

Analysis of allele-specific expression reveals *cis*-regulatory changes associated with a recent mating system shift and floral adaptation in *Capsella*

Kim A. Steige¹, Johan Reimegård², Daniel Koenig³, Douglas G. Scofield¹, Tanja Slotte^{1,4,*}

¹Department of Ecology and Genetics, Uppsala University, Uppsala, Sweden

²Science for Life Laboratory, Uppsala University, Uppsala, Sweden

³Max Planck Institute for Developmental Biology, Tübingen, Germany

⁴Department of Ecology, Environment and Plant Sciences, Science for Life Laboratory, Stockholm University, Stockholm, Sweden

*Corresponding author:

Email: tanja.slotte@su.se (TS)

Abstract

Cis-regulatory changes have long been suggested to contribute to organismal adaptation. While *cis*-regulatory changes can now be identified on a transcriptome-wide scale, in most cases the adaptive significance and mechanistic basis of rapid *cis*-regulatory divergence remains unclear. Here, we have characterized *cis*-regulatory changes associated with recent adaptive floral evolution in the selfing plant *Capsella rubella*, which diverged from the outcrosser *Capsella grandiflora* less than 200 kya. We assessed allele-specific expression (ASE) in leaves and flower buds at a total of 18,452 genes in three interspecific F1 *C. grandiflora* x *C. rubella* hybrids. After accounting for technical variation and read-mapping biases using genomic reads, we estimate that an average of 44% of these genes show evidence of ASE, however only 6% show strong allelic expression biases. Flower buds, but not leaves, show an enrichment of genes with ASE in genomic regions responsible for phenotypic divergence between *C. rubella* and *C. grandiflora*. We further detected an excess of heterozygous transposable element (TE) insertions in the vicinity of genes with ASE, and TE insertions targeted by uniquely mapping 24-nt small RNAs were associated with reduced allelic expression of nearby genes. Our results suggest that *cis*-regulatory changes have been important for recent adaptive floral evolution in *Capsella* and that differences in TE dynamics between selfing and outcrossing species could be an important mechanism underlying rapid regulatory divergence.

Author Summary

The role of regulatory changes for adaptive evolution has long been debated. *Cis*-regulatory changes have been proposed to be especially likely to contribute to phenotypic adaptation, because they are expected to have fewer negative side effects than protein-coding changes. Here we assess the regulatory divergence between two closely related plant species that differ in their mating system and floral traits. We directly assess *cis*-regulatory divergence by quantifying the expression levels of both alleles in F1 hybrids of these species, and we find that genes showing *cis*-regulatory divergence are enriched in genomic regions that are responsible for floral and reproductive differences between the species. In combination with information on gene function for genes with *cis*-regulatory divergence in flower buds, this suggests that *cis*-regulatory changes might have been important for morphological differentiation between these species. Additionally we discover that transposable elements, which accumulate differently depending on mating system, might be involved in rapid regulatory divergence. These findings are an important step towards a better understanding of the role and the mechanisms of rapid regulatory divergence between plant species.

Introduction

The molecular nature of genetic changes that contribute to adaptation is a topic of long-standing interest in evolutionary biology. Ever since the discovery of regulatory sequences by Jacob and Monod in the early 1960s [1], there has been a strong focus on the role of regulatory changes for organismal adaptation (*e.g.* [2-7]). This work has mostly centered on changes in *cis*-regulatory elements (CREs), regulatory regions such as promoters or enhancers that are linked to a focal gene.

Due to the modular nature of CREs, *cis*-regulatory changes can alter the expression of the focal gene in a very specific manner, affecting only a particular tissue, cell type, or developmental stage. These changes therefore potentially have fewer negative pleiotropic effects than nonsynonymous mutations in coding regions [3]. For this reason, *cis*-regulatory changes have been suggested to contribute disproportionately to organismal adaptation ([3, 4, 5, 8, 9] but see [10]).

Numerous detailed investigations of single genes have identified causal *cis*-regulatory changes responsible for changes in animal form and color (*e.g.* *Drosophila* wing pigmentation [11]; pelvic reduction [12]; pigmentation [13] and tooth number in stickleback [14]. In yeast, the molecular mechanisms for and mode of selection on *cis*-regulatory variation have begun to be clarified in detail [15, 16]. *Cis*-regulatory changes in individual genes contributing to phenotypic evolution have also been identified in plants, with perhaps the most well-known example being an insertion of a transposable element (TE) into the regulatory region in the *teosinte branched 1* (*tb1*) gene causing increased apical dominance in maize [17, 18]. Other examples include increased tolerance of heavy-metal polluted soils in *Arabidopsis halleri* due to a combination of copy number expansion and *cis*-regulatory changes at the gene *HMA4* [19], *cis*-regulatory variation at the *RCO-A* gene conferring a change in leaf

morphology in *Capsella* [20], and *cis*-regulatory variation at *FLC* conferring variation in vernalization response in *A. thaliana* [21].

With the advent of high-throughput methods for assessing gene expression, the prospects for identifying *cis*-regulatory changes on a transcriptome-wide scale have greatly improved [22]. Genes with *cis*-regulatory changes can be identified based on mapping local expression QTL (*cis*-eQTL) or by assessing allele-specific expression (ASE). Whereas map-based approaches can identify QTL for all genes with expression data, resolution is typically limited. In contrast, ASE studies require the presence of transcribed polymorphisms as well as rigorous bioinformatic approaches, but have greater resolution and can identify individual genes with *cis*-regulatory changes [23].

