Col-0 genome  
556 (TAIR10), and SNPs were called using SHORE. The SNPs from pooled  
557 populations were used as markers to identify regions with an excess of  
558 homozygous alleles across the genome using SHOREmap (Schneeberger et al.,  
559 2009).

#### 560 **Statistical analysis**

561 Statistical analysis was performed using R (<http://www.r-project.org/>). Population  
562 structure-corrected ANOVAs were performed according to a previous study (Li et  
563 al., 2014). Significance was calculated by ANOVA with phenotypic differences  
564 between haplotypes on the residuals after subtracting the best linear unbiased  
565 predictors (BLUPs). EMMAX was used to estimate the BLUPs with a kinship  
566 matrix calculated using an SNP matrix available based on mixed linear models  
567 from the 1001 Genomes Project (Kang et al., 2010).

#### 568 **Accession numbers**

569 Sequence data from this article can be found in the GenBank/EMBL libraries  
570 under the following accession numbers: The *LAZY1* (AT5G14090) sequences of

571 Col-0, 3-2, and 29-8, under accession numbers MF621897–MF621899. Other  
572 genes studied in this article can be found in the Arabidopsis TAIR database  
573 (<https://www.arabidopsis.org>) under the following accession numbers: *CRK36*  
574 (AT4G04490), *KUK* (AT1G60370), *PRR5* (AT5G24470). The genome data for  
575 the pooled Col-0 × 29-8 F<sub>2</sub> populations have been deposited in the NCBI  
576 Sequence Read Archive (SRA) under accession number: SRR5947598.

577

## 578 **Supplemental Data**

579 **Supplemental Figure 1.** Accumulated false positive LoF rate.

580 **Supplemental Figure 2.** Variation in TE density (TE), nucleotide diversity ( $\pi$ ), and the  
581 density of genes with stop-gain (stop-gain), frameshift (Frameshift), splice site (Splice),  
582 and LoF variants across the chromosomes.

583 **Supplemental Figure 3.** Distribution of stop-gain and frameshift variants at gene  
584 positions in all 1,071 accessions, distributed in 5% bins.

585 **Supplemental Figure 4.** Gene families with multiple genes have more LoF  
586 variants than gene families with one or two genes based on amino acid  
587 sequence similarity in the Col-0 genome.

588 **Supplemental Figure 5.** The number of observed LoF variants based on GWAS  
589 compared with the distribution from permutations based on non-LoF.

590 **Supplemental Figure 6.** Histogram distributions for each bioclimatic trait used  
591 for GWAS.

592 **Supplemental Figure 7.** The narrow sense of estimated heritability in GWAS  
593 analysis.

594 **Supplemental Figure 8.** Genome-wide association study of LoF variants and 10  
595 climate parameters.

596 **Supplemental Figure 9.** Expression level variation of 30,000 random selected  
597 non-LoF SNP variants in 541 accessions based on RNA-seq data.

598 **Supplemental Figure 10.** Genome-wide association study of stop-gain variants and  
599 10 climate parameters.

600 **Supplemental Figure 11.** Examples of the origin of LoF variants in the *KUK* gene  
601 based on 129 accessions with phenotypic data.

602 **Supplemental Figure 12.** Geographic distribution of LoF variants for two genes  
603 in 1,071 accessions.

604 **Supplemental Figure 13.** Cellular root trait variation for *KUK* alleles with and  
605 without LoF.

606 **Supplemental Figure 14.** Alignment of Col-0, 3-2 and 29-8 LAZY1 amino acid  
607 sequences.

608 **Supplemental Figure 15.** The results of niche identity tests (*I*) for niches  
609 between non-LoF and LoF groups.

610 **Supplemental Data Set 1.** Summary of the 893 samples used in this study  
611 from the 1001 Genomes Project.

612 **Supplemental Data Set 2.** Summary of the 61 samples from African  
613 sequenced genomes.

614 **Supplemental Data Set 3.** Summary of the 117 samples from our genome  
615 project used in this study.

616 **Supplemental Data Set 4.** LoF mutation validation based on 10 assembled  
617 genomes.

618 **Supplemental Data Set 5.** Environmental variables used in GWAS and  
619 ecological niche modeling.

620 **Supplemental Data Set 6.** LoF variants with significant signals in GWAS.

621 **Supplemental Data Set 7.** Stop-gain variants with significant signals in GWAS.

622 **Supplemental Data Set 8.** The 129 samples used in the root length analysis  
623 from a previous study.

624 **Supplemental Data Set 9.** Cellular root trait measurements for 10 individuals  
625 per accession.

626 **Supplemental Data Set 10.** The 812 samples used in the flowering time  
627 analysis.

628 **Supplemental Data Set 11.** Markers used in the analysis of plants with large  
629 branch angles of the F<sub>2</sub> generation of Col-0 × 29-8.

630 **Supplemental Data Set 12.** Markers used in the analysis of plants with large  
631 branch angles of the F<sub>2</sub> generation of 3-2 × 29-8.

632 **Supplemental Data Set 13.** 530 Genes located in the 48 regions under positive  
633 selection in the Yangtze River basin population.

634 **Supplemental Data Set 14.** Genes with LoF variants located within the regions  
635 under positive selection in the Yangtze River population.

636 **Supplemental Data Set 15.** Primers used for the sequencing and amplification  
637 of *LAZY1*.

638 **Supplemental File 1.** ANOVA/test tables.

639

## 640 **ACKNOWLEDGEMENTS**

641 We would like to thank Magnus Nordborg and Huijing Ma for helpful suggestions  
642 about the study, and members of the Guo lab for suggestions and comments  
643 about this work. Especially we thank the anonymous reviewers for their help  
644 improving the manuscript. This work was supported by the Strategic Priority  
645 Research Program of the Chinese Academy of Sciences (XDB27010305), the  
646 Innovative Academy of Seed Design, Chinese Academy of Sciences, the  
647 National Natural Science Foundation of China (91731306 to Y.-L.G.), and start-  
648 up funds from the Salk Institute for Biological Research (W.B.).

649

## 650 **AUTHOR CONTRIBUTIONS**

651 Y.-L.G. conceived the study. Y.-C.X., J.-F.C., Y.-P.Z., Q.W., Y.E.Z., and Y.-L.G.  
652 analyzed the data. X.-M.N. and X.-X.L. conducted *LAZY1* experiments; W.H.  
653 and W.B. measured cellular traits for natural accessions. Y.-C.X. and Y.-L.G.  
654 wrote the paper with contributions from all authors.

655

## 656 **REFERENCES**

657 **Abdel-Samad, R., Zalzali, H., Rammah, C., Giraud, J., Naudin, C.,**  
658 **Dupasquier, S., Poulat, F., Boizet-Bonhoure, B., Lumbroso, S., Mouzat,**  
659 **K., Bonnans, C., Pignodel, C., Raynaud, P., Fort, P., Quittau-Prevostel,**  
660 **C., and Blache, P.** (2011). MiniSOX9, a dominant-negative variant in colon  
661 cancer cells. *Oncogene* **30**, 2493-2503.

662 **Albalat, R., and Canestro, C.** (2016). Evolution by gene loss. *Nat Rev Genet*  
663 **17**, 379-391.

664 **Amrad, A., Moser, M., Mandel, T., de Vries, M., Schuurink, R.C., Freitas, L.,**  
665 **and Kuhlemeier, C.** (2016). Gain and loss of floral scent production through  
666 changes in structural genes during pollinator-mediated speciation. *Curr Biol*  
667 **26**, 3303-3312.

668 **Benjamini, Y., and Hochberg, Y.** (1995). Controlling the false discovery rate - a  
669 practical and powerful approach to multiple testing. *J R Stat Soc B* **57**, 289-  
670 300.

