Adaptation and Phenotypic Diversification in Arabidopsis through Loss-of-Function Mutations in Protein-Coding Genes

Yong-Chao Xu,^{a,b} Xiao-Min Niu,^{a,b} Xin-Xin Li,^{a,b} Wenrong He,^c Jia-Fu Chen,^{a,b} Yu-Pan Zou,^{a,b} Qiong Wu,^a Yong E. Zhang,^{b,d} Wolfgang Busch,^c and Ya-Long Guo^{a,b,1}

^a State Key Laboratory of Systematic and Evolutionary Botany, Institute of Botany, Chinese Academy of Sciences, Beijing 100093, China

^b University of Chinese Academy of Sciences, Beijing 100049, China

^o Salk Institute for Biological Studies, Plant Molecular and Cellular Biology Laboratory, La Jolla, California 92037

^d 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

ORCID IDs: 0000-0002-1378-4049 (Y.-C.X.); 0000-0001-6568-5336 (X.-M.N.); 0000-0003-0967-6413 (X.-X.L.); 0000-0003-1408-5327 (W.H.); 0000-0003-0090-1979 (J.-F.C.); 0000-0003-3033-0273 (Y.-P.Z.); 0000-0001-6519-3526 (Q.W.); 0000-0003-3770-2383 (Y.E.Z.); 0000-0003-2042-7290 (W.B.); 0000-0002-4643-4889 (Y.-L.G.)

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 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 1071 Arabidopsis (*Arabidopsis thaliana*) genomes and found that 34% of Arabidopsis protein-coding genes annotated in the Columbia-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.

INTRODUCTION

Variation in gene copy number plays important roles in adaptation and diversification. The two major processes that result in the gain of genes or gene copies are de novo gene formation and gene duplication (Chen et al., 2010; Carvunis et al., 2012; Guo, 2013; Palmieri et al., 2014; Zhao et al., 2014; McLysaght and Guerzoni, 2015; Li et al., 2016). However, the less-is-more hypothesis proposes that gene loss also contributes to evolutionary change in the context of gene copy number variation (Olson, 1999). The dosage balance hypothesis suggests that altering the stoichiometry of members of macromolecular complexes can affect the function of the whole complex and could ultimately affect phenotype and evolutionary fitness (Edger and Pires, 2009; Birchler and Veitia, 2012; Hao et al., 2018). Therefore, the dosage balance hypothesis provides a more compelling explanation of the advantages and disadvantages of removing gene copies from a genome and explains the less-is-more hypothesis in

^[OPEN]Articles can be viewed without a subscription. www.plantcell.org/cgi/doi/10.1105/tpc.18.00791 a mechanistic way. There are two major mechanisms underlying gene loss: the physical removal of a gene by recombination and gene inactivation by loss-of-function (LoF) mutations (Albalat and Cañestro, 2016). The latter mechanism can occur through the gain of a premature stop codon, splice site disruption, or disruption of a transcript reading frame. LoF mutations are frequent and have recently gained attention owing to the improved detection power of advanced sequencing techniques (MacArthur and 1000 Genomes Project Consortium et al., 2012; Narasimhan et al., 2016).

Consistent with the less-is-more hypothesis, many cases of gene loss have been reported in diverse organisms, including unicellular organisms such as bacteria (Will et al., 2010; Hottes et al., 2013), and multicellular organisms, such as plants (Zufall and Rausher, 2004; Hoballah et al., 2007; Song et al., 2007; Tang et al., 2007; Gujas et al., 2012; Amrad et al., 2016; Sas et al., 2016; Wu et al., 2017) and animals (Greenberg et al., 2003; Hodgson et al., 2014; Goldman-Huertas et al., 2015), including humans (MacArthur and 1000 Genomes Project Consortium et al., 2012; Narasimhan et al., 2016). LoF mutations have been shown to affect important biological processes, such as development and stress resistance in plants (Gujas et al., 2012; Wu et al., 2017) and intellectual disabilities in humans (Green et al., 2017). However, the extent to which LoF mutations are correlated with adaptation and phenotypic diversification is largely unknown (Albalat and Cañestro, 2016). Even though LoF mutations that produce large

¹Address correspondence to yalong.guo@ibcas.ac.cn.

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) is: Ya-Long Guo (yalong. guo@ibcas.ac.cn).

IN A NUTSHELL

Background: The gain of genes is thought to play important roles in the adaptation and diversification of plants. Loss-of-function mutations lead to pseudogenization or gene loss, which might contribute to adaptation and diversification as well, but this topic has been largely neglected. The "less-is-more" hypothesis proposes that gene loss functions as an engine for evolutionary change. Most loss-of-function mutations are neutral or deleterious, but the extent to which loss-of-function mutations contribute to adaptation and phenotypic diversification is largely unknown.

Question: Taking advantage of more than 1,000 resequenced *Arabidopsis thaliana* genomes, we investigated the effect of natural loss-of-function mutations on the adaptation and diversification of this model plant.

Findings: We systematically identified 60,819 loss-of-function mutations across 1,071 Arabidopsis genomes worldwide. Intriguingly, 34% of *A. thaliana* protein-coding genes do not have any loss-of-function mutations, implying that loss-of-function mutations in these genes would reduce plant fitness in natural environments. Furthermore, in the Yangtze River basin population, 1% of genes with loss-of-function mutations were under positive selection, suggesting these loss-of-function mutations are closely related to adaptation.

Next steps: The mechanisms by which loss-of-function mutations produce functional effects are complicated. We plan to perform in-depth studies investigating how loss-of-function mutations affect plant adaptation.

and detrimental effects on fitness will be depleted in inbreeding species compared with outcrossing plant and animal species, inbreeders are an excellent system to study the effect of natural homozygous knockouts, as inbreeding leads to a high frequency of homozygous LoF mutations (Saleheen et al., 2017). The best-characterized selfing species is Arabidopsis (*Arabidopsis thaliana*), which is naturally occurring in almost all parts of the world.

In this study, to investigate the evolutionary patterns of LoF mutations in natural populations, we explored 1071 genomes of Arabidopsis from the 1001 Genomes Project (2016), the African genomes project (Durvasula et al., 2017), and our own sequencing project (Zou et al., 2017). Across these 1071 Arabidopsis genomes, we identified 60,819 LoF variants, including 17,453 stop codon-introducing (stop-gain) variants; 37,935 disruptions of a transcript reading frame (frameshift); and 5431 splice sitedisrupting single nucleotide variants ([SNVs]; splice site). We found that the presence of LoF mutations is correlated with the level of nucleotide diversity, the density of transposable elements (TEs), and gene family size. Intriguingly, 1% of genes with LoF mutations are under positive selection in the Yangtze River basin population, suggesting that these genes with LoF mutations are crucial for adaptation. Overall, this study highlights the importance of LoF mutations for adaptation and phenotypic diversification.