In *Drosophila* and yeast, transcriptome-wide studies have found that *cis*-regulatory changes or concordant *cis*- and *trans*-regulatory changes may be disproportionately fixed between lineages, which implies the action of directional selection on gene expression during divergence (e.g. [25-27]). Evidence for positive selection on *cis*-regulatory changes has also been found in crop plants, including rice [27] and maize [28].

Recent transcriptome-scale studies have begun to shed light on the mechanistic basis of *cis*-regulatory variation in plants. Studies in *Arabidopsis* have shown that silencing of transposable elements through the RNA-directed methylation pathway may be particularly important, as silencing of TEs through targeting by 24-nt small interfering RNA (siRNA) and subsequent methylation also affects the expression of nearby genes [29, 30]. Transcriptional gene silencing through the RNA-directed methylation pathway has been suggested to be an important mechanism by which regulatory variation is generated both within [31] and between species [30, 32,

33]. While analyses of the population frequencies and age distribution of methylated TE insertions suggest that most methylated TE insertions near genes are deleterious [29, 30], it has been suggested that some proportion of TE insertions might also contribute to organismal adaptation [34]. TE insertions have been selected for during domestication (e.g. maize [18]; domesticated silkworm [35]), and patterns of population differentiation suggest that TEs have contributed to adaptation to temperate environments in *Drosophila* [36]. Studies in *Arabidopsis* [37], maize [38], and rice [39] have also shown that TE insertions can influence stress-induced expression of nearby genes. However, the extent to which TEs contribute to adaptation in the wild is currently not clear for most species.

The crucifer genus *Capsella* is a promising system for assessing the role of *cis*-regulatory changes in association with plant mating system shifts and adaptation. In *Capsella*, genetic and genomic studies are greatly facilitated by the availability of the sequenced reference genome of *Capsella rubella* [40] and because it is feasible to generate crosses among closely related species. *Capsella* harbors four closely related species that vary in both mating system and ploidy: the self-incompatible outcrossing diploid *Capsella grandiflora*, the self-compatible diploids *Capsella rubella* and *Capsella orientalis*, and finally the allopolyploid *Capsella bursa-pastoris* [41].

In *C. rubella*, the transition to selfing occurred relatively recently (~100 kya), and was associated with speciation from an outcrossing progenitor similar to present-day *C. grandiflora* [40, 42-45]. Despite the recent shift to selfing, *C. rubella* already exhibits a derived reduction in petal size and an elevated pollen-ovule ratio, as well as a reduction of the degree of flower opening [46, 47]. *C. rubella* therefore exhibits floral characteristics typical of self-fertilizing plants, a so-called "selfing syndrome". The selfing syndrome of *C. rubella* is associated with improved efficacy of

autonomous self-pollination [46], and regions with quantitative trait loci for floral divergence between *C. rubella* and *C. grandiflora* exhibit an excess of fixed differences and reduced polymorphism in *C. rubella* [47]. Together, these observations suggest that the rapid evolution of the selfing syndrome in *C. rubella* was driven by positive selection. While the molecular genetic basis of the selfing syndrome in *C. rubella* has not been identified, it has been suggested that *cis*-regulatory changes could be involved, and a previous study found many flower and pollen development genes to be differentially expressed in flower buds of *C. grandiflora* and *C. rubella* [40]. As the two species differ in their genomic distribution of TEs, with *C. rubella* harboring fewer TEs close to genes than *C. grandiflora* [48], it is possible that TE silencing through the RNA-directed methylation pathway could constitute a mechanism for *cis*-regulatory divergence in this system.

In this study we assess *cis*-regulatory divergence between *C. grandiflora* and *C. rubella* and investigate the role of *cis*-regulatory changes for floral and reproductive trait divergence in *C. rubella*. We conduct deep sequencing of transcriptomes as well as genomes of *C. grandiflora* x *C. rubella* F1 hybrids to identify genes with *cis*-regulatory divergence in flower buds and leaves, and test whether *cis*-regulatory changes in flowers are overrepresented in genomic regions responsible for adaptive phenotypic divergence. We further conduct small RNA sequencing and test whether TE insertions targeted by uniquely mapping 24-nt siRNAs are associated with *cis*-regulatory divergence. Our results provide insight into the mechanisms and adaptive significance of *cis*-regulatory divergence in association with recent adaptation and phenotypic divergence in a wild plant system.

Results

Many genes exhibit allele-specific expression in interspecific F1 hybrids

In order to quantify ASE between *C. grandiflora* and *C. rubella*, we generated deep whole transcriptome RNAseq data from flower buds and leaves of three *C. grandiflora* x *C. rubella* F1 hybrids (total 52.1 vs 41.8 Gbp with Q \geq 30 for flower buds and leaves, respectively). We included three technical replicates for one F1 in order to examine the reliability of our expression data. For all F1s and their *C. rubella* parents, we also generated deep (38-68x) whole genome resequencing data in order to reconstruct parental haplotypes and account for read mapping biases.