- 671 **Birchler, J.A., and Veitia, R.A.** (2007). The gene balance hypothesis: from  
672 classical genetics to modern genomics. *Plant Cell* **19**, 395-402.
- 673 **Birchler, J.A., and Veitia, R.A.** (2012). Gene balance hypothesis: connecting  
674 issues of dosage sensitivity across biological disciplines. *Proceedings of the*  
675 *National Academy of Sciences of the United States of America* **109**, 14746-  
676 14753.
- 677 **Carvunis, A.R., Rolland, T., Wapinski, I., Calderwood, M.A., Yildirim, M.A.,**  
678 **Simonis, N., Charlotteaux, B., Hidalgo, C.A., Barbette, J., Santhanam, B.,**  
679 **Brar, G.A., Weissman, J.S., Regev, A., Thierry-Mieg, N., Cusick, M.E.,**  
680 **and Vidal, M.** (2012). Proto-genes and de novo gene birth. *Nature* **487**, 370-  
681 374.
- 682 **Chen, S., Zhang, Y.E., and Long, M.** (2010). New genes in *Drosophila* quickly  
683 become essential. *Science* **330**, 1682-1685.
- 684 **de Valles-Ibanez, G., Hernandez-Rodriguez, J., Prado-Martinez, J., Luisi, P.,**  
685 **Marques-Bonet, T., and Casals, F.** (2016). Genetic load of loss-of-function  
686 polymorphic variants in great apes. *Genome Biol Evol* **8**, 871-877.
- 687 **DePristo, M.A., Banks, E., Poplin, R., Garimella, K.V., Maguire, J.R., Hartl,**  
688 **C., Philippakis, A.A., del Angel, G., Rivas, M.A., Hanna, M., McKenna, A.,**  
689 **Fennell, T.J., Kernytsky, A.M., Sivachenko, A.Y., Cibulskis, K., Gabriel,**  
690 **S.B., Altshuler, D., and Daly, M.J.** (2011). A framework for variation  
691 discovery and genotyping using next-generation DNA sequencing data. *Nat*  
692 *Genet* **43**, 491-498.
- 693 **Derrien, T., Estelle, J., Marco Sola, S., Knowles, D.G., Raineri, E., Guigo, R.,**  
694 **and Ribeca, P.** (2012). Fast computation and applications of genome  
695 mappability. *PLoS One* **7**, e30377.
- 696 **Durvasula, A., Fulgione, A., Gutaker, R.M., Alacakaptan, S.I., Flood, P.J.,**  
697 **Neto, C., Tsuchimatsu, T., Burbano, H.A., Pico, F.X., Alonso-Blanco, C.,**  
698 **and Hancock, A.M.** (2017). African genomes illuminate the early history and  
699 transition to selfing in *Arabidopsis thaliana*. *Proceedings of the National*  
700 *Academy of Sciences of the United States of America* **114**, 5213-5218.
- 701 **Edgar, R.C.** (2004). MUSCLE: a multiple sequence alignment method with  
702 reduced time and space complexity. *BMC Bioinformatics* **5**, 113.
- 703 **Edger, P.P., and Pires, J.C.** (2009). Gene and genome duplications: the impact  
704 of dosage-sensitivity on the fate of nuclear genes. *Chromosome Res* **17**,  
705 699-717.
- 706 **Gan, X., Stegle, O., Behr, J., Steffen, J.G., Drewe, P., Hildebrand, K.L.,**  
707 **Lyngsoe, R., Schultheiss, S.J., Osborne, E.J., Sreedharan, V.T., Kahles,**  
708 **A., Bohnert, R., Jean, G., Derwent, P., Kersey, P., Belfield, E.J., Harberd,**  
709 **N.P., Kemen, E., Toomajian, C., Kover, P.X., Clark, R.M., Ratsch, G., and**  
710 **Mott, R.** (2011). Multiple reference genomes and transcriptomes for  
711 *Arabidopsis thaliana*. *Nature* **477**, 419-423.

- 712 **Goldman-Huertas, B., Mitchell, R.F., Lapoint, R.T., Faucher, C.P.,**  
713 **Hildebrand, J.G., and Whiteman, N.K.** (2015). Evolution of herbivory in  
714 *Drosophilidae* linked to loss of behaviors, antennal responses, odorant  
715 receptors, and ancestral diet. *Proceedings of the National Academy of*  
716 *Sciences of the United States of America* **112**, 3026-3031.
- 717 **Green, C., Willoughby, J., Study, D.D.D., and Balasubramanian, M.** (2017).  
718 *De novo SETD5* loss-of-function variant as a cause for intellectual disability  
719 in a 10-year old boy with an aberrant blind ending bronchus. *Am J Med*  
720 *Genet A* **173**, 3165-3171.
- 721 **Greenberg, A.J., Moran, J.R., Coyne, J.A., and Wu, C.I.** (2003). Ecological  
722 adaptation during incipient speciation revealed by precise gene replacement.  
723 *Science* **302**, 1754-1757.
- 724 **Gujas, B., Alonso-Blanco, C., and Hardtke, C.S.** (2012). Natural *Arabidopsis*  
725 *brx* loss-of-function alleles confer root adaptation to acidic soil. *Curr Biol* **22**,  
726 1962-1968.
- 727 **Guo, Y.L.** (2013). Gene family evolution in green plants with emphasis on the  
728 origination and evolution of *Arabidopsis thaliana* genes. *Plant J* **73**, 941-951.
- 729 **Hao, Y., Washburn, J.D., Rosenthal, J., Nielsen, B., Lyons, E., Edger, P.P.,**  
730 **Pires, J.C., and Conant, G.C.** (2018). Patterns of population variation in two  
731 paleopolyploid eudicot lineages suggest that dosage-based selection on  
732 homeologs is long-lived. *Genome Biol Evol* **10**, 999-1011.
- 733 **Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G., and Jarvis, A.** (2005).  
734 Very high resolution interpolated climate surfaces for global land areas. *Int J*  
735 *Climatol* **25**, 1965-1978.
- 736 **Hoballah, M.E., Gubitz, T., Stuurman, J., Broger, L., Barone, M., Mandel, T.,**  
737 **Dell'Olivo, A., Arnold, M., and Kuhlemeier, C.** (2007). Single gene-  
738 mediated shift in pollinator attraction in *Petunia*. *Plant Cell* **19**, 779-790.
- 739 **Hodgson, J.A., Pickrell, J.K., Pearson, L.N., Quillen, E.E., Prista, A., Rocha,**  
740 **J., Soodyall, H., Shriver, M.D., and Perry, G.H.** (2014). Natural selection for  
741 the Duffy-null allele in the recently admixed people of Madagascar. *Proc Biol*  
742 *Sci* **281**, 20140930.
- 743 **Hottes, A.K., Freddolino, P.L., Khare, A., Donnell, Z.N., Liu, J.C., and**  
744 **Tavazoie, S.** (2013). Bacterial adaptation through loss of function. *Plos*  
745 *Genet* **9**, e1003617.
- 746 **Hou, J., Shi, X., Chen, C., Islam, M.S., Johnson, A.F., Kanno, T., Huettel, B.,**  
747 **Yen, M.R., Hsu, F.M., Ji, T., Chen, P.Y., Matzke, M., Matzke, A.J.M.,**  
748 **Cheng, J., and Birchler, J.A.** (2018). Global impacts of chromosomal  
749 imbalance on gene expression in *Arabidopsis* and other taxa. *Proceedings of*  
750 *the National Academy of Sciences of the United States of America* **115**,  
751 E11321-E11330.