RESULTS

The Presence of LoF Variants Is Correlated with Nucleotide Diversity, TE Density, and Gene Family Size

We scanned the genomes of 1071 accessions for LoF variants (Figures 1A and 1B), including 893 genomes downloaded from the 1001 Genomes Project (Supplemental Data Set 1; 1001 Genomes Consortium, 2016), 61 from the African genomes project (Supplemental Data Set 2; Durvasula et al., 2017), and 117 from our own genome project (Supplemental Data Set 3; Zou et al., 2017). Given that LoF mutations annotated by mapping sequences to a reference genome will have high false positive rates,

a series of stringent filtering steps were implemented that were similar to those used in previous studies of human genomes (MacArthur and 1000 Genomes Project Consortium et al., 2012; Narasimhan et al., 2016). First, SNVs and indels (insertions/deletions) were removed if they matched multiple times to the reference genome or were located in short tandem repeats or if the SNVs were in close proximity (3 bp or less) to an indel. Second, stop-gain variants were excluded if they were within a codon that was linked to other SNVs, and frameshift variants were removed if they were occurring together with other frameshifts that resulted in the restoration of the reading frame. Third, SNVs and indels were removed if the inferred LoF mutations occurred within the last 5% of the transcript or were observed in the ancestral states inferred from Arabidopsis Ivrata and Capsella rubella (see "Methods"). Fourth, LoF variants that affect all known transcripts of the proteincoding gene in the reference genome (Columbia-0 [Col-0]) were retained for subsequent analysis. After filtering according to these rules, we detected 17,453 (214.7 per accession) premature stop codons that were caused by SNVs (stop-gain), 37,935 (512.2 per accession) insertion/deletion-induced frameshift variants (frameshift) leading to the disruption of a transcript reading frame, and 5431 (103.0) splice site-disrupting SNVs (splice site; Table 1). Taking all of these stop-gain, frameshift, and splice site variants together, we identified 60,819 (829.9 per accession) LoF variants across the 1071 accessions (Table 1).

To assess the power of our detection method and to estimate the false discovery rate (FDR) for LoF variants, we used RNA sequencing (RNA-seq) data for eight de novo-assembled genomes (Gan et al., 2011), the Landsberg erecta genome assembled with PacBio Sequel data (Zapata et al., 2016), and the KBS-Mac-74 genome assembled with Nanopore data (Michael et al., 2018). The results suggest that our method has a very low FDR (2.6% in each accession on average; Supplemental Data Set 4). Given that the accumulated FDR across 1071 accessions might be higher than that in a single genome, we calculated the accumulated false positive rates based on validation data in these 10 accessions. While with the addition of genomes to the analysis the accumulated false positive LoF mutations as well as the total



Figure 1. Identification of LoF Mutations in 1071 Arabidopsis Genomes.

(A) Geographic distribution of the 1071 accessions used in this study.

(B) Process used to identify LoF variants. SnpEff annotation indicates SNVs/indels annotated by SnpEff software (see "Methods"). (C) Variation in TE density (TE), nucleotide diversity (π), and the density of genes with stop-gain (Stop-gain), frameshift (Frameshift), splice site (Splice), and LoF variants across chromosome 1. The values for the four other chromosomes are shown in Supplemental Figure 2. Chr1, chromosome 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; guanine-cytosine (GC) content in Col-0; TE density in Col-0; nucleotide diversity (π); 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).

number of LoF mutations increased, the accumulated FDR was largely stable (Supplemental Figure 1). Therefore, we estimated that in each Arabidopsis accession, an average of 808.2 genes have LoF mutations (3.0% of all protein-coding genes in the reference Col-0) with an FDR of 2.6%. In particular, 80 false

positive LoF mutations (22 stop-gains, 56 frameshifts, and 2 splice sites) identified in the validation process were excluded from the subsequent analysis.

The density of genes with stop-gain, frameshift, splice site, and LoF variants (summed up for all the stop-gain, frameshift, and

Variant type	Before filtering ^a		After filtering ^a	
	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)

splice site variants), as calculated by dividing the number of genes with LoF variants (genes with LoF variants detected in any of the 1071 accessions) by the gene number in each 200-kb window in the Col-0 reference genome, varied across the genome (Figure 1C; Supplemental Figure 2). By contrast, the density of stop-gain and frameshift variants within gene bodies was roughly evenly distributed, with a slightly higher frequency at the beginning and at the end of genes (Supplemental Figure 3). This pattern is similar to the distribution in human and great ape genomes (MacArthur and 1000 Genomes Project Consortium et al., 2012; de Valles-Ibáñez et al., 2016).

To infer mechanisms that favor LoF, we analyzed the correlations between LoF and various gene features. The density of genes with LoF mutations (stop-gain, frameshift, splice site, and LoF variants) across the Col-0 genome showed a strong correlation with nucleotide diversity (π) across the 1071 genomes and the TE density in the Col-0 genome (Figure 1D). Large gene families tend to have high gene redundancy (Wagner, 2005), and gene loss in large gene families may be an advantage according to the dosage balance hypothesis (Edger and Pires, 2009; Birchler and Veitia, 2012; Hao et al., 2018). We therefore evaluated whether these genes were prone to acquiring LoF mutations. We classified the 27,206 protein-coding genes of the reference Col-0 genome into 7430 gene families based on amino acid sequence similarity. We then estimated the correlation between gene family size and the proportion of genes with common LoF variants (variants displaying a minor allele frequency [MAF] larger than 5% across the 1071 accessions). For all gene families, there was a positive correlation between gene family size and the LoF ratio (the percentage of genes with LoF variants in each gene family; Spearman r = 0.21, P = 2.2 × 10^{-16}). This suggests that genes belonging to larger gene families are more likely to have LoF variants. Similarly, the median sizes of gene families for genes with and without LoF variants were 11 and 6, respectively (Wilcoxon sum test, P < 2.2×10^{-16} ; Figure 2A). This indicates that genes belonging to larger gene families are more likely to acquire LoF mutations. Furthermore, analogous to previous studies (Prachumwat and Li, 2008; Guo, 2013), we divided gene families into three classes: one gene or singleton (class 1: 3825 families), two genes (class 2: 1353 families), and three genes or more (class 3: 2252 families). The genes in class 3 were more likely to contain LoF variants than the others (class 3 versus class 1 or 2; Wilcoxon sum test, all P < 0.01; Supplemental Figure 4). Therefore, the presence of LoF variants is correlated with nucleotide diversity, TE density, and gene family size.

A Large Fraction of Arabidopsis Genes Could Affect Fitness in Natural Environments

A genome-wide study of LoF variants can provide insight into gene lethality by analyzing the frequency of all null alleles for a given gene in natural environments (population gene lethality; Albalat and Cañestro, 2016). We found that 9249 protein-coding genes (34.0% of all protein-coding genes in the Col-0 reference genome) do not have any LoF variant within our panel of 1071 accessions. To ask whether known essential genes are depleted of LoF variants, we used a database of 358 lethal genes identified in laboratory studies (Meinke et al., 2008); 88.3% of the lethal genes do not have any natural LoF mutation, and 7.3% of the lethal genes display a LoF allele frequency less than 0.5% across the 1071 accessions (Figure 2B). This suggests that lethal genes in Arabidopsis are significantly depleted of natural LoF mutations (Pearson's chi-squared test, $P < 2.2 \times 10^{-16}$). We then performed gene ontology (GO) enrichment analysis of the four different gene classes based on the LoF allele frequency in the 1071 accessions (genes without LoF, with LoF, with a LoF allele frequency of <0.05 and \geq 0.05). Although both the gene sets (with and without LoF variants) were enriched for the GO term 'unknown function', the only gene set highly enriched for many other GO terms was the set without LoF variants (Figure 2C). These results suggest a model in which natural knockouts of 9249 Arabidopsis genes without LoF variants would reduce fitness in natural environments, and individuals with LoF mutations in these nonessential genes can likely be removed from populations by purifying selection.