F1 RNAseq reads were mapped with high stringency to reconstructed parental haplotypes specific for each F1, i.e. reconstructed reference genomes containing whole-genome haplotypes for both the *C. grandiflora* and the *C. rubella* parent of each F1 (see Methods). We conducted stringent filtering of genomic regions where SNPs were deemed unreliable for ASE analyses due to e.g. high repeat content, copy number variation, or a high proportion of heterozygous genotypes in an inbred *C. rubella* line (for details, see Methods and S1 text); this mainly resulted in removal of pericentromeric regions (S2 Fig - S5 Fig). After filtering, we identified ~18,200 genes with ~274,000 transcribed heterozygous SNPs that were amenable to ASE analysis in each F1 (Table 1). The mean allelic ratio of genomic read counts at these SNPs was 0.5 (S6 Fig), suggesting that our bioinformatic procedures efficiently minimized read mapping biases. Furthermore, technical reliability of our RNAseq data was high, as indicated by a mean Spearman's ρ between replicates of 0.98 (range 0.94-0.99).

We assessed ASE using a Bayesian statistical method with a reduced false positive rate compared to the standard binomial test [49]. The method uses genomic

read counts to model technical variation in ASE and estimates the global proportion of genes with ASE, independent of specific significance cutoffs, and also yields gene-specific estimates of the ASE ratio and the posterior probability of ASE. The model also allows for and estimates the degree of variability in ASE along the gene, through the inclusion of a dispersion parameter.

Based on this method, we estimate that on average, the proportion of assayed genes with ASE is as high as 44.6% (S6 Table). In general, most allelic expression biases were moderate, and only 5.9% of assayed genes showed ASE ratios greater than 0.8 or less than 0.2 (Figs. 1 and 2). There was little variation in ASE ratios along genes, as indicated by the distribution of the dispersion parameter estimates having a mode close to zero and a narrow range (Figs. 1 and 2). This suggests that unequal expression of differentially spliced transcripts is not a major contributor to regulatory divergence between *C. rubella* and *C. grandiflora* (Figs. 1 and 2). It also suggests limits to ASE patterns arising as stochastic artifacts, which might also tend to create variation in ASE ratios within genes.

For genes with evidence for ASE (hereafter defined as posterior probability of $\text{ASE} \geq 0.95$), there was a moderate shift toward higher expression of the *C. rubella* allele (mean ratio *C. rubella*/total=0.56; Figs. 1 and 2). This shift was present for all F1s, for both leaves and flowers (Figs. 1 and 2). No such shift was apparent for genomic reads, and ratios of genomic read counts for SNPs in genes with ASE were very close to 0.5 (mean ratio *C. rubella*/total=0.51; Figs 1 and 2). Furthermore, qPCR with allele-specific probes for five genes validated our ASE results empirically (S8 Table). This suggests that *C. rubella* alleles are on average expressed at a slightly higher level than *C. grandiflora* alleles in our F1s.

The mean ASE proportion, as well as the absolute number of genes with ASE was greater for leaves (49%; 6010 genes) than for flower buds (40%; 5216 genes), although this difference was largely driven by leaf samples from one of our F1s (Table 1). Most instances of ASE were specific to either leaves or flower buds, and on average, only 15% of genes expressed in both leaves and flower buds showed consistent ASE in both organs (Fig. 3). Many cases of ASE were also specific to a particular F1, and across all three F1s, there were 1305 genes that showed consistent ASE in flower buds, and 1663 in leaves (Fig. 3).

Enrichment of genes with ASE in genomic regions responsible for phenotypic divergence

We used permutation tests to check for an excess of genes showing ASE within previously-identified narrow (<2 Mb) QTL regions responsible for floral and reproductive trait divergence [47]. As the selfing syndrome seems to have a shared genetic basis in independent *C. rubella* accessions [46, 47], we reasoned that genes with consistent ASE across all F1s would be most likely to represent candidate *cis*-regulatory changes underlying QTL. Out of the 1305 genes with ASE in flower buds of all F1s, 85 were found in narrow QTL regions, and this overlap was significantly greater than expected by chance (permutation test, $P=0.03$; Fig. 4; see Methods for details). In contrast, for leaves, there was no significant excess of genes showing ASE in narrow QTL (permutation test, $P=1$; Fig. 4). Thus, the association between QTL and ASE in flower buds is unlikely to be an artifact of locally elevated heterozygosity facilitating both ASE and QTL detection, which should affect analyses of both leaf and flower samples.

List enrichment analyses reveal floral candidate genes with ASE

We conducted list enrichment analyses to characterize the functions of genes showing ASE. There was an enrichment of Gene Ontology (GO) terms involved in defense and stress responses for genes with ASE in flower buds and in leaves (S9 Table). GO terms related to hormonal responses, including brassinosteroid and auxin biosynthetic processes, were specifically enriched among genes with ASE in flower buds (S9 Table). We further identified nineteen genes involved in floral and reproductive development in *A. thaliana*, which are located in QTL regions (see above), and show ASE in flower buds (Table 2). These genes are of special interest as candidate genes for detailed studies of the genetic basis of the selfing syndrome in *C. rubella*.