- 752 **Kang, H.M., Sul, J.H., Service, S.K., Zaitlen, N.A., Kong, S.Y., Freimer, N.B.,**  
753 **Sabatti, C., and Eskin, E.** (2010). Variance component model to account for  
754 sample structure in genome-wide association studies. *Nat Genet* **42**, 348-  
755 354.
- 756 **Kawakatsu, T., Huang, S.C., Jupe, F., Sasaki, E., Schmitz, R.J., Urich, M.A.,**  
757 **Castanon, R., Nery, J.R., Barragan, C., He, Y., Chen, H., Dubin, M., Lee,**  
758 **C.R., Wang, C., Bemm, F., Becker, C., O'Neil, R., O'Malley, R.C.,**  
759 **Quarless, D.X., Genomes, C., Schork, N.J., Weigel, D., Nordborg, M., and**  
760 **Ecker, J.R.** (2016). Epigenomic diversity in a global collection of *Arabidopsis*  
761 *thaliana* accessions. *Cell* **166**, 492-505.
- 762 **Kremling, K.A.G., Chen, S.Y., Su, M.H., Lepak, N.K., Romay, M.C., Swarts,**  
763 **K.L., Lu, F., Lorant, A., Bradbury, P.J., and Buckler, E.S.** (2018).  
764 Dysregulation of expression correlates with rare-allele burden and fitness  
765 loss in maize. *Nature* **555**, 520-523.
- 766 **Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell,**  
767 **T., O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B.,**  
768 **Tukiainen, T., Birnbaum, D.P., Kosmicki, J.A., Duncan, L.E., Estrada, K.,**  
769 **Zhao, F., Zou, J., Pierce-Hoffman, E., Berghout, J., Cooper, D.N.,**  
770 **Deflaux, N., DePristo, M., Do, R., Flannick, J., Fromer, M., Gauthier, L.,**  
771 **Goldstein, J., Gupta, N., Howrigan, D., Kiezun, A., Kurki, M.I.,**  
772 **Moonshine, A.L., Natarajan, P., Orozco, L., Peloso, G.M., Poplin, R.,**  
773 **Rivas, M.A., Ruano-Rubio, V., Rose, S.A., Ruderfer, D.M., Shakir, K.,**  
774 **Stenson, P.D., Stevens, C., Thomas, B.P., Tiao, G., Tusie-Luna, M.T.,**  
775 **Weisburd, B., Won, H.H., Yu, D., Altshuler, D.M., Ardissino, D., Boehnke,**  
776 **M., Danesh, J., Donnelly, S., Elosua, R., Florez, J.C., Gabriel, S.B., Getz,**  
777 **G., Glatt, S.J., Hultman, C.M., Kathiresan, S., Laakso, M., McCarroll, S.,**  
778 **McCarthy, M.I., McGovern, D., McPherson, R., Neale, B.M., Palotie, A.,**  
779 **Purcell, S.M., Saleheen, D., Scharf, J.M., Sklar, P., Sullivan, P.F.,**  
780 **Tuomilehto, J., Tsuang, M.T., Watkins, H.C., Wilson, J.G., Daly, M.J.,**  
781 **MacArthur, D.G., and Exome Aggregation, C.** (2016). Analysis of protein-  
782 coding genetic variation in 60,706 humans. *Nature* **536**, 285-291.
- 783 **Li, P., Filiault, D., Box, M.S., Kerdaffrec, E., van Oosterhout, C., Wilczek,**  
784 **A.M., Schmitt, J., McMullan, M., Bergelson, J., Nordborg, M., and Dean,**  
785 **C.** (2014). Multiple *FLC* haplotypes defined by independent cis-regulatory  
786 variation underpin life history diversity in *Arabidopsis thaliana*. *Genes Dev* **28**,  
787 1635-1640.
- 788 **Li, Z.W., Chen, X., Wu, Q., Hagmann, J., Han, T.S., Zou, Y.P., Ge, S., and**  
789 **Guo, Y.L.** (2016). On the origin of de novo genes in *Arabidopsis thaliana*  
790 populations. *Genome Biol Evol* **8**, 2190-2202.
- 791 **Lipka, A.E., Tian, F., Wang, Q., Peiffer, J., Li, M., Bradbury, P.J., Gore, M.A.,**  
792 **Buckler, E.S., and Zhang, Z.** (2012). GAPIT: genome association and  
793 prediction integrated tool. *Bioinformatics* **28**, 2397-2399.
- 794 **MacArthur, D.G., Balasubramanian, S., Frankish, A., Huang, N., Morris, J.,**  
795 **Walter, K., Jostins, L., Habegger, L., Pickrell, J.K., Montgomery, S.B.,**

- 796 **Albers, C.A., Zhang, Z.D., Conrad, D.F., Lunter, G., Zheng, H., Ayub, Q.,**  
797 **DePristo, M.A., Banks, E., Hu, M., Handsaker, R.E., Rosenfeld, J.A.,**  
798 **Fromer, M., Jin, M., Mu, X.J., Khurana, E., Ye, K., Kay, M., Saunders, G.I.,**  
799 **Suner, M.M., Hunt, T., Barnes, I.H., Amid, C., Carvalho-Silva, D.R.,**  
800 **Bignell, A.H., Snow, C., Yngvadottir, B., Bumpstead, S., Cooper, D.N.,**  
801 **Xue, Y., Romero, I.G., Genomes Project, C., Wang, J., Li, Y., Gibbs, R.A.,**  
802 **McCarroll, S.A., Dermitzakis, E.T., Pritchard, J.K., Barrett, J.C., Harrow,**  
803 **J., Hurles, M.E., Gerstein, M.B., and Tyler-Smith, C. (2012). A systematic**  
804 **survey of loss-of-function variants in human protein-coding genes. Science**  
805 **335, 823-828.**
- 806 **Marco-Sola, S., Sammeth, M., Guigo, R., and Ribeca, P. (2012). The GEM**  
807 **mapper: fast, accurate and versatile alignment by filtration. Nat Methods 9,**  
808 **1185-1188.**
- 809 **McLaren, W., Pritchard, B., Rios, D., Chen, Y., Flicek, P., and Cunningham,**  
810 **F. (2010). Deriving the consequences of genomic variants with the Ensembl**  
811 **API and SNP Effect Predictor. Bioinformatics 26, 2069-2070.**
- 812 **McLysaght, A., and Guerzoni, D. (2015). New genes from non-coding**  
813 **sequence: the role of de novo protein-coding genes in eukaryotic**  
814 **evolutionary innovation. Philos Trans R Soc Lond B Biol Sci 370, 20140332.**
- 815 **Meijon, M., Satbhai, S.B., Tsuchimatsu, T., and Busch, W. (2014). Genome-**  
816 **wide association study using cellular traits identifies a new regulator of root**  
817 **development in Arabidopsis. Nat Genet 46, 77-81.**
- 818 **Meinke, D., Muralla, R., Sweeney, C., and Dickerman, A. (2008). Identifying**  
819 **essential genes in Arabidopsis thaliana. Trends Plant Sci 13, 483-491.**
- 820 **Meng, D., Dubin, M., Zhang, P., Osborne, E.J., Stegle, O., Clark, R.M., and**  
821 **Nordborg, M. (2016). Limited contribution of DNA methylation variation to**  
822 **expression regulation in Arabidopsis thaliana. Plos Genet 12, e1006141.**
- 823 **Michael, T.P., Jupe, F., Bemm, F., Motley, S.T., Sandoval, J.P., Lanz, C.,**  
824 **Loudet, O., Weigel, D., and Ecker, J.R. (2018). High contiguity Arabidopsis**  
825 **thaliana genome assembly with a single nanopore flow cell. Nat Commun 9,**  
826 **541.**
- 827 **Nakamichi, N., Takao, S., Kudo, T., Kiba, T., Wang, Y., Kinoshita, T., and**  
828 **Sakakibara, H. (2016). Improvement of Arabidopsis biomass and cold,**  
829 **drought and salinity stress tolerance by modified circadian clock-associated**  
830 **PSEUDO-RESPONSE REGULATORS. Plant Cell Physiol 57, 1085-1097.**
- 831 **Narasimhan, V.M., Hunt, K.A., Mason, D., Baker, C.L., Karczewski, K.J.,**  
832 **Barnes, M.R., Barnett, A.H., Bates, C., Bellary, S., Bockett, N.A., Giorda,**  
833 **K., Griffiths, C.J., Hemingway, H., Jia, Z., Kelly, M.A., Khawaja, H.A., Lek,**  
834 **M., McCarthy, S., McEachan, R., O'Donnell-Luria, A., Paigen, K.,**  
835 **Parisinos, C.A., Sheridan, E., Southgate, L., Tee, L., Thomas, M., Xue, Y.,**  
836 **Schnall-Levin, M., Petkov, P.M., Tyler-Smith, C., Maher, E.R., Trembath,**  
837 **R.C., MacArthur, D.G., Wright, J., Durbin, R., and van Heel, D.A. (2016).**