LoF Variants Are Correlated with Climate Variables

LoF variants could be functionally important and associated with climate variables. To investigate this notion, we performed a genome-wide association study (GWAS) of 14,270 LoF variants with MAF larger than 0.5% and 20 environmental variables, including latitude and 19 climate variables (Supplemental Data Set 5). To validate that the associations were not random, we compared the total number of significantly associated LoF variants with the number obtained from 1000 GWAS permutations based on 14,270 randomly resampled single-nucleotide polymorphisms ([SNPs]; MAF > 0.5%) in the coding regions of protein-coding genes across the whole genome for each of the 20 environmental variables. If LoF variants were randomly associated with environmental variables, we would expect the number of associated variants obtained in the permutation analysis and the observed value for each environmental variable to be similar. For 10 of the 20 environmental variables, the observed number was significantly higher



Figure 2. Functional Analysis 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 Arabidopsis. <0.005 indicates MAF less than 0.5%, <0.05 indicates MAF < 5%, and >0.05 indicates MAF > 5%.

(C) GO enrichment analysis of genes with LoF, with different LoF allele frequencies, and those without LoF variants across the 1071 genomes. Colors indicate 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. NS, insignificant difference.

(E) Expression level variation between alleles with and without LoF variants (6675 genes) based on RNA-seq data for 541 accessions. NS, insignificant difference.

(F) GWAS of stop-gain variants of KUK with regard to precipitation in the wettest month. Chr, chromosome.

than the number identified by permutation analysis (one-sample Wilcoxon signed rank test, P < 0.05; Supplemental Figure 5), suggesting that LoF variants are more prone to associate with climate variables. Also, 125 of the 14,270 LoF variants (0.9%) were significantly associated with at least 1 of the 10 environmental variables, involving 124 genes (Supplemental Figures 6 to 8; Supplemental Data Set 6). Of these genes, some are known to be functionally important, such as *Cysteine-rich receptor-like protein kinase 36* (*CRK36*; Supplemental Data Set 6). *CRK36* encodes an abiotic stress-inducible receptor-like protein kinases that negatively regulates abscisic acid signaling (Tanaka et al., 2012).

We used RNA-seq data for 541 accessions to compare the expression levels of alleles with and without LoF mutations for each of these environmentally associated genes (Figure 2D; Kawakatsu and 1001 Genomes Consortium et al., 2016). Among the 124 associated genes, 49 genes have both LoF and non-LoF alleles in the 541 accessions for which RNA-seq data were available. There were no LoF variants of the other 75 genes in this set of accessions. We therefore focused on the 49 genes with both LoF and non-LoF alleles (different LoF variants of the same gene were combined). The transcript levels of 16.4% of these genes were significantly altered (8.2% upregulated and 8.2% downregulated) in accessions harboring the LoF variants (Wilcoxon sum test, FDR corrected P < 0.05; Figure 2D). Furthermore, genes with LoF variants associated with

environmental variables exhibited greater rates of upregulation (8.2 versus 3.0%) or downregulation (8.2 versus 6.5%) than genes with LoF variants at the genome level (Figure 2E). To determine whether the expression changes were caused by the different genetic backgrounds of various natural accessions, we calculated the upand downregulated gene numbers based on 30,000 randomly selected alleles with non-LoF variants in coding regions across the 541 accessions. A higher proportion of genes with LoF variants associated with environmental variables are up- (8.2 versus 1.7%) or downregulated (8.2 versus 2.4%) compared with those detected when only the genetic background was taken into account (Supplemental Figure 9). The higher proportion of upregulated genes with LoF variants associated with environmental variables could result from diverse mechanisms, such as a dominant gain-offunction effect, in which LoF mutations are effectively gain-offunction mutations, thereby mimicking increased gene function (Schild et al., 1995; Xie et al., 1998).

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

Common LoF alleles are nearly always less deleterious than rare LoF alleles, but some could have an impact on phenotypes (MacArthur and 1000 Genomes Project Consortium et al., 2012).

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

To independently test the hypothesis that truncated versions of KUK are associated with increased meristem and mature cell sizes, we phenotyped accessions containing specific LoF alleles that were not evaluated in the previous study (with one exception: Bur-0; Meijón et al., 2014). We chose three accessions that did not contain a stop codon (control); three accessions with the most common stop codon located in the first half of KUK (111 G/A: TGG-TGA, SG111), the only two accessions in the 1001 Genomes panel that had a distal stop codon (652 C/T: CGA-TGA, SG652), including one of the accessions with the largest meristem and mature cell traits (Bur-0) according to a previous study (Meijón et al., 2014); and two accessions with a distal frameshift (738 T/-, FS738; Supplemental Figure 13; Supplemental Data Set 9). Compared with the control, all three classes showed significant increases (pairwise Wilcoxon rank sum test, P < 0.05) in meristem length, with the most common stop codon associated with the smallest increase (Supplemental Figure 13; Supplemental Data Set 9). Overall, these data suggest that truncated alleles of KUK



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.

influence cellular traits, most likely via a dominant mechanism. Further transgenic studies should be performed to confirm the effects of these genetic variants in the same genetic background.

In addition to the LoF alleles of genes that were associated with environmental variables based on the GWAS, we studied *PSEUDO-RESPONSE REGULATOR5 (PRR5)*, which is involved in regulating various circadian-associated biological events, such as flowering time (Nakamichi et al., 2016). *PRR5* alleles with and without LoF variants are widely distributed in our set of 1071 accessions (Supplemental Figure 12B). Of the 812 accessions with flowering time data (Figure 3C; Supplemental Data Set 10; 1001 Genomes Consortium, 2016), the flowering time of accessions with *PRR5* LoF variants (frameshift, or total LoF; no putative splice variation was observed) was significantly delayed

compared with accessions without LoF variants at either 10 or 16°C (population structure–corrected ANOVA, P < 0.05; Figure 3D; Supplemental File). We note that while flowering time for the stop-gain allele was delayed relative to the non-LoF allele, the difference was not significant, probably owing to the small number of accessions.

Given the above-mentioned two cases of genes having common LoF variants, we further studied *LAZY1* as an example for a rare LoF variant. Only one accession from the Yangtze River basin population (29-8) showed a frameshift mutation in *LAZY1*, which that was caused by a 20-bp deletion in exon 3 and resulted in a premature stop codon (Figure 4D; Supplemental Figure 14). *LAZY1* mutant plants display a much larger branch angle (81°) than the wild type (42°; Yoshihara et al., 2013). As expected, we found



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) shows the distribution of branch angles in different accessions. The borders of the boxes 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 F2 individuals from the Col- 0×29 -8 cross. The average and range of branch angles of two natural accessions are indicated.

(C) SHORE map 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. chr., chromosome.

(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 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.

that accession 29-8 has a large branch angle (74°; Figure 4A). To validate that this frameshift mutation was the causal mutation, we crossed accession 29-8 with accessions displaying smaller branch angles: Col-0 (30°) and 3-2 (30°). In the Col-0 imes29-8 F2 population, the segregation ratio for branch angles was roughly 3:1, suggesting that a major, recessively acting gene was responsible for the large branch angle of 29-8 (Figure 4B). Using a mapping-by-sequencing approach based on 57 F2 plants from the Col-0 \times 29-8 F2 population with large branch angles (>85°), we identified a causal region on chromosome 5 (Figure 4C). To refine the target interval, we used F2 plants in both Col-0 \times 29-8 and 3-2 \times 29-8 populations, and finally narrowed down the causal locus to a 240-kb interval containing LAZY1 (Figure 4D; Supplemental Data Sets 11 and 12). The truncated transcript of LAZY1 in the 29-8 accession caused the loss of the only functional domain of the nuclear localization signal, which is located between conserved regions 3 and 4. This nuclear localization signal is a key component during the nuclear import process (Supplemental Figure 14; Yoshihara et al., 2013). Overall, these results provide evidence that the frameshift variant of LAZY1 is likely responsible for the larger branch angle in accession 29-8.