Intergenic divergence is elevated near genes with ASE

To assess the role of polymorphisms in regulatory regions for ASE, we assessed levels of heterozygosity in intergenic regions within 1 kb of genes that likely contain an elevated proportion of *cis*-regulatory elements, and in previously identified conserved noncoding regions [50] within 5 kb and 10 kb of genes. Genes with ASE were not significantly more likely to be associated with conserved noncoding regions with heterozygous SNPs than genes without ASE. However, levels of intergenic heterozygosity were slightly but significantly higher for genes with ASE than for those without ASE (median heterozygosity values 1 kb upstream of genes of 0.016 vs. 0.014, respectively, S10 Table), suggesting that polymorphisms in regulatory regions upstream of genes might have contributed to *cis*-regulatory divergence.

Enrichment of TEs near genes with ASE

To test whether differences in TE content might contribute to *cis*-regulatory divergence between *C. rubella* and *C. grandiflora*, we examined whether heterozygous TE insertions near genes were associated with ASE. We identified TE insertions specific to the *C. grandiflora* or *C. rubella* parents of our F1s using genomic read data, as in Ågren et al (2014) [48] (Table 3; see Methods). Consistent with their results [48], we found that *C. rubella* harbored fewer TE insertions close to genes than *C. grandiflora* (on average, 482 vs 1154 insertions within 1 kb of genes in *C. rubella* and *C. grandiflora*, respectively). Among heterozygous TE insertions, *Gypsy* insertions were the most frequent (Table 3). There was a significant association between heterozygous TE insertions within 1 kb of genes and ASE, for both leaves and flower buds, and the strength of the association was greater for TE insertions closer to genes (Table 4; Fig. 5). This was true for individual F1s, as well as for all F1s collectively (Table 4; Fig. 5; S11 Table).

TEs targeted by uniquely mapping 24-nt small RNAs are associated with reduced allelic expression of nearby genes

To test whether siRNA-based silencing of TEs might be responsible for the association between TE insertions and ASE in *Capsella*, we analyzed data for flower buds from one of our F1s, for which we had matching small RNA data (see Methods). We selected only those 24-nt siRNA reads that mapped uniquely, without mismatch, to one site within each of our F1s, because uniquely mapping siRNAs have been shown to have a more marked association with gene expression in *Arabidopsis* [29]. For each gene, we then assessed the ASE ratio of the allele on the same chromosome as a TE insertion (i.e. ASE ratios were polarized such that relative ASE was equal to the ratio of the expression of the allele with a TE insertion on the same chromosome

over the total expression of both alleles), and then further examined the influence of nearby siRNAs. Overall, the mean relative ASE was reduced for genes with nearby TE insertions (Fig. 6) with a more pronounced effect for TE insertions within 1 kb (within the gene: Wilcoxon rank sum test, $W = 1392103$, $p\text{-value} = 8.76 \times 10^{-3}$; within 200 bp: Wilcoxon rank sum test, $W = 1903047$, $p\text{-value} = 7.17 \times 10^{-3}$; within 1 kb: Wilcoxon rank sum test, $W = 3687972$, $p\text{-value} = 8.19 \times 10^{-3}$). The magnitude of the effect on ASE was more pronounced for genes near TE insertions targeted by uniquely mapping 24-nt siRNAs (Fig. 6; for genes with a TE insertion within the gene: Wilcoxon rank sum test, $W = 423369$, $p\text{-value} = 1.36 \times 10^{-4}$; within 200 bp: $W = 540926$, $p\text{-value} = 1.82 \times 10^{-5}$; within 1 kb: $W = 983938$, $p\text{-value} = 3.13 \times 10^{-3}$). In contrast, no significant effect on ASE was apparent for genes near TE insertions that were not targeted by uniquely mapping 24-nt siRNAs (Fig. 6). Thus, uniquely mapping siRNAs targeting TE insertions appear to be responsible for the association we observe between ASE and TE insertions.

Discussion

Understanding the causes and consequences of *cis*-regulatory divergence is a long-standing aim in evolutionary genetics. In this study, we have quantified allele-specific expression in order to understand the mechanisms and adaptive significance of *cis*-regulatory changes in association with a recent plant mating system shift.

Our results indicate that many genes, on average over 40%, harbor *cis*-regulatory changes between *C. rubella* and *C. grandiflora*. The proportion of genes with ASE may seem high given the recent divergence (~100 kya) between *C. rubella* and *C. grandiflora* [40, 45]. However, the majority of genes with ASE showed relatively mild allelic expression biases, and while our estimates are higher than those

in a recent microarray-based study of interspecific *Arabidopsis* hybrids (<10%) [32], our results are consistent with recent analyses of RNAseq data from intraspecific F1 hybrids of *Arabidopsis* accessions (~30%) [51]. Somewhat higher levels of ASE were found in a recent study of maize and teosinte (~70% of genes showed ASE in at least one tissue and F1 individual [28]), and using RNAseq data and the same hierarchical Bayesian analysis that we employed, Skelly et al (2011) [49] estimated that a substantially higher proportion, >70% of assayed genes, showed ASE among two strains of *Saccharomyces cerevisiae*. Thus, our estimates of the proportion of genes with ASE fall within the range commonly observed for recently diverged accessions or lines based on RNAseq data.