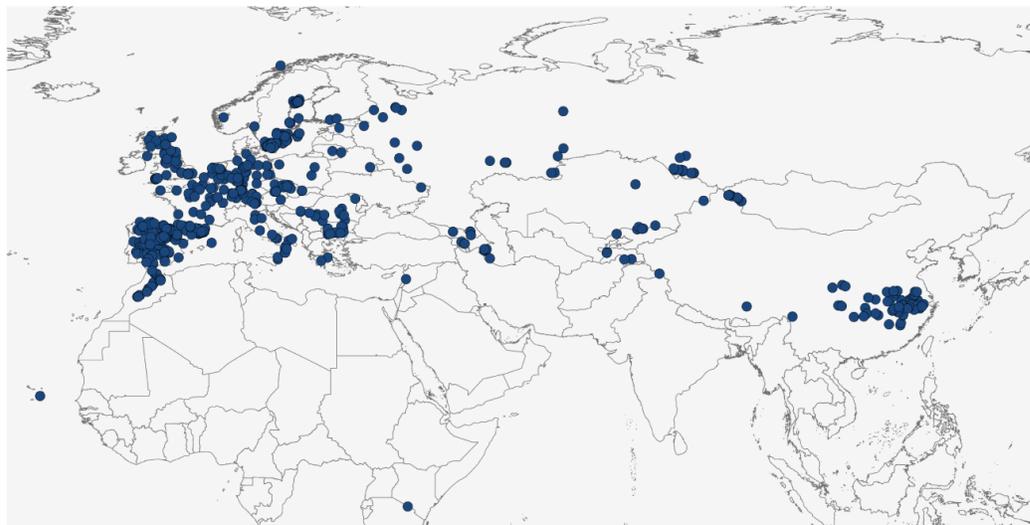
- 838 Health and population effects of rare gene knockouts in adult humans with  
839 related parents. *Science* **352**, 474-477.
- 840 **Olson, M.V.** (1999). When less is more: gene loss as an engine of evolutionary  
841 change. *Am J Hum Genet* **64**, 18-23.
- 842 **Palmieri, N., Kosiol, C., and Schlötterer, C.** (2014). The life cycle of  
843 *Drosophila* orphan genes. *eLife* **3**, e01311.
- 844 **Phillips, S.J., and Dudik, M.** (2008). Modeling of species distributions with  
845 Maxent: new extensions and a comprehensive evaluation. *Ecography* **31**,  
846 161-175.
- 847 **Phillips, S.J., Anderson, R.P., and Schapire, R.E.** (2006). Maximum entropy  
848 modeling of species geographic distributions. *Ecol Model* **190**, 231-259.
- 849 **Prachumwat, A., and Li, W.H.** (2008). Gene number expansion and contraction  
850 in vertebrate genomes with respect to invertebrate genomes. *Genome Res*  
851 **18**, 221-232.
- 852 **Prieto-Godino, L.L., Rytz, R., Bargeton, B., Abuin, L., Arguello, J.R., Peraro,  
853 M.D., and Benton, R.** (2016). Olfactory receptor pseudo-pseudogenes.  
854 *Nature* **539**, 93-97.
- 855 **Saleheen, D., Natarajan, P., Armean, I.M., Zhao, W., Rasheed, A., Khetarpal,  
856 S.A., Won, H.H., Karczewski, K.J., O'Donnell-Luria, A.H., Samocha, K.E.,  
857 Weisburd, B., Gupta, N., Zaidi, M., Samuel, M., Imran, A., Abbas, S.,  
858 Majeed, F., Ishaq, M., Akhtar, S., Trindade, K., Mucksavage, M., Qamar,  
859 N., Zaman, K.S., Yaqoob, Z., Saghir, T., Rizvi, S.N., Memon, A., Hayyat  
860 Mallick, N., Ishaq, M., Rasheed, S.Z., Memon, F.U., Mahmood, K.,  
861 Ahmed, N., Do, R., Krauss, R.M., MacArthur, D.G., Gabriel, S., Lander,  
862 E.S., Daly, M.J., Frossard, P., Danesh, J., Rader, D.J., and Kathiresan, S.**  
863 (2017). Human knockouts and phenotypic analysis in a cohort with a high  
864 rate of consanguinity. *Nature* **544**, 235-239.
- 865 **Sas, C., Muller, F., Kappel, C., Kent, T.V., Wright, S.I., Hilker, M., and  
866 Lenhard, M.** (2016). Repeated inactivation of the first committed enzyme  
867 underlies the loss of benzaldehyde emission after the selfing transition in  
868 *Capsella*. *Curr Biol* **26**, 3313-3319.
- 869 **Schild, L., Canessa, C.M., Shimkets, R.A., Gautschi, I., Lifton, R.P., and  
870 Rossier, B.C.** (1995). A mutation in the epithelial sodium channel causing  
871 Liddle disease increases channel activity in the *Xenopus laevis* oocyte  
872 expression system. *Proceedings of the National Academy of Sciences of the  
873 United States of America* **92**, 5699-5703.
- 874 **Schindelin, J., Arganda-Carreras, I., Frise, E., Kaynig, V., Longair, M.,  
875 Pietzsch, T., Preibisch, S., Rueden, C., Saalfeld, S., Schmid, B., Tinevez,  
876 J.Y., White, D.J., Hartenstein, V., Eliceiri, K., Tomancak, P., and  
877 Cardona, A.** (2012). Fiji: an open-source platform for biological-image  
878 analysis. *Nat Methods* **9**, 676-682.

- 879 **Schneeberger, K., Ossowski, S., Lanz, C., Juul, T., Petersen, A.H., Nielsen,**  
880 **K.L., Jorgensen, J.E., Weigel, D., and Andersen, S.U.** (2009).  
881 SHOREmap: simultaneous mapping and mutation identification by deep  
882 sequencing. *Nat Methods* **6**, 550-551.
- 883 **Slater, G.S., and Birney, E.** (2005). Automated generation of heuristics for  
884 biological sequence comparison. *BMC Bioinformatics* **6**, 31.
- 885 **Song, X.J., Huang, W., Shi, M., Zhu, M.Z., and Lin, H.X.** (2007). A QTL for rice  
886 grain width and weight encodes a previously unknown RING-type E3  
887 ubiquitin ligase. *Nat Genet* **39**, 623-630.
- 888 **Swets, J.A.** (1988). Measuring the accuracy of diagnostic systems. *Science*  
889 **240**, 1285-1293.
- 890 **Tanaka, H., Osakabe, Y., Katsura, S., Mizuno, S., Maruyama, K., Kusakabe,**  
891 **K., Mizoi, J., Shinozaki, K., and Yamaguchi-Shinozaki, K.** (2012). Abiotic  
892 stress-inducible receptor-like kinases negatively control ABA signaling in  
893 *Arabidopsis*. *Plant J* **70**, 599-613.
- 894 **Tang, C., Toomajian, C., Sherman-Broyles, S., Plagnol, V., Guo, Y.L., Hu,**  
895 **T.T., Clark, R.M., Nasrallah, J.B., Weigel, D., and Nordborg, M.** (2007).  
896 The evolution of selfing in *Arabidopsis thaliana*. *Science* **317**, 1070-1072.
- 897 **The 1001 Genomes Consortium.** (2016). 1,135 genomes reveal the global  
898 pattern of polymorphism in *Arabidopsis thaliana*. *Cell* **166**, 481-491.
- 899 **Tian, T., Liu, Y., Yan, H., You, Q., Yi, X., Du, Z., Xu, W., and Su, Z.** (2017).  
900 agriGO v2.0: a GO analysis toolkit for the agricultural community, 2017  
901 update. *Nucleic Acids Res* **45**, W122-W129.
- 902 **Wagner, A.** (2005). Distributed robustness versus redundancy as causes of  
903 mutational robustness. *Bioessays* **27**, 176-188.
- 904 **Warren, D.L., Glor, R.E., and Turelli, M.** (2010). ENMTools: A toolbox for  
905 comparative studies of environmental niche models. *Ecography* **33**, 607-611.
- 906 **Will, J.L., Kim, H.S., Clarke, J., Painter, J.C., Fay, J.C., and Gasch, A.P.**  
907 (2010). Incipient balancing selection through adaptive loss of aquaporins in  
908 natural *Saccharomyces cerevisiae* populations. *Plos Genet* **6**, e1000893.
- 909 **Wu, W., Liu, X., Wang, M., Meyer, R.S., Luo, X., Ndjiondjop, M.N., Tan, L.,**  
910 **Zhang, J., Wu, J., Cai, H., Sun, C., Wang, X., Wing, R.A., and Zhu, Z.**  
911 (2017). A single-nucleotide polymorphism causes smaller grain size and loss  
912 of seed shattering during African rice domestication. *Nat Plants* **3**, 17064.
- 913 **Xie, J., Murone, M., Luoh, S.M., Ryan, A., Gu, Q., Zhang, C., Bonifas, J.M.,**  
914 **Lam, C.W., Hynes, M., Goddard, A., Rosenthal, A., Epstein, E.H., Jr., and**  
915 **de Sauvage, F.J.** (1998). Activating *Smoothed* mutations in sporadic  
916 basal-cell carcinoma. *Nature* **391**, 90-92.