A Subset of LoF Variants Exhibits a Signature of Natural Selection

Compared with other SNP variants across the whole genome in the 1071 accessions, LoF variants are biased toward low frequencies (Figure 5A). This suggests that many natural LoF variants might be deleterious in Arabidopsis and are under purifying selection. Similarly, compared with 1000 permutation results of non-LoF variants associated with environmental variables (GWAS non-LoF variants), LoF variants associated with environmental variables (GWAS LoF variants) were biased toward low allele frequencies (Figure 5A). This suggests that a large fraction of LoF variants that are associated with environmental variables are likely under purifying selection as well.

Although LoF mutations are usually deleterious, adaptive LoF mutations can occur and spread rapidly in small populations (Olson, 1999). To determine whether genes with LoF mutations were influenced by natural selection, we focused on the Yangtze River basin population (86 accessions; Supplemental Data Set 3), a population that recently adapted to this region (Zou, 2017). In total, 2709 genes had LoF variants (3349 LoF variants) in this population; among them, 54 genes (64 LoF variants) were under





(A) Distribution of allele frequencies of all SNP variants and LoF variants, LoF variants (GWAS LoF variants), and permutated non-LoF variants (GWAS non-LoF variants) associated with environmental variables at the genome level across 1071 accessions.

(B) Genes under positive selection in the Yangtze River basin population based on selective sweep analysis. Gray lines indicate allele frequencies of 60,819 LoF variants across 1071 accessions; blue lines indicate genes with 3349 LoF variants in the Yangtze River basin population; and 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%. Chr, chromosome. Leucine-rich repeat like protein (domain of unknown function, *DUF567*), *Increased Salt Tolerance 12 (ISTL12), cation/H+exchanger 16* (*CHX16*), *two Thiolation of Cytidine tRNA biosynthesis protein (TtcA), No Apical Meristem 45* (*NAC45*).

(C) Ecological niche modeling of KUK. The niche overlap between the LoF and non-LoF haplotype groups was assessed by Warren's *I* similarity statistics (lo = 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.

positive selection based on a selection sweep analysis in our recent study (Figure 5B; for details, see Supplemental Data Sets 13 and 14; Zou et al., 2017). Of the 54 genes in these selected regions, 27 genes (1.0% of the genes with LoF variants) had a LoF allele frequency exceeding 95% and were more likely under positive selection, including *KUK*, in which one LoF allele (111 G/A: TGG-TGA type; Supplemental Figure 11) was fixed in the Yangtze River basin population (Supplemental Data Set 14).

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

DISCUSSION

The less-is-more hypothesis proposes an evolutionary process involving adaptive gene loss (Olson, 1999). LoF mutations have gained increasing attention (MacArthur and 1000 Genomes Project Consortium et al., 2012; Yang et al., 2015; Narasimhan et al., 2016). For example, recent studies have demonstrated the functional importance of natural LoF mutations (Gujas et al., 2012; Lek and Exome Aggregation Consortium et al., 2016; Saleheen et al., 2017; Wu et al., 2017). However, our understanding of the evolutionary pattern of LoF mutations at the genome level and the adaptive effects of LoF variants at the population level is highly limited. In this study, we investigated the evolutionary pattern of LoF mutations in Arabidopsis and found that a high level of nucleotide diversity, high TE density, and high gene redundancy (large gene family size) are associated with LoF mutations. The latter two factors were also observed in a previous study of LoF mutations in human populations (MacArthur and 1000 Genomes Project Consortium et al., 2012). Given that these three factors themselves are largely correlated with rates of evolution, rapid sequence evolution could be associated with a high frequency of LoF mutations. More importantly, we hypothesize that 34% of Arabidopsis genes can affect fitness in natural environments, and furthermore, 1% of genes with LoF variants are under positive selection in the Yangtze River basin population.

The mechanisms by which LoF mutations can produce functional effects are complicated. Some of the profound effects of LoF can be explained by the dosage balance hypothesis (Birchler and Veitia, 2007; Hou et al., 2018; Kremling et al., 2018). A simple mechanism for a beneficial effect of a null mutation is the removal of a protein that is detrimental in the current environment (Hottes et al., 2013). For example, *brevis radix* LoF alleles in Arabidopsis help roots adapt to acidic soil (Gujas et al., 2012). However, it is also possible that LoF variants can act as dominant mutations. For instance, stop-gain variants of the Sex-determining region Y-box transcription factors 9 (MiniSOX9) act in a dominant-negative manner, thereby counteracting the activity of the wild-type variant (Abdel-Samad et al., 2011). Other stop-gain and missense variants can act as dominant gain-of-function mutations, thereby mimicking increased gene function (Schild et al., 1995; Xie et al., 1998). In the case of *KUK*, a gene that acts as a positive regulator of meristem and cell size (Meijón et al., 2014), a hypermorphic gain-of-function is a more likely explanation than a LoF. Although most premature stop codons result in a LoF or dominant effects, more complex possibilities exist. For instance, in the fruitfly (*Drosophila sechellia*), an olfactory receptor pseudogene encodes a functional receptor as a result of translation read-through of the premature termination codon (Prieto-Godino et al., 2016).

Of the many LoF variants (or pseudogenes) that we identified in more than 1000 natural accessions, some are presumably functional, with a minor or altered function, rather than nonfunctional. More efforts are needed to understand their functional effects. For example, transgenic analyses are needed to clarify the functional differences between the predicted pseudogenes and their ancestral genes. While more function research is needed, LoF mutational variation also provides a way of investigating the phenotypic consequences of the loss of specific genes within diverse genomes while establishing their role in the evolutionary process. Overall, our study highlights the importance of natural knockouts for adaptation and phenotypic diversification and emphasizes the need for more in-depth studies of LoF variants.

METHODS

Data Analyzed in This Study

The genomes of 1071 accessions were used, including 893 representative genomes of Arabidopsis (*Arabidopsis thaliana*) from the 1001 Genomes Project (Supplemental Data Set 1; 1001 Genomes Consortium, 2016), 61 from the African sequenced genomes project (Supplemental Data Set 2; Durvasula et al., 2017), and 117 genomes from our own sequencing project (Supplemental Data Set 3; Zou, 2017). The 541 transcriptomes used in this study were obtained from a previous study (Kawakatsu and 1001 Genomes Consortium et al., 2016), as indicated in Supplemental Data Set 1. All 1071 accessions with clean data were mapped against the Col-0 reference genome, and SNPs and indels were called using Genome Analysis Toolkit (GATK v2.1.8) with default parameters (DePristo et al., 2011).