One of the key motivations for this study was to investigate whether *cis*-regulatory changes contributed to floral and reproductive adaptation to selfing in *C. rubella*. Two lines of evidence support this hypothesis; first, we find an excess of genes with ASE in flower buds within previously identified narrow QTL regions for floral and reproductive traits that harbor a signature of selection [47]. In contrast, no such excess is present for genes with ASE in leaves, suggesting that this observation is not simply a product of higher levels of divergence among *C. rubella* and *C. grandiflora* in certain genomic regions facilitating both QTL delimitation and ASE analysis. Second, we find that genes involved in hormonal responses, including brassinosteroid biosynthesis, are overrepresented among genes with ASE in flower buds, but not in leaves. Based on a study of differential expression and functional information from *Arabidopsis thaliana*, regulatory changes in this pathway were previously suggested to be important for the selfing syndrome in *C. rubella* [40]. While we do not find evidence for ASE the specific genes detected as differentially expressed in [40], our work nonetheless provides additional support for regulatory

changes in the brassinosteroid pathway contributing to the selfing syndrome of *C. rubella*. Future studies should conduct fine-scale mapping and functional validation to fully explore this hypothesis. To facilitate this work, we have identified a set of candidate genes with ASE that are located in genomic regions harboring QTL for floral and reproductive trait divergence between *C. rubella* and *C. grandiflora*. Of particular interest in this list is the gene *JAGGED* (*JAG*). In *A. thaliana*, this gene is involved in determining petal growth and shape by promoting cell proliferation in the distal part of the petal [52, 53]. As *C. rubella* has reduced petal size due to a shortened period of proliferative growth [46], and the *C. rubella* allele is expressed at a lower level than the *C. grandiflora* allele, this gene is a very promising candidate gene for detailed studies of the genetic basis of the selfing syndrome.

Many instances of ASE were specific to a particular individual or tissue, an observation also supported by recent studies (e.g. [28, 32]). This suggests that there is substantial variation in ASE depending on genotype and developmental stage, consistent with the reasoning that *cis*-regulatory changes can have very specific effects, but expression noise is probably also a contributing factor. In our analyses, we took several steps to model and account for technical variation in order to reduce the incidence of false positives. However, it is difficult to completely rule out the possibility that some cases of subtle ASE may not represent biologically meaningful *cis*-regulatory variation. We also cannot fully rule out imprinting effects as potential causes of ASE, because generating reciprocal F1 hybrids was not possible due to seed abortion in our *C. rubella* x *C. grandiflora* crosses. However, we do not expect these effects to make a major contribution to the patterns we observed; in *Arabidopsis*, imprinting effects are only prevalent in endosperm tissue, and are rare in more

advanced stage tissues such as those analyzed here [51, 54, 55], which suggests that imprinting is not likely to be responsible for the patterns we observe.

One somewhat unexpected finding was the subtle global shift in expression levels toward higher relative expression of the *C. rubella* allele in our F1 hybrids. While it is difficult to completely rule out systematic biases in ASE estimation as the cause for this shift, no marked bias was present for the same SNPs and genes in our genomic data, suggesting that if systematic bioinformatic biases are the cause, the effect is specific to transcriptomic reads. While this remains a possibility, it seems unlikely to completely explain the shift in expression that we observe, as we made considerable effort to avoid reference mapping bias, including high stringency mapping of transcriptomic reads to reconstructed parental haplotypes specific to each F1. Similar global shifts toward higher expression of the alleles from one parent have also been observed in F1s of maize and teosinte [28] and *Drosophila* [56]. An even stronger bias toward higher expression of the *A. lyrata* allele was recently observed in F1s of *A. thaliana* and *A. lyrata* [32], and was attributed to interspecific differences in gene silencing.

To investigate potential mechanisms for *cis*-regulatory divergence, we first examined heterozygosity in regulatory regions and conserved noncoding regions close to genes. While genes with ASE in general showed slightly elevated levels of heterozygosity in putatively regulatory regions, there was no enrichment of conserved noncoding regions with heterozygous SNPs close to genes with ASE. It thus seems likely that divergence in regulatory regions in the proximity of genes, but not specifically in conserved noncoding regions, has contributed to global *cis*-regulatory divergence between *C. rubella* and *C. grandiflora*.

To examine biological explanations for the shift toward a higher relative expression of *C. rubella* alleles, we examined the relationship between TE insertions and ASE. As *C. rubella* harbors a lower number of TE insertions near genes than *C. grandiflora*, we reasoned that TE silencing might contribute to the global shift in expression toward higher relative expression of the *C. rubella* allele, with *C. grandiflora* alleles being preferentially silenced due to targeted methylation of nearby TEs, through transcriptional gene silencing mediated by 24-nt siRNAs. Our results are consistent with this hypothesis. Not only is there is an association between genes with TEs and heterozygous TE insertions in our F1s, there is also reduced expression of alleles that reside on the same haplotype as a nearby TE insertion, and the reduction is particularly strong for TEs that are targeted by uniquely mapping siRNAs. In contrast, no effect on ASE is apparent for TEs that are not targeted by uniquely mapping siRNAs. Moreover, the relatively limited spatial scale over which siRNA-targeted TE insertions are associated with reduced expression of nearby genes (<1 kb) is consistent with previous results from *Arabidopsis* [29-31]. We did not directly assess methylation patterns in this study, but it has been shown that data on siRNA targeting is a reliable proxy for TE methylation [29]. While other factors have probably also contributed, these findings suggest that TEs have been important for global *cis*-regulatory divergence between *C. rubella* and *C. grandiflora*.