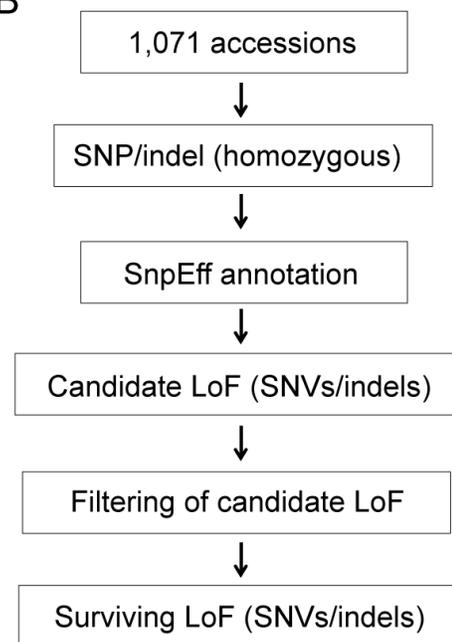
- 917 **Yang, H., He, B.Z., Ma, H., Tsauro, S.C., Ma, C., Wu, Y., Ting, C.T., and Zhang,**  
918 **Y.E.** (2015). Expression profile and gene age jointly shaped the genome-  
919 wide distribution of premature termination codons in a *Drosophila*  
920 *melanogaster* population. *Mol Biol Evol* **32**, 216-228.
- 921 **Yoshihara, T., Spalding, E.P., and Iino, M.** (2013). *AtLAZY1* is a signaling  
922 component required for gravitropism of the *Arabidopsis thaliana*  
923 inflorescence. *Plant J* **74**, 267-279.
- 924 **Zapata, L., Ding, J., Willing, E.M., Hartwig, B., Bezdán, D., Jiao, W.B., Patel,**  
925 **V., Velikkakam James, G., Koornneef, M., Ossowski, S., and**  
926 **Schneeberger, K.** (2016). Chromosome-level assembly of *Arabidopsis*  
927 *thaliana* Ler reveals the extent of translocation and inversion polymorphisms.  
928 *Proceedings of the National Academy of Sciences of the United States of*  
929 *America* **113**, E4052-4060.
- 930 **Zhang, Z., Ersoz, E., Lai, C.Q., Todhunter, R.J., Tiwari, H.K., Gore, M.A.,**  
931 **Bradbury, P.J., Yu, J., Arnett, D.K., Ordovas, J.M., and Buckler, E.S.**  
932 (2010). Mixed linear model approach adapted for genome-wide association  
933 studies. *Nat Genet* **42**, 355-360.
- 934 **Zhao, L., Saelao, P., Jones, C.D., and Begun, D.J.** (2014). Origin and spread  
935 of de novo genes in *Drosophila melanogaster* populations. *Science* **343**, 769-  
936 772.
- 937 **Zou, Y.P., Hou, X.H., Wu, Q., Chen, J.F., Li, Z.W., Han, T.S., Niu, X.M., Yang,**  
938 **L., Xu, Y.C., Zhang, J., Zhang, F.M., Tan, D., Tian, Z., Gu, H., and Guo,**  
939 **Y.L.** (2017). Adaptation of *Arabidopsis thaliana* to the Yangtze River basin.  
940 *Genome Biol* **18**, 239.
- 941 **Zufall, R.A., and Rausher, M.D.** (2004). Genetic changes associated with floral  
942 adaptation restrict future evolutionary potential. *Nature* **428**, 847-850.  
943  
944  
945  
946  
947 Table 1. Total number of different LoF variants in the 1,071 genomes. Numbers  
948 in parentheses indicate the gene number; some genes have multiple LoF  
949 variants.

Variant type	Before filtering		After filtering	
	Total	Per sample	Total	Per sample
Stop-gain	28,944 (10,350)	432.1 (403.5)	17,453 (7,264)	214.7 (214.7)
Frameshift	72,335 (14,994)	1,448.4 (1,178.9)	37,935 (10,800)	512.2 (512.2)
Splice site	10,414 (5,769)	277.5 (264.9)	5,431 (3,449)	103.0 (103.0)
Total	111,693 (17,957)	2,157.9 (1,732.8)	60,819 (12,918)	829.9 (829.9)