Identification of Candidate LoF Variants

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

Filtering of Candidate LoF Variants

To remove false positive LoF variants caused by sequencing and mapping errors, annotation errors, and the effect of nearby variants, a series of stringent filters were used. First, a sequence context filter was applied. Variants that could be mapped to multiple regions of the reference genome (Col-0) were removed using Genome Multitool (Derrien et al., 2012; Marco-Sola et al., 2012). SNVs (stop-gain and splice site) and indels (frameshift) overlapping with tandem repeats and SNVs in close proximity (3 bp or less)

to a known indel were excluded. Second, multi-nucleotide polymorphism filters were applied. To filter incorrect premature stop codon annotation, stop-gain variants that were found within the same codon linked to other SNVs were removed. Frameshift variants were filtered out if two frameshift variants resulted in the restoration of the reading frame. Three or more frameshift variants occurring in each gene of the same accession were excluded; most resulted in the restoration of the reading frame. Third, annotation filters were applied. SNVs and indels were filtered out if the inferred LoF allele was also the ancestral state, as this implies that such variants were gain-of-function or the sequencing errors of the gene model at this location in the reference.

Two related species, *Arabidopsis lyrata* (MN47) and *Capsella rubella* (MTE), were used to infer the ancestral state of SNVs and indels. For SNVs and indels, ancestral states were determined using alignments of Col-0 with the *A. lyrata* and *C. rubella* reference genomes, as described in our previous study (Li et al., 2016). SNVs and indels were filtered out if either of the two outgroups showed the same stop-gain or frameshift mutations observed in the 1071 accessions.

The disrupted fraction of the coding sequence of the longest transcript caused by stop-gain and frameshift variants was calculated, and stop-gain or frameshift mutations occurring within the last 5% of the transcript were removed. Because we did not determine the position of splice site disruption on the final transcript (MacArthur and 1000 Genomes Project Consortium et al., 2012), we were not able to perform this analysis for splice site variation. Finally, if there were two or more LoF variants in the same allele, only the LoF variant closest to the start codon (ATG) was considered.

Validation of LoF Mutations

We estimated the FDR using the assembled transcripts of eight accessions from a previous study (Gan et al., 2011), as well as gene sequences that were annotated using two published long reads assembled genomes: Landsberg erecta assembled from PacBio Sequel data (Zapata et al., 2016) and KBS-Mac-74 assembled from Nanopore data (Michael et al., 2018). The genes of KBS-Mac-74 and Landsberg erecta were annotated using exonerate v2.2.0 (Slater and Birney, 2005) with coding sequences (CDS) of Col-0. All assembled transcript sequences for genes with LoF mutations in the eight accessions, and annotated gene sequences of KBS-Mac-74, Landsberg erecta, and Col-0 gene sequences were aligned using MUSCLE v3.8.31 (Edgar, 2004). All LoF mutations were manually checked based on the aligned sequences. Stop-gain and frameshift variants of the 10 accessions were used to assess the curve of accumulated false positive rate (all false positive LoF variants found in each accession combined). Splice site variants were not used in the analysis of accumulated false positive rate, because assembled transcripts of eight accessions were assembled based on short reads, which are not robust enough to validate splice sites. Nevertheless, the FDRs for splice sites in the two long reads assembled genomes are similar or even lower than either stop-gain or frameshift. The accumulated false positive LoF number, accumulated total LoF number (the number of LoF variants with transcriptional or long reads data in each accession), and accumulated false positive LoF ratio were calculated with all probable accumulated groups.

LoF Polymorphism Variation

All parameters were calculated using a 200-kb window size and 10-kb step size along the Col-0 reference genome. The density of genes with stop-gain variants was calculated by dividing the number of genes with stop-gains by the number of genes in each 200-kb window in the Col-0 reference genome. The same rationale was used to calculate the density of frameshift, splice site, and total LoF variants (summed over all the three LoF variants). The π was calculated based on the SNP matrix of 1071 accessions in each 200-kb window. The guanine-cytosine content was calculated based on

the Col-0 genome in each 200-kb window. Exon number, CDS length, and gene expression level were calculated for each gene based on the Col-0 genome in each 200-kb window. Gene expression levels in Col-0 were obtained from a previous study (Kawakatsu and 1001 Genomes Consortium et al., 2016). TE density was calculated as the number of TEs in each 200-kb window in the Col-0 reference genome with TE annotations downloaded from The Arabidopsis Information Resource (TAIR10).

Gene Family, Sequence, and Expression Analysis

Gene family analysis was performed according to our previous study (Guo, 2013). GO term enrichment analysis was performed using agriGO (Tian et al., 2017). All genes (5050) with candidate LoF variants were removed from the non-LoF category in the GO analysis. PCR was performed to obtain the *LAZY1* upstream region (\sim 1.5 kb), gene body, and downstream regions (\sim 0.5 kb) from Col-0, 3-2, and 29-8 (marked with an asterisk in Supplemental Data Set 3) using Q5 polymerase (New England Biolabs). Sequences were aligned using Lasergene Seqman (DNASTAR). Markers used in the analysis of plants with large branch angles are listed in Supplemental Data Sets 11 and 12. Primers used for the sequencing and amplification are listed in Supplemental Data Set 15.

In total, 541 transcriptomes were available from the 1001 Genomes Project (Supplemental Data Set 1; Kawakatsu and 1001 Genomes Consortium et al., 2016), and 6675 genes with LoF variants whose mean expression levels (fragments per kilobase of exon per million fragments mapped) were larger than 3 across 541 accessions were kept for subsequent analysis (Meng et al., 2016). For each of the 6675 genes with LoF variants, 541 alleles were divided into two groups: alleles with and without LoF variants. The Wilcoxon sum test was used to evaluate the expression data. All P-values were adjusted for multiple testing by computing their FDR (Benjamini and Hochberg, 1995). Up- or downregulation was calculated using the mean expression data in the two groups (alleles with and without LoF variants).

Genome-Wide Association Study

GWAS was performed by using the compressed mixed linear model (Zhang et al., 2010) from the GAPIT R package (Lipka et al., 2012) to associate data for 20 environmental variables (latitude and 19 climate variables from the WORLDCLIM database; Supplemental Data Set 5) with the 14,270 LoF variants (MAF > 0.5%) of the 1071 accessions. Principal components (Q matrix) and a kinship matrix (K matrix) were included to account for population structure. Bayesian information criterion-based model selection as implemented in the GAPIT R package was used to find the optimal number of principal components for each trait. The whole-genome significance cutoff for associations was set to 0.01/total LoF variants [-log10 (p) = 6.12], p indicates p-value estimated by F-test for the association between each SNP and phenotype value. For the GWAS permutation analysis, the SNP matrix that contained 14,270 SNPs (MAF > 0.5%) was randomly resampled 1000 times for those SNPs in the coding region of protein-coding genes across the whole genome. Whether the actual value is larger than the permutation results was tested using the one-sample Wilcoxon signed rank test. GWAS analysis was also performed using 3822 stop-gain variants (MAF >0.5%) with these 20 environmental variables within the 1071 accessions. The significance cutoff in this analysis was set to 0.05/total LoF variants $[-\log_{10} (p) = 4.88]$.

Ecological Niche Modeling Analysis

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

ENMTools 1.3 (Warren et al., 2010) was used to perform niche identity test by calculating Warren's *I*, which ranged from 0 (no niche overlap) to 1 (identical niches), with 100 pseudo-replicates, between the LoF and non-LoF *KUK* groups using the 10 environmental variables that were used in MaxEnt 3.3.3. The *Io* was calculated using the ecological niche overlap function implemented in ENMTools 1.3.