Why then do *C. rubella* and *C. grandiflora* differ with respect to silenced TEs near genes? In *Arabidopsis*, methylated TE insertions near genes appear to be predominantly deleterious, and exhibit a signature of purifying selection [29]. It is tempting to speculate that the reduced prevalence of TE insertions near genes in *C. rubella* could be due to purging of recessive deleterious alleles that have become exposed to purifying selection due to increased homozygosity in this self-fertilizing

species. Indeed, a recent simulation study has shown that such purging can occur rapidly upon the shift to selfing [57]. However, we prefer the alternative interpretation that deleterious alleles that were rare in the outcrossing ancestor were preferentially lost in *C. rubella*, mainly as a consequence of the reduction in effective population size associated with the shift to selfing in this species. The latter interpretation is more in line with analyses of polymorphism and divergence at nonsynonymous sites, for which *C. rubella* exhibits patterns consistent with a general relaxation of purifying selection [40]. We also note that none of the genes in narrow QTL regions that show ASE in all three F1s harbor nearby heterozygous TE insertions. Our study thus provides no evidence for a contribution of TE silencing to putatively adaptive *cis*-regulatory divergence.

If TE dynamics are generally important for *cis*-regulatory divergence in association with plant mating system shifts, we might expect different effects on *cis*-regulatory divergence depending not only on the genome-wide distribution of TEs, but also on the efficacy of silencing mechanisms in the host [29, 30, 58]. For instance, He et al (2012) [32] found a shift toward higher relative expression of alleles from the outcrosser *A. lyrata*, which harbors a higher TE content, a fact which they attributed to differences in silencing efficacy between *A. thaliana* and *A. lyrata*; indeed, TEs also showed upregulation of the *A. lyrata* allele [33] and *A. lyrata* TEs were targeted by a lower fraction of uniquely mapping siRNAs [30]. In contrast, we found no evidence for a difference in silencing efficacy between *C. rubella* and *C. grandiflora*, which harbor similar fractions of uniquely mapping siRNAs (12% vs 10% uniquely mapping/total 24-nt RNA reads for *C. rubella* and *C. grandiflora*, respectively). Thus, in the absence of strong divergence in silencing efficacy, differences in the spatial distribution of TEs might be more important for *cis*-regulatory divergence. More

studies of ASE in F1s of selfers of different ages and their outcrossing relatives are needed to assess the general contribution of differences in silencing efficacy versus genomic distribution of TE insertions for *cis*-regulatory divergence in association with mating system shifts.

Conclusions We have shown that many genes exhibit *cis*-regulatory changes between *C. rubella* and *C. grandiflora* and that there is an enrichment of genes with floral ASE in genomic regions responsible for phenotypic divergence. In combination with analyses of the function of genes with floral ASE, this suggests that *cis*-regulatory changes might have contributed to the evolution of the selfing syndrome in *C. rubella*. We further observe a general shift toward higher relative expression of the *C. rubella* allele, an observation that can at least in part be explained by elevated TE content close to genes in *C. grandiflora* and reduced expression of *C. grandiflora* alleles due to silencing of nearby TEs. These results support the idea that TE dynamics and silencing are of general importance for *cis*-regulatory divergence in association with plant mating system shifts.

Methods

Plant Material

We generated three interspecific *C. grandiflora* x *C. rubella* F1s by crossing two accessions of the selfer *C. rubella* with three different accessions of the outcrosser *C. grandiflora*, from different populations (S12 Table). All crosses had *C. grandiflora* as the seed parent and *C. rubella* as the pollen donor, as no viable seeds were obtained from reciprocal crosses [47]. Seeds from F1s and their *C. rubella* parental lines were

surface-sterilized and germinated on 0.5 x Murashige-Skoog medium. We transferred one-week old seedlings to soil in pots that were placed in randomized order in a growth chamber (16 h light: 8 h dark; 20° C: 14° C). After four weeks, but prior to bolting, we sampled young leaves for RNA sequencing. Mixed-stage flower buds were sampled 3 weeks later, when all F1s were flowering. To assess data reliability, we collected three separate samples of leaves and flower buds from one F1 individual, and three biological replicates of one *C. rubella* parental line. For genomic DNA extraction, we sampled leaves from all three F1 individuals as well as from their *C. rubella* parents. For small RNA sequencing, we germinated six F2 offspring from one of our F1 individuals and sampled flower buds as described above.