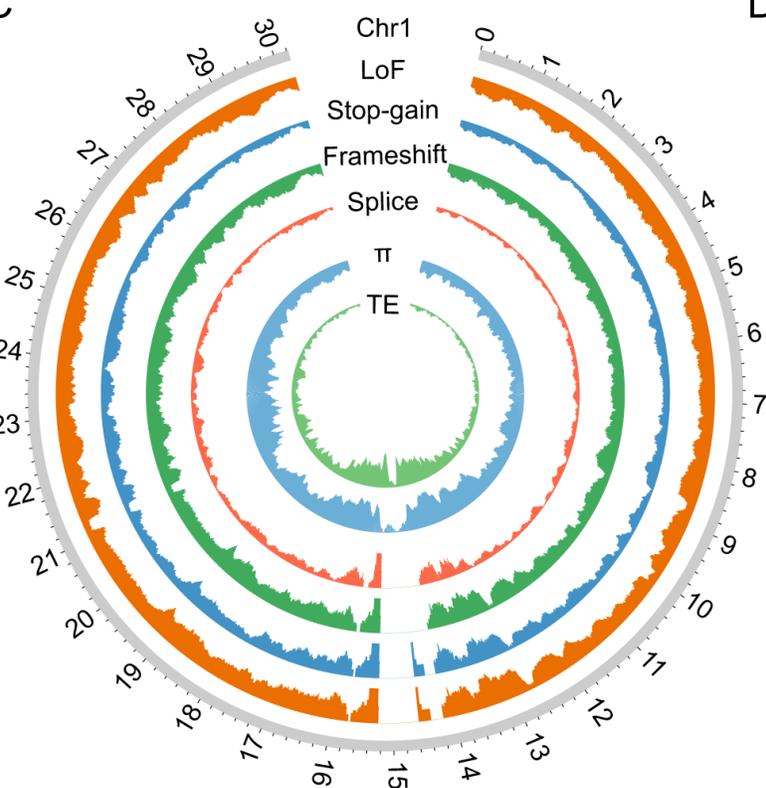
A



B



C



D

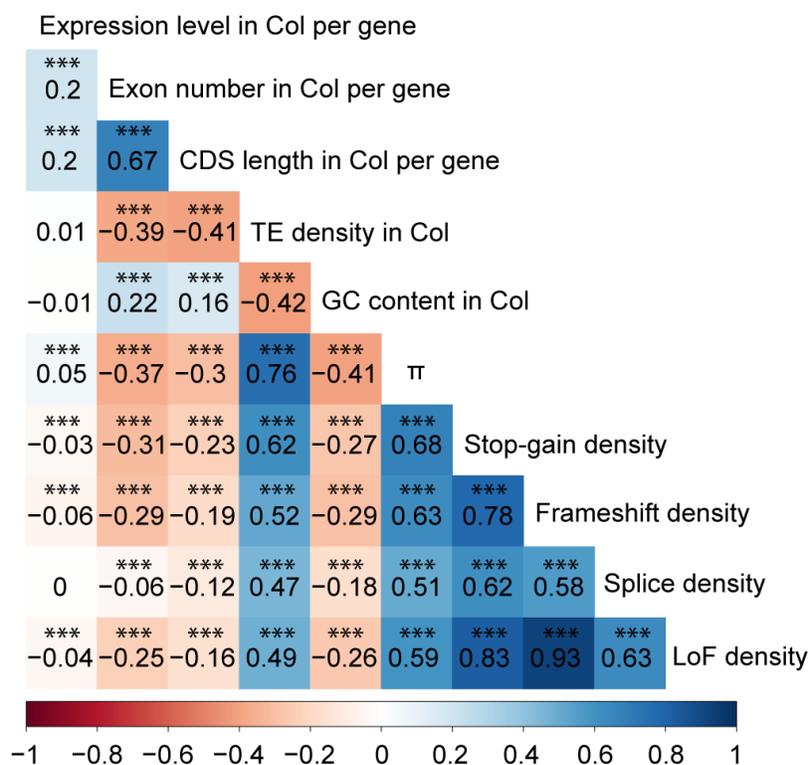


Figure 1. Identification of LoF mutations in 1,071 *Arabidopsis thaliana* genomes. (A) Geographic distribution of the 1,071 accessions. (B) Process to identify LoF variants. (C) Variation in TE density (TE), nucleotide diversity ( $\pi$ ), and the density of genes with stop-gain (Stop-gain), frameshift (Frameshift), splice site (Splice), and LoF variants across chromosome 1. The other four chromosomes are shown in Supplemental Figure 1. (D) Spearman correlations adjusted for multiple testing between different gene features across the whole genome, including expression level, exon number, and CDS length per gene in Col-0; GC content in Col-0; TE density in Col-0; nucleotide diversity ( $\pi$ ); and the density of genes with stop-gain (Stop-gain density), frameshift (Frameshift density), splice site (Splice density), and LoF (LoF density). Values marked with asterisks indicate strong correlation coefficients (\*\*\*)  $p < 0.001$ .

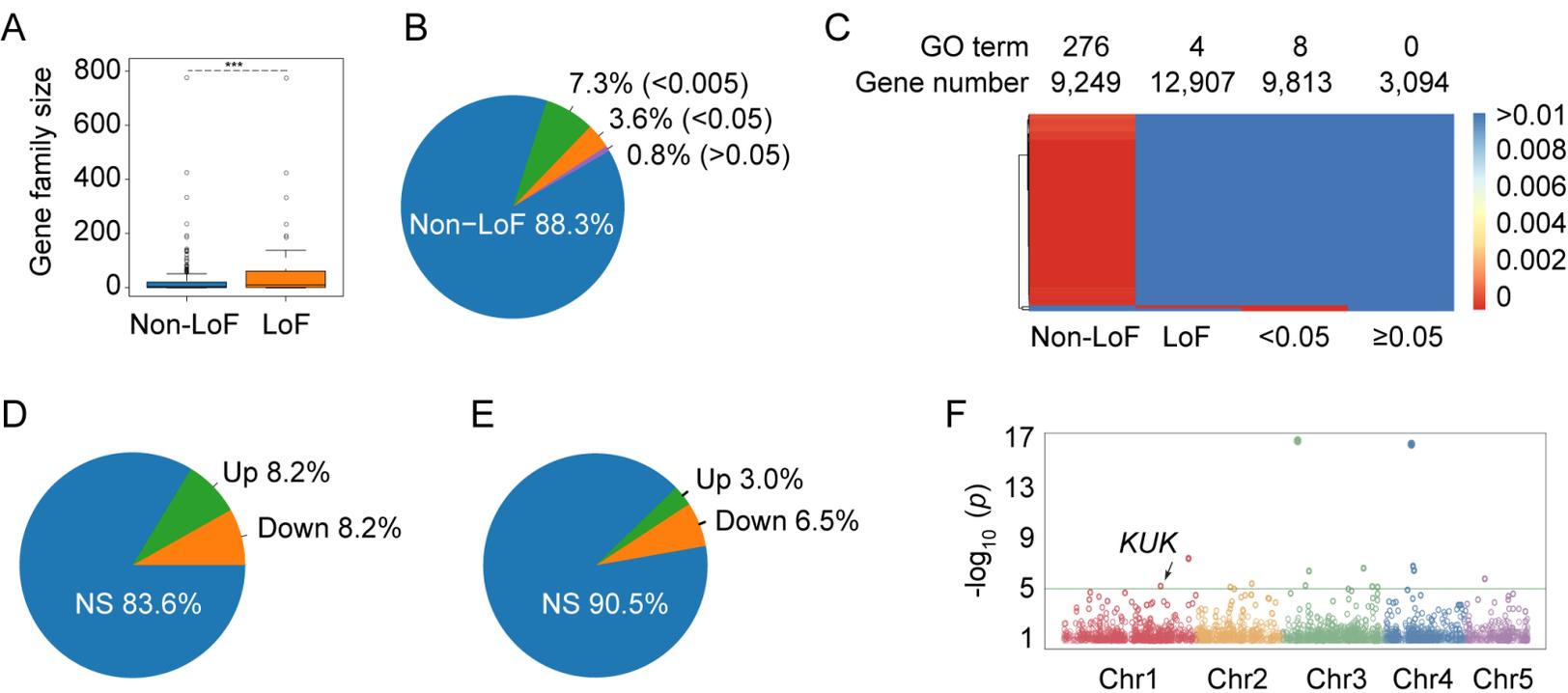


Figure 2. Functional analyses of genes with and without LoF variants. (A) Genes from larger gene families tend to have LoF variants based on the Col-0 genome. In the boxplot, the box equals the difference between the 75th and 25th percentiles. The thick line indicates the median. The two whiskers indicate 1.5 interquartile range of the lower quartile or the upper quartile. Outliers are plotted as individual points. \*\*\*  $p < 0.001$ . (B) LoF mutations are under-represented in lethal genes in *A. thaliana*. “<0.005” indicates a minor allele frequency (MAF) less than 0.5%, “<0.05” indicates MAF < 5%, and “>0.05” indicates MAF > 5%. (C) Gene Ontology enrichment analysis of genes with LoF, with different LoF allele frequencies, and those without LoF variants across the 1,071 genomes. Color indicates p values of GO enrichment analysis. (D) Expression level variation between alleles with and without LoF variants of genes with significant GWAS signals based on RNA-seq data for 541 accessions. (E) Expression level variation between alleles with and without LoF variants (6,675 genes) based on RNA-seq data for 541 accessions. NS indicates insignificant difference. (F) Genome-wide association study of stop-gain variants with regard to precipitation in the wettest month.