Measurements of Root Traits

Measurements were performed using 3.5-d-old seedlings grown on 1 imesMurashige and Skoog medium with 1% Suc under a 16-h-light/8-h-dark cycle at 22°C (120 μ mol m⁻² s⁻¹ light intensity). The photon flux density was 120 μ mol m⁻² s⁻¹ (five cold white fluorescent Philips T5 28W/840 light bulbs and one warm white and yellow Philips T5 28W/830 light bulb). For root length measurement, roots were imaged using charge-coupled device scanners and measured using Fiji software (Schindelin et al., 2012). More than 15 seedlings were measured per line. To measure cellular traits, each seedling was treated with 10 $\mu g/mL$ propidium iodide solution for 2 min and washed with water twice. Each root tip was subsequently imaged under a Zeiss 710 confocal microscope with a 20× objective. Both meristematic zone and mature cell sizes were measured using Fiji software (Schindelin et al., 2012). The length of the meristematic zone of each seedling was determined by measuring the distance from the quiescent center to the root cortex cell below the first cortex cell that was twice as long as the cell below. The values for the left and right side of root longitudinal sections were averaged to obtain the length of the meristematic zone. The length of mature cells of each seedling was determined by averaging cell lengths of three mature cortex cells in the zone in which the xylem strands became visible. Five seedlings were measured per accession.

Mapping-by-Sequencing of Branch Angle

Plants of three natural accessions (Col-0, 3-2, and 29-8), and F2 populations, were grown under long day conditions (16-h-light/8-h-dark, $120 \,\mu$ mol m⁻²s⁻¹ light intensity) at 20°C. When the lateral shoot was longer than 5 cm, a protractor was used to measure the angle between the lateral branch and the main stem. The branch angles were measured in 2016 and 2017.

DNA was extracted from pooled leaves of 57 plants with a large branch angle (>85°) selected from a group of 684 F2 plants of Col-0 \times 29-8 with a large branch angle (>85°). The pooled DNA sample was sequenced using Illumina HiSeq X Ten system (150-bp pair-end reads, insert size of 450 bp). In total, 41,757,716 total reads (\sim 34.5-fold coverage) were mapped to the Col-0 genome (TAIR10), and SNPs were called using SHORE. The SNPs from pooled populations were used as markers to identify regions with an excess of homozygous alleles across the genome using SHOREmap (Schneeberger et al., 2009).

Statistical Analysis

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

Accession Numbers

Sequence data from this article can be found in the GenBank/EMBL libraries under the following accession numbers: the *LAZY1* (AT5G14090) sequences of Col-0, 3-2, and 29-8 under accession numbers MF621897 to MF621899. Other genes studied in this article can be found in The Arabidopsis Information Resource database (https://www.arabidopsis.org) under the following accession numbers: *CRK36* (AT4G04490), *KUK* (AT1G60370), and *PRR5* (AT5G24470). The genome data for the pooled Col-0 \times 29-8 F2 populations have been deposited in the National Center for Biotechnology Information Sequence Read Archive under accession number SRR5947598.

Supplemental Data

Supplemental Figure 1. Accumulated false positive LoF rate.

Supplemental Figure 2. Variation in TE density (TE), nucleotide diversity (π), and the density of genes with stop-gain (Stop-gain), frameshift (Frameshift), splice site (Splice), and LoF variants across the chromosomes.

Supplemental Figure 3. Distribution of stop-gain and frameshift variants at gene positions in all 1071 accessions, distributed in 5% bins.

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

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

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

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

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

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

Supplemental Figure 10. Genome-wide association study of stopgain variants and 10 climate parameters.

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

Supplemental Figure 12. Geographic distribution of LoF variants for two genes in 1071 accessions.

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

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

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

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

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

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

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

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

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

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

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

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

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

Supplemental Data Set 11. Markers used in the analysis of plants with large branch angles of the F_2 generation of Col-0 \times 29-8.

Supplemental Data Set 12. Markers used in the analysis of plants with large branch angles of the F_2 generation of 3-2 \times 29-8.

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

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

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

Supplemental File. ANOVA/test tables.

ACKNOWLEDGMENTS

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

AUTHOR CONTRIBUTIONS

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

Received October 18, 2018; revised February 25, 2019; accepted March 17, 2019; published March 18, 2019.

REFERENCES

1001 Genomes Consortium (2016). 1,135 Genomes reveal the global pattern of polymorphism in *Arabidopsis thaliana*. Cell **166**: 481–491.

- Abdel-Samad, R., et al. (2011) MiniSOX9, a dominant-negative variant in colon cancer cells. Oncogene 30: 2493–2503.
- Albalat, R., and Cañestro, C. (2016). Evolution by gene loss. Nat. Rev. Genet. 17: 379–391.
- Amrad, A., Moser, M., Mandel, T., de Vries, M., Schuurink, R.C., Freitas, L., and Kuhlemeier, C. (2016). Gain and loss of floral scent production through changes in structural genes during pollinatormediated speciation. Curr. Biol. 26: 3303–3312.
- Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. J. R. Stat. Soc. B 57: 289–300.
- **Birchler, J.A., and Veitia, R.A.** (2007). The gene balance hypothesis: From classical genetics to modern genomics. Plant Cell **19:** 395–402.
- Birchler, J.A., and Veitia, R.A. (2012). Gene balance hypothesis: Connecting issues of dosage sensitivity across biological disciplines. Proc. Natl. Acad. Sci. USA 109: 14746–14753.
- Carvunis, A.R., et al. (2012) Proto-genes and de novo gene birth. Nature 487: 370–374.
- Chen, S., Zhang, Y.E., and Long, M. (2010). New genes in *Drosophila* quickly become essential. Science **330**: 1682–1685.
- DePristo, M.A., Banks, E., Poplin, R., Garimella, K.V., Maguire, J.R., Hartl, C., Philippakis, A.A., del Angel, G., Rivas, M.A., Hanna, M., McKenna, A., and Fennell, T.J., et al. (2011). A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nat. Genet. 43: 491–498.
- Derrien, T., Estellé, J., Marco Sola, S., Knowles, D.G., Raineri, E., Guigó, R., and Ribeca, P. (2012). Fast computation and applications of genome mappability. PLoS One 7: e30377.
- de Valles-Ibáñez, G., Hernandez-Rodriguez, J., Prado-Martinez, J., Luisi, P., Marquès-Bonet, T., and Casals, F. (2016). Genetic load of loss-of-function polymorphic variants in great apes. Genome Biol. Evol. 8: 871–877.
- Durvasula, A., Fulgione, A., Gutaker, R.M., Alacakaptan, S.I., Flood, P.J., Neto, C., Tsuchimatsu, T., Burbano, H.A., Picó, F.X., Alonso-Blanco, C., and Hancock, A.M. (2017). African genomes illuminate the early history and transition to selfing in *Arabidopsis thaliana*. Proc. Natl. Acad. Sci. USA **114**: 5213–5218.
- Edgar, R.C. (2004). MUSCLE: A multiple sequence alignment method with reduced time and space complexity. BMC Bioinformatics 5: 113.
- Edger, P.P., and Pires, J.C. (2009). Gene and genome duplications: The impact of dosage-sensitivity on the fate of nuclear genes. Chromosome Res. 17: 699–717.
- Gan, X., et al. (2011) Multiple reference genomes and transcriptomes for *Arabidopsis thaliana*. Nature **477**: 419–423.
- Goldman-Huertas, B., Mitchell, R.F., Lapoint, R.T., Faucher, C.P., Hildebrand, J.G., and Whiteman, N.K. (2015). Evolution of herbivory in Drosophilidae linked to loss of behaviors, antennal responses, odorant receptors, and ancestral diet. Proc. Natl. Acad. Sci. USA 112: 3026–3031.
- Green, C., Willoughby, J., Study, D., and Balasubramanian, M. (2017). De novo *SETD5* loss-of-function variant as a cause for intellectual disability in a 10-year old boy with an aberrant blind ending bronchus. Am. J. Med. Genet. A. **173**: 3165–3171.
- Greenberg, A.J., Moran, J.R., Coyne, J.A., and Wu, C.I. (2003). Ecological adaptation during incipient speciation revealed by precise gene replacement. Science **302**: 1754–1757.
- Gujas, B., Alonso-Blanco, C., and Hardtke, C.S. (2012). Natural Arabidopsis *brx* loss-of-function alleles confer root adaptation to acidic soil. Curr. Biol. 22: 1962–1968.
- **Guo, Y.L.** (2013). Gene family evolution in green plants with emphasis on the origination and evolution of *Arabidopsis thaliana* genes. Plant J. **73**: 941–951.