Sample Preparation and Sequencing

We extracted total RNA for whole transcriptome sequencing with the RNEasy Plant Mini Kit (Qiagen, Hilden, Germany), according to the manufacturer's instructions. For small RNA sequencing, we extracted total RNA using the mirVana kit (Life Technologies). For whole genome sequencing, we used a modified CTAB DNA extraction [59] to obtain predominantly nuclear DNA. RNA sequencing libraries were prepared using the TruSeq RNA v2 protocol (Illumina, San Diego, CA, USA). DNA sequencing libraries were prepared using the TruSeq DNA v2 protocol. Sequencing was performed on an Illumina HiSeq 2000 instrument (Illumina, San Diego, CA, USA) to gain 100bp paired end reads, except for small RNA samples for which single end 50 bp reads were obtained. Sequencing was done at the Uppsala SNP & SEQ Technology Platform, Uppsala University, except for accession *C. rubella* Cr39.1 where genomic DNA sequencing was done at the Max Planck Institute of Developmental Biology, Tübingen. In total, we obtained 93.9 Gbp (Q≥30) of RNAseq

data, with an average of 9.3 Gbp per sample. In addition we obtained 45.6 Gbp (Q \geq 30) of DNaseq data, corresponding to a mean expected coverage per individual of 52x, and 106,110,000 high-quality (Q \geq 30) 50 bp small RNA reads. All sequence data has been submitted to the European Bioinformatics Institute (www.ebi.ac.uk), with study accession number: PRJEB9020.

Sequence Quality and Trimming

We merged read pairs from fragment spanning less than 185 nt (this also removes potential adapter sequences) in SeqPrep (<https://github.com/jstjohn/SeqPrep>) and trimmed reads based on sequence quality (phred cutoff of 30) in CutAdapt 1.3 [60]. For DNA and RNAseq reads, we removed all read pairs where either of the reads was shorter than 50 nt. We then analyzed each sample individually using fastQC v. 0.10.1 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) to identify potential errors that could have occurred in the process of amplifying DNA and RNA. We assessed RNA integrity by analyzing the overall depth of coverage over annotated coding genes, using geneBody_coverage.py that is part of the RSeQC package v. 2.3.3 [61]. For DNA reads we analyzed the genome coverage using bedtools v.2.17.0 [62] and removed all potential PCR duplicates using Picard v.1.92 (<http://picard.sourceforge.net>). Small RNA reads were trimmed using custom scripts and CutAdapt 1.3 and filtered to retain only reads of 24 nt length.

Read Mapping and Variant Calling

We mapped both genomic reads and RNAseq reads to the v1.0 reference *C. rubella* assembly [40] (<http://www.phytozome.net/capsella>) using STAR v.2.3.0.1 [63] with default parameters. For genomic reads we modified the default STAR settings to

avoid splitting up reads, and for mapping 24-nt small RNA we used STAR with settings modified to require perfect matches to the parental haplotypes of the F1s as well as to a TE library based on multiple Brassicaceae species and previously used in Slotte et al (2013) [40].

Variant calling was done in GATK v. 2.5-2 [64] according to GATK best practices [65, 66]. Briefly, after duplicate marking, local realignment around indels was undertaken, and base quality scores were recalibrated, using a set of 1,538,085 SNPs identified in *C. grandiflora* [50] as known variants. Only SNPs considered high quality by GATK were kept for further analysis. Variant discovery was done jointly on all samples using the UnifiedGenotyper, and for each F1, genotypes were phased by transmission, by reference to the genotype of its highly inbred *C. rubella* parental accession.

We validated our procedure for calling variants in genomic data by comparing our calls for the inbred line *C. rubella* 1GR1 at 176,670 sites sequenced in a different individual from the same line by Sanger sequencing [67]. Overall, we found 29 calls that differed among the two sets, resulting in an error rate of 0.00016, considerably lower than the level of divergence among *C. rubella* and *C. grandiflora* (0.02 [45]).

Reconstruction of parental haplotypes of interspecific F1s

We reconstructed genome-wide parental haplotype sequences for each interspecific F1 and used these as a reference sequence for mapping genomic and transcriptomic reads for ASE analyses. The purpose of this was to reduce effects of read mapping biases on our analyses of ASE by increasing the number of mapped reads and reducing mismapping that can result when masking heterozygous SNPs in F1s [68].

To reconstruct parental genomes for each F1, we first conducted genomic read

mapping, variant calling and phasing by reference to the inbred *C. rubella* parent as described in the section "Read Mapping and Variant Calling" above. The resulting phased vcf files were used in conjunction with the *C. rubella* reference genome sequence to create a new reference for each F1, containing both of its parental genome-wide haplotypes. Read mapping of both genomic and RNA reads from each F1 was then redone to its specific parental haplotype reference genome, and read counts at all reliable SNPs (see section "Filtering" below) were obtained using Samtools mpileup and a custom software written in javascript by Johan Reimegård. The resulting files with allele counts for genomic and transcriptomic data were used in all downstream analyses of allelic expression biases (see section "Analysis of Allele-Specific Expression" below).

Filtering

We used two approaches to filter the genome assembly to identify regions where we have high confidence in our SNP calls. Genomic regions with evidence for large-scale copy number variation were identified using Control-FREEC [69], and repeats and selfish genetic elements were identified using RepeatMasker (<http://www.repeatmasker.org>). Additionally, we identified genomic regions with unusually high proportions of heterozygous genotype calls in a lab-inbred *C. rubella* line, which is expected to be highly homozygous. Regions with evidence for high proportions of repeats, copy number variation or high proportion of heterozygous calls in the inbred line mainly corresponded to centromeric and pericentromeric regions, and these were removed from consideration in further analyses of allele-specific expression (S2 Fig. - S5 Fig.).