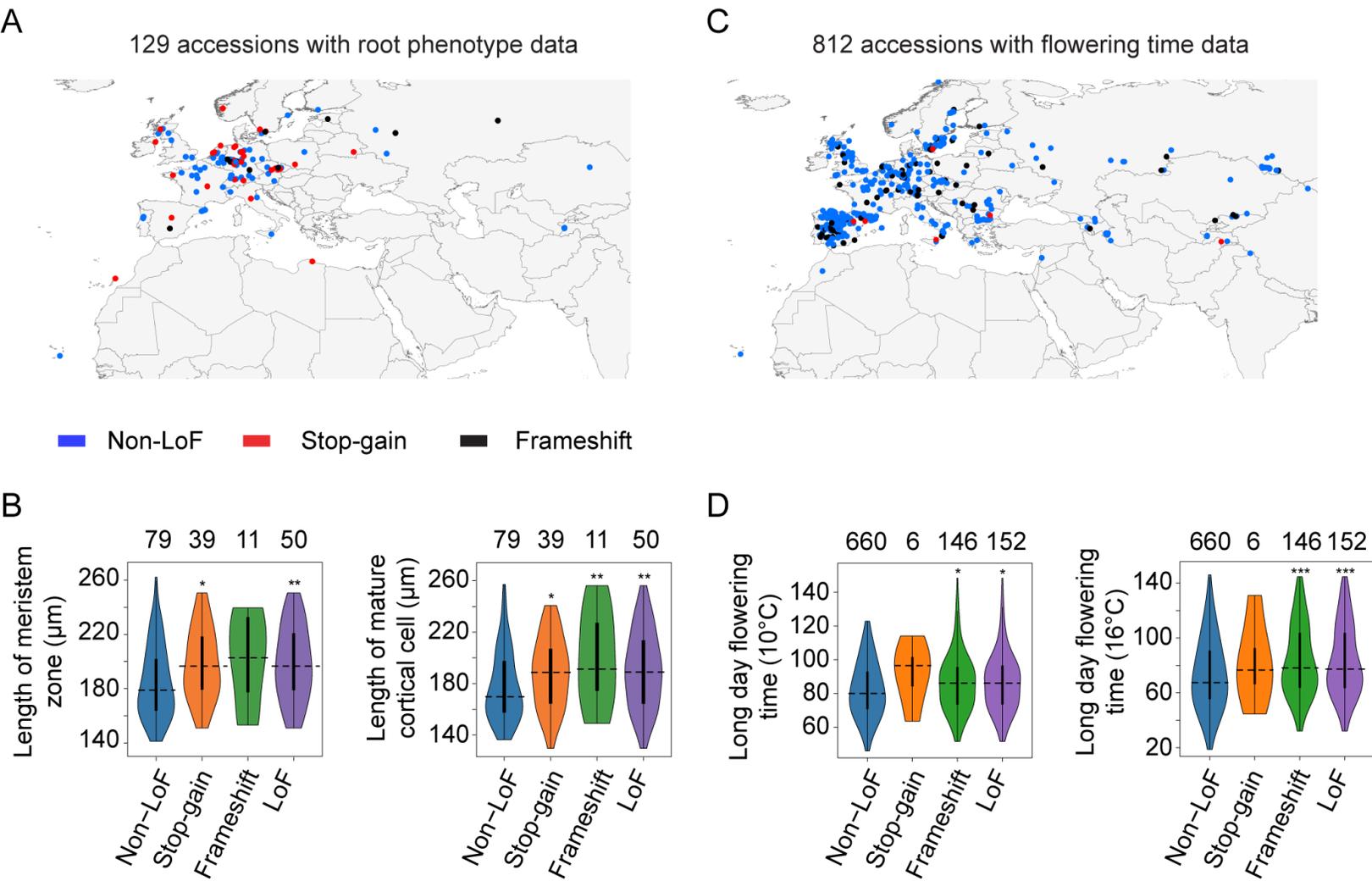
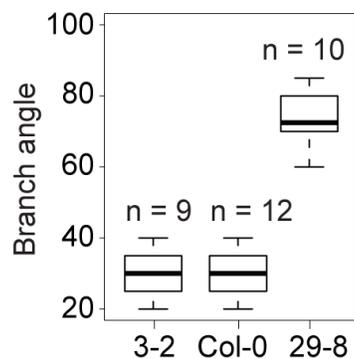
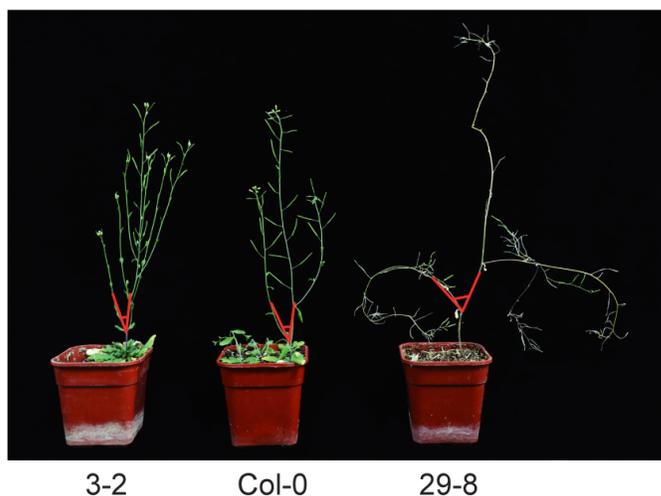
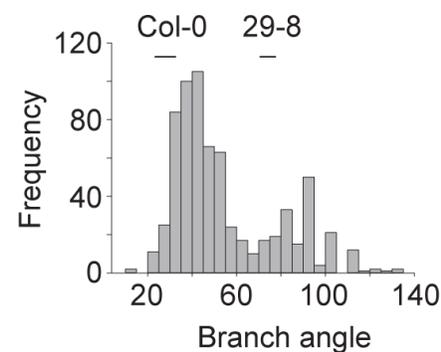


Figure 3. Functional differentiation of alleles with or without LoF. (A) Geographic distribution of different KUK alleles in 129 accessions with phenotypic data. (B) Phenotypic variation for KUK alleles with and without LoF. (C) Geographic distribution of PRR5 alleles in 812 accessions with flowering time data. (D) Flowering time variation for PRR5 alleles with and without LoF. The numbers above each violin plot indicate the number of accessions containing the particular variant type. The dashed line indicates the median value of the phenotypic trait. Population structure-corrected ANOVA, \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

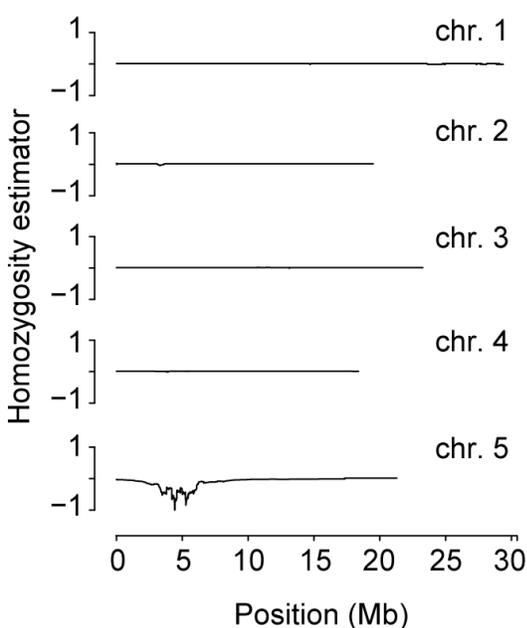
A



B



C



D

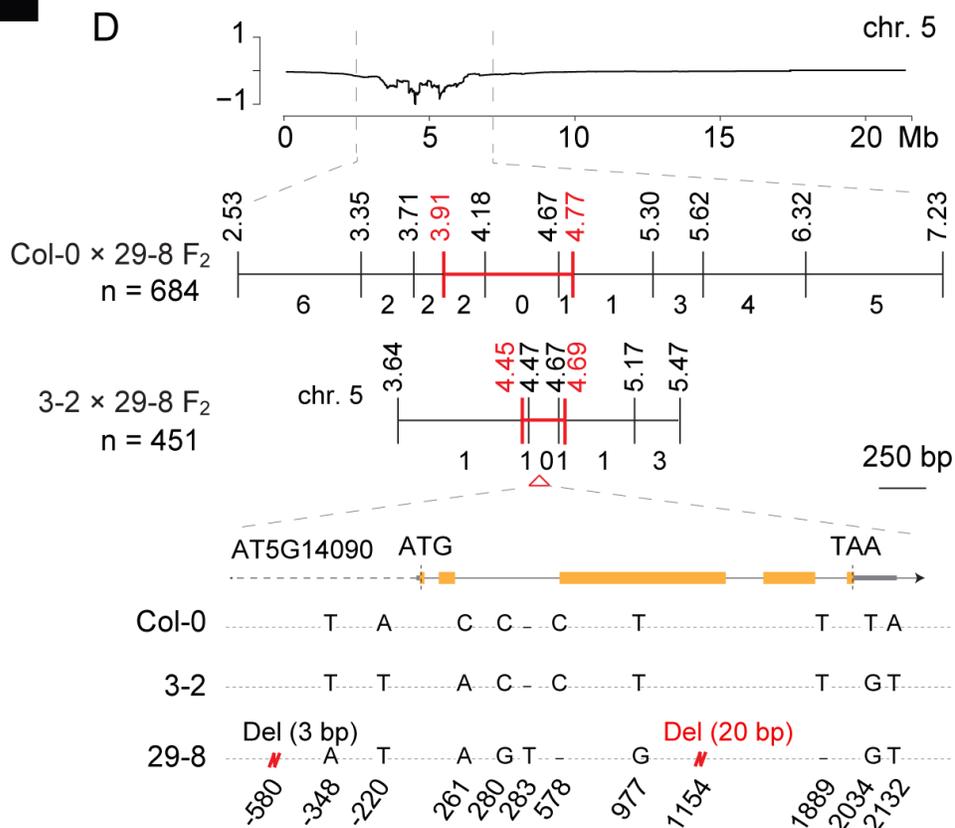
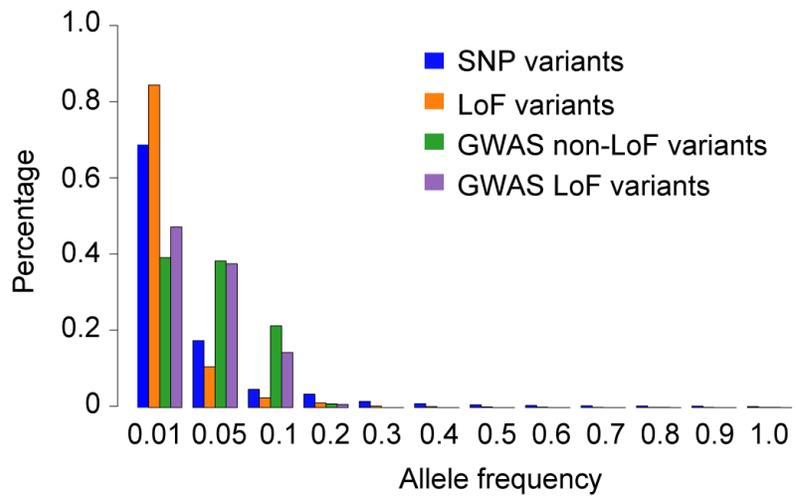
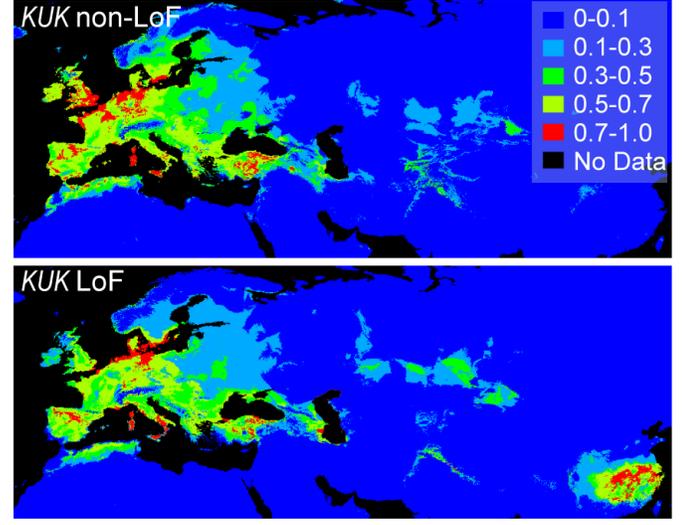