- Hao, Y., Washburn, J.D., Rosenthal, J., Nielsen, B., Lyons, E., Edger, P.P., Pires, J.C., and Conant, G.C. (2018). Patterns of population variation in two paleopolyploid eudicot lineages suggest that dosage-based selection on homeologs is long-lived. Genome Biol. Evol. 10: 999–1011.
- Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G., and Jarvis,
 A. (2005). Very high resolution interpolated climate surfaces for global land areas. Int. J. Climatol. 25: 1965–1978.
- Hoballah, M.E., Gübitz, T., Stuurman, J., Broger, L., Barone, M., Mandel, T., Dell'Olivo, A., Arnold, M., and Kuhlemeier, C. (2007). Single gene-mediated shift in pollinator attraction in *Petunia*. Plant Cell **19**: 779–790.
- Hodgson, J.A., Pickrell, J.K., Pearson, L.N., Quillen, E.E., Prista, A., Rocha, J., Soodyall, H., Shriver, M.D., and Perry, G.H. (2014). Natural selection for the Duffy-null allele in the recently admixed people of Madagascar. Proc. Biol. Sci. 281: 20140930.
- Hottes, A.K., Freddolino, P.L., Khare, A., Donnell, Z.N., Liu, J.C., and Tavazoie, S. (2013). Bacterial adaptation through loss of function. PLoS Genet. 9: e1003617.
- Hou, J., et al. (2018) Global impacts of chromosomal imbalance on gene expression in *Arabidopsis* and other taxa. Proc. Natl. Acad. Sci. USA 115: E11321–E11330.
- Kang, H.M., Sul, J.H., Service, S.K., Zaitlen, N.A., Kong, S.Y., Freimer, N.B., Sabatti, C., and Eskin, E. (2010). Variance component model to account for sample structure in genome-wide association studies. Nat. Genet. 42: 348–354.
- Kawakatsu, T., et al.; 1001 Genomes Consortium (2016). Epigenomic diversity in a global collection of *Arabidopsis thaliana* accessions. Cell 166: 492–505.
- Kremling, K.A.G., Chen, S.Y., Su, M.H., Lepak, N.K., Romay, M.C., Swarts, K.L., Lu, F., Lorant, A., Bradbury, P.J., and Buckler, E.S. (2018). Dysregulation of expression correlates with rare-allele burden and fitness loss in maize. Nature 555: 520–523.
- Lek, M., et al.; Exome Aggregation Consortium (2016). Analysis of protein-coding genetic variation in 60,706 humans. Nature 536: 285–291.
- Li, P., Filiault, D., Box, M.S., Kerdaffrec, E., van Oosterhout, C., Wilczek, A.M., Schmitt, J., McMullan, M., Bergelson, J., Nordborg, M., and Dean, C. (2014). Multiple *FLC* haplotypes defined by independent cis-regulatory variation underpin life history diversity in *Arabidopsis thaliana*. Genes Dev. 28: 1635–1640.
- Li, Z.W., Chen, X., Wu, Q., Hagmann, J., Han, T.S., Zou, Y.P., Ge, S., and Guo, Y.L. (2016). On the origin of de novo genes in *Arabidopsis thaliana* populations. Genome Biol. Evol. 8: 2190–2202.
- Lipka, A.E., Tian, F., Wang, Q., Peiffer, J., Li, M., Bradbury, P.J., Gore, M.A., Buckler, E.S., and Zhang, Z. (2012). GAPIT: Genome association and prediction integrated tool. Bioinformatics 28: 2397–2399.
- MacArthur, D.G., et al.; 1000 Genomes Project Consortium (2012). A systematic survey of loss-of-function variants in human proteincoding genes. Science 335: 823–828.
- Marco-Sola, S., Sammeth, M., Guigó, R., and Ribeca, P. (2012). The GEM mapper: Fast, accurate and versatile alignment by filtration. Nat. Methods 9: 1185–1188.
- McLaren, W., Pritchard, B., Rios, D., Chen, Y., Flicek, P., and Cunningham, F. (2010). Deriving the consequences of genomic variants with the Ensembl API and SNP Effect Predictor. Bioinformatics 26: 2069–2070.
- McLysaght, A., and Guerzoni, D. (2015). New genes from non-coding sequence: The role of de novo protein-coding genes in eukaryotic evolutionary innovation. Philos. Trans. R. Soc. Lond. B Biol. Sci. 370: 20140332.
- Meijón, M., Satbhai, S.B., Tsuchimatsu, T., and Busch, W. (2014). Genome-wide association study using cellular traits identifies a new

regulator of root development in Arabidopsis. Nat. Genet. 46: 77-81.

- Meinke, D., Muralla, R., Sweeney, C., and Dickerman, A. (2008). Identifying essential genes in *Arabidopsis thaliana*. Trends Plant Sci. 13: 483–491.
- Meng, D., Dubin, M., Zhang, P., Osborne, E.J., Stegle, O., Clark, R.M., and Nordborg, M. (2016). Limited contribution of DNA methylation variation to expression regulation in *Arabidopsis thaliana*. PLoS Genet. 12: e1006141.
- Michael, T.P., Jupe, F., Bemm, F., Motley, S.T., Sandoval, J.P., Lanz, C., Loudet, O., Weigel, D., and Ecker, J.R. (2018). High contiguity *Arabidopsis thaliana* genome assembly with a single nanopore flow cell. Nat. Commun. 9: 541.
- Nakamichi, N., Takao, S., Kudo, T., Kiba, T., Wang, Y., Kinoshita, T., and Sakakibara, H. (2016). Improvement of Arabidopsis biomass and cold, drought and salinity stress tolerance by modified circadian clock-associated PSEUDO-RESPONSE REGULATORs. Plant Cell Physiol. 57: 1085–1097.
- Narasimhan, V.M., et al. (2016) Health and population effects of rare gene knockouts in adult humans with related parents. Science **352**: 474–477.
- Olson, M.V. (1999). When less is more: Gene loss as an engine of evolutionary change. Am. J. Hum. Genet. **64:** 18–23.
- Palmieri, N., Kosiol, C., and Schlötterer, C. (2014). The life cycle of Drosophila orphan genes. eLife 3: e01311.
- Phillips, S.J., and Dudik, M. (2008). Modeling of species distributions with Maxent: New extensions and a comprehensive evaluation. Ecography 31: 161–175.
- Phillips, S.J., Anderson, R.P., and Schapire, R.E. (2006). Maximum entropy modeling of species geographic distributions. Ecol. Modell. 190: 231–259.
- Prachumwat, A., and Li, W.H. (2008). Gene number expansion and contraction in vertebrate genomes with respect to invertebrate genomes. Genome Res. 18: 221–232.
- Prieto-Godino, L.L., Rytz, R., Bargeton, B., Abuin, L., Arguello, J.R., Peraro, M.D., and Benton, R. (2016). Olfactory receptor pseudopseudogenes. Nature 539: 93–97.
- Saleheen, D., et al. (2017) Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity. Nature **544:** 235– 239.
- Sas, C., Müller, F., Kappel, C., Kent, T.V., Wright, S.I., Hilker, M., and Lenhard, M. (2016). Repeated inactivation of the first committed enzyme underlies the loss of benzaldehyde emission after the selfing transition in *Capsella*. Curr. Biol. 26: 3313–3319.
- Schild, L., Canessa, C.M., Shimkets, R.A., Gautschi, I., Lifton, R.P., and Rossier, B.C. (1995). A mutation in the epithelial sodium channel causing Liddle disease increases channel activity in the *Xenopus laevis* oocyte expression system. Proc. Natl. Acad. Sci. USA 92: 5699–5703.
- Schindelin, J., et al. (2012) Fiji: An open-source platform for biological-image analysis. Nat. Methods 9: 676–682.
- Schneeberger, K., Ossowski, S., Lanz, C., Juul, T., Petersen, A.H., Nielsen, K.L., Jørgensen, J.E., Weigel, D., and Andersen, S.U. (2009). SHOREmap: Simultaneous mapping and mutation identification by deep sequencing. Nat. Methods 6: 550–551.
- Slater, G.S., and Birney, E. (2005). Automated generation of heuristics for biological sequence comparison. BMC Bioinformatics 6: 31.
- Song, X.J., Huang, W., Shi, M., Zhu, M.Z., and Lin, H.X. (2007). A QTL for rice grain width and weight encodes a previously unknown RING-type E3 ubiquitin ligase. Nat. Genet. **39**: 623–630.
- Swets, J.A. (1988). Measuring the accuracy of diagnostic systems. Science 240: 1285–1293.

- Tanaka, H., Osakabe, Y., Katsura, S., Mizuno, S., Maruyama, K., Kusakabe, K., Mizoi, J., Shinozaki, K., and Yamaguchi-Shinozaki, K. (2012). Abiotic stress-inducible receptor-like kinases negatively control ABA signaling in Arabidopsis. Plant J. 70: 599–613.
- Tang, C., Toomajian, C., Sherman-Broyles, S., Plagnol, V., Guo, Y.L., Hu, T.T., Clark, R.M., Nasrallah, J.B., Weigel, D., and Nordborg, M. (2007). The evolution of selfing in *Arabidopsis thaliana*. Science **317**: 1070–1072.
- Tian, T., Liu, Y., Yan, H., You, Q., Yi, X., Du, Z., Xu, W., and Su, Z. (2017). agriGO v2.0: A GO analysis toolkit for the agricultural community, 2017 update. Nucleic Acids Res. 45): W122–W129.
- Wagner, A. (2005). Distributed robustness versus redundancy as causes of mutational robustness. BioEssays 27: 176–188.
- Warren, D.L., Glor, R.E., and Turelli, M. (2010). ENMTools: A toolbox for comparative studies of environmental niche models. Ecography 33: 607–611.
- Will, J.L., Kim, H.S., Clarke, J., Painter, J.C., Fay, J.C., and Gasch, A.P. (2010). Incipient balancing selection through adaptive loss of aquaporins in natural *Saccharomyces cerevisiae* populations. PLoS Genet. 6: e1000893.
- Wu, W., et al. (2017) A single-nucleotide polymorphism causes smaller grain size and loss of seed shattering during African rice domestication. Nat. Plants 3: 17064.
- Xie, J., et al. (1998) Activating *Smoothened* mutations in sporadic basal-cell carcinoma. Nature **391**: 90–92.

- Yang, H., He, B.Z., Ma, H., Tsaur, S.C., Ma, C., Wu, Y., Ting, C.T., and Zhang, Y.E. (2015). Expression profile and gene age jointly shaped the genome-wide distribution of premature termination codons in a *Drosophila melanogaster* population. Mol. Biol. Evol. 32: 216–228.
- Yoshihara, T., Spalding, E.P., and lino, M. (2013). AtLAZY1 is a signaling component required for gravitropism of the Arabidopsis thaliana inflorescence. Plant J. 74: 267–279.
- Zapata, L., Ding, J., Willing, E.M., Hartwig, B., Bezdan, D., Jiao, W.B., Patel, V., Velikkakam James, G., Koornneef, M., Ossowski, S., and Schneeberger, K. (2016). Chromosome-level assembly of *Arabidopsis thaliana* Ler reveals the extent of translocation and inversion polymorphisms. Proc. Natl. Acad. Sci. USA **113**: E4052– E4060.
- Zhang, Z., Ersoz, E., Lai, C.Q., Todhunter, R.J., Tiwari, H.K., Gore, M.A., Bradbury, P.J., Yu, J., Arnett, D.K., Ordovas, J.M., and Buckler, E.S. (2010). Mixed linear model approach adapted for genome-wide association studies. Nat. Genet. 42: 355–360.
- Zhao, L., Saelao, P., Jones, C.D., and Begun, D.J. (2014). Origin and spread of de novo genes in *Drosophila melanogaster* populations. Science 343: 769–772.
- Zou, Y.P., (2017). Adaptation of *Arabidopsis thaliana* to the Yangtze River basin. Genome Biol. **18:** 239.
- Zufall, R.A., and Rausher, M.D. (2004). Genetic changes associated with floral adaptation restrict future evolutionary potential. Nature 428: 847–850.

Adaptation and Phenotypic Diversification in Arabidopsis through Loss-of-Function Mutations in Protein-Coding Genes

Yong-Chao Xu, Xiao-Min Niu, Xin-Xin Li, Wenrong He, Jia-Fu Chen, Yu-Pan Zou, Qiong Wu, Yong E. Zhang, Wolfgang Busch and Ya-Long Guo *Plant Cell* 2019;31;1012-1025; originally published online March 18, 2019; DOI 10.1105/tpc.18.00791

This information is current as of January 5, 2020

Supplemental Data	/content/suppl/2019/03/18/tpc.18.00791.DC1.html /content/suppl/2019/03/18/tpc.18.00791.DC2.html
References	This article cites 73 articles, 17 of which can be accessed free at: /content/31/5/1012.full.html#ref-list-1
Permissions	https://www.copyright.com/ccc/openurl.do?sid=pd_hw1532298X&issn=1532298X&WT.mc_id=pd_hw1532298X
eTOCs	Sign up for eTOCs at: http://www.plantcell.org/cgi/alerts/ctmain
CiteTrack Alerts	Sign up for CiteTrack Alerts at: http://www.plantcell.org/cgi/alerts/ctmain
Subscription Information	Subscription Information for <i>The Plant Cell</i> and <i>Plant Physiology</i> is available at: http://www.aspb.org/publications/subscriptions.cfm