Figure 4. Frameshift in LAZY1 increases branch angle. (A) Representative accessions with different branch angles. The red triangles show the tangent lines drawn for the measurement. The boxplot (right) indicates distribution of branch angles in different accessions. The borders of the box indicate the 75th and 25th percentiles. The thick line marks the median. The two whiskers extend to the most extreme data points.  $n$  indicates the number of individuals measured for each accession. (B) Distribution of branch angle frequency of the 684 F<sub>2</sub> individuals from the Col-0  $\times$  29-8 cross. Average and range of branch angles of two natural accessions are indicated. (C) SHOREmap analysis of branch angle. The homozygosity estimator is 0 at even allele frequencies for both natural accessions, 1 when homozygous for the small angle accession Col-0, and -1 when homozygous for the large angle accession 29-8. (D) Fine mapping by genetic-linkage analysis and sequence variation in the candidate gene LAZY1. The causal locus was narrowed down to the region between markers 4.45 M and 4.69 M. The number of recombinants between the markers and the causal locus is indicated on the bottom of the linkage map. Schematic illustration of sequence variation in LAZY1 is shown for Col-0, 3-2, and 29-8. LAZY1 is shown at the top, and the position of ATG is defined as +1. Thin lines indicate introns; thick yellow lines indicate protein-coding sequences; the frameshift induced by the 20 bp deletion is indicated in red.

A



C



B

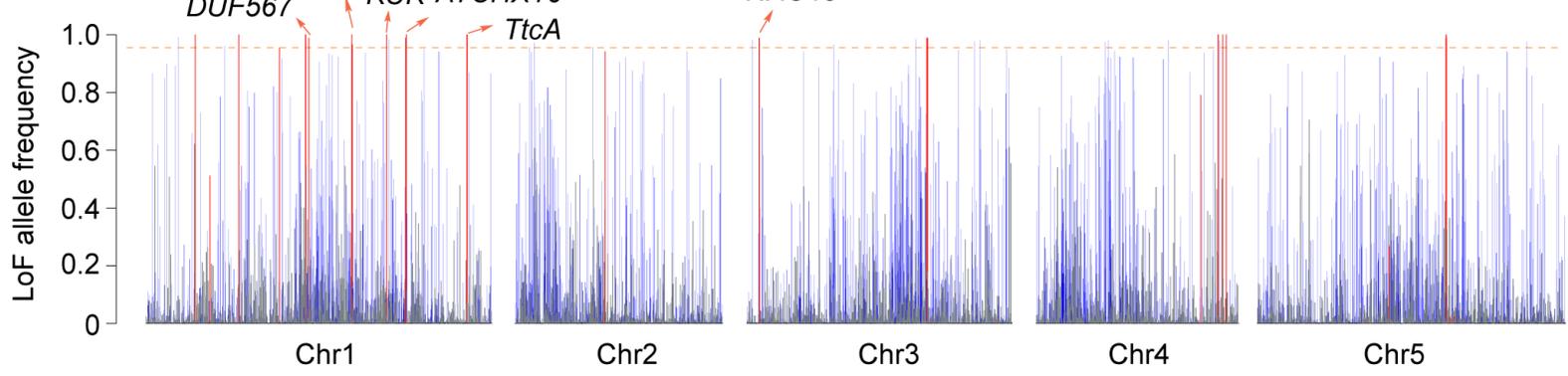


Figure 5. LoF variants are under natural selection. (A) Distribution of allele frequencies of all SNP variants and LoF variants, LoF variants (GWAS LoF variants) and permuted non-LoF variants (GWAS non-LoF variants) associated with environmental variables, at the genome level across 1,071 accessions. (B) Genes under positive selection in the Yangtze River population based on a selective sweep analysis. Gray lines indicate allele frequencies of 60,819 LoF variants across 1,071 accessions; blue lines indicate genes with 3,349 LoF variants in the Yangtze River basin population; red lines indicate genes with LoF variants in regions with positive selection in the Yangtze River basin population. The horizontal dashed line indicates an allele frequency of 95%. (C) Ecological niche modeling of *KUK*. The niche overlap between the LoF and non-LoF haplotype groups assessed by Warren's I similarity statistics ( $I_o = 0.877$ ,  $p < 0.001$ ); areas of suitable habitats are marked in different colors corresponding to a scale from 0 to 1, based on 100 pseudo-replicates.

**Adaptation and phenotypic diversification through loss-of-function mutations in Arabidopsis protein-coding genes**

Yong-Chao Xu, Xiao-Min Niu, Xin-Xin Li, Wenrong He, Jia-Fu Chen, Yu-Pan Zou, Qiong Wu, Yong Edward Zhang, Wolfgang Busch and Ya-Long Guo  
*Plant Cell*; originally published online March 18, 2019;  
DOI 10.1105/tpc.18.00791

This information is current as of March 18, 2019

<b>Supplemental Data</b>	<a href="/content/suppl/2019/03/18/tpc.18.00791.DC1.html">/content/suppl/2019/03/18/tpc.18.00791.DC1.html</a> <a href="/content/suppl/2019/03/18/tpc.18.00791.DC2.html">/content/suppl/2019/03/18/tpc.18.00791.DC2.html</a>
<b>Permissions</b>	<a href="https://www.copyright.com/ccc/openurl.do?sid=pd_hw1532298X&amp;issn=1532298X&amp;WT.mc_id=pd_hw1532298X">https://www.copyright.com/ccc/openurl.do?sid=pd_hw1532298X&amp;issn=1532298X&amp;WT.mc_id=pd_hw1532298X</a>
<b>eTOCs</b>	Sign up for eTOCs at: <a href="http://www.plantcell.org/cgi/alerts/ctmain">http://www.plantcell.org/cgi/alerts/ctmain</a>
<b>CiteTrack Alerts</b>	Sign up for CiteTrack Alerts at: <a href="http://www.plantcell.org/cgi/alerts/ctmain">http://www.plantcell.org/cgi/alerts/ctmain</a>
<b>Subscription Information</b>	Subscription Information for <i>The Plant Cell</i> and <i>Plant Physiology</i> is available at: <a href="http://www.aspb.org/publications/subscriptions.cfm">http://www.aspb.org/publications/subscriptions.cfm</a>