Analysis of allele-specific expression

Analyses of allele-specific expression (ASE) were done using a hierarchical Bayesian method developed by Skelly et al (2011) [49]. This analysis method has a reduced rate of false positives and naturally incorporates replicated data. The method requires read counts at heterozygous coding SNPs for both genomic and transcriptomic data. Genomic read counts are used to fit the parameters of a beta-binomial distribution, in order to obtain an empirical estimate of the distribution of variation in allelic ratios due to technical variation (as there is no true ASE for genomic data on read counts for heterozygous SNPs). This distribution is then used in analyses of RNAseq data where genes are assigned posterior probabilities of exhibiting ASE. Ultimately, this results in an estimate of the posterior probability of ASE for each gene, the mean level of ASE, and the degree of variability of ASE along the gene.

We conducted ASE analyses using the method of Skelly et al (2011) [49] for each of our three F1 individuals. Prior to analyses, we filtered the genomic data to only retain read counts for heterozygous SNPs in coding regions that did not overlap with neighboring genes, and following Skelly et al (2011) [49], we also removed SNPs that were the most strongly biased in the genomic data (specifically, in the 1% tails of a beta-binomial distribution fit to all heterozygous SNPs in each sample), as such highly biased SNPs may result in false inference of variable ASE if retained. The resulting data set showed very little evidence for read mapping bias affecting allelic ratios: the mean ratio of *C. rubella* alleles to total was 0.507 (S6 Fig).

All analyses were run in triplicate and MCMC convergence was checked by comparing parameter estimates across independent runs from different starting points, and by assessing the degree of mixing of chains. For all analyses of RNA counts, we used median estimates of the parameters of the beta-binomial distribution from

analyses of genomic data for all three F1s (S7 Table). Runs were completed on a high-performance computing cluster at Uppsala University (UPPMAX) using the pqr implementation of R (<http://www.pqr-project.org>), for 200,000 generations or a maximum runtime of 10 days. We discarded the first 10% of each run as burn-in prior to obtaining parameter estimates.

ASE Validation by qPCR

We validated ASE results by performing qPCR with TaqMan® Reverse Transcription Reagents (LifeTechnologies, Carlsbad, CA, USA) using oligo(dT)₁₆s to convert mRNA into cDNA using the manufacturers protocol and performed qPCR with the Custom TaqMan® Gene Expression Assay (LifeTechnologies, Carlsbad, CA, USA) with the colors FAM and VIC using manufacturers protocol. The qPCR for both alleles was multiplexed in one well to directly compare the two alleles using a Bio-Rad CFX96 Touch™ Real-Time PCR Detection System (Bio-Rad, Hercules, CA, USA). For further details see S1 Text. To exclude color bias, we tested 5 genes using reciprocal probes with VIC and FAM colorant (S12 Table). The expression difference between the *C. rubella* and *C. grandiflora* allele was quantified using the difference in relative expression between the two alleles, as well as the Quantification Cycle (Cq value). A lower Cq value correlates with a higher amount of starting material in the sample. If the direction of allelic imbalance inferred by qPCR was the same as for ASE inferred by the method by Skelly et al (2011) [49], we considered that the qPCR supported the ASE results.

Enrichment of genes with ASE in genomic regions responsible for phenotypic divergence

We tested whether there was an excess of genes with evidence for ASE (posterior probability of ASE ≥ 0.95 in all three F1 hybrids) in previously identified genomic regions harboring QTL for phenotypic divergence between *C. rubella* and *C. grandiflora* [47]. For this purpose, we concentrated on narrow QTL regions defined in a previous study [47] (i.e. QTL regions with 1.5-LOD confidence intervals < 2 Mb). Significance was based on a permutation test (1000 permutations) in R 3.1.2.

List enrichment tests of GO terms

We tested for enrichment of GO biological process terms among genes with ASE in all of our F1s using Fisher exact tests in the R module TopGO [70]. GO terms were downloaded from TAIR (<http://www.arabidopsis.org>) on September 3rd, 2013, for all *A. thaliana* genes that have orthologs in the *C. rubella* v1.0 annotation, and we only considered GO terms with at least two annotated members in the background set. Separate tests were conducted for leaf and flower bud samples, and background sets consisted of all genes where we could assess ASE.

Intergenic heterozygosity in regulatory and conserved noncoding regions

We quantified intergenic heterozygosity 1 kb upstream of genes using VCFTools [71], and compared levels of polymorphism among genes with and without ASE using a Wilcoxon rank sum test. We further assessed whether there was an enrichment of conserved noncoding elements (identified in Williamson et al (2014) [50]) with heterozygous SNPs within 5 kb of genes with ASE, using Fisher exact tests. Separate tests were conducted for each F1.

Identification of TE insertions and association with ASE

We used PoPoolationTE [72] to identify transposable elements in our F1 parents. While intended for pooled datasets, this method can also be used on genomic reads from single individuals [48]. For this purpose we used a library of TE sequences based on several Brassicaceae species [40]. We used the default pipeline for PoPoolationTE, modified to require a minimum of 5 reads to call a TE insertion, and the procedure in Ågren et al (2014) [48] to determine heterozygosity or homozygosity of TE insertions. Parental origins of TE insertions were inferred by combining information from runs on F1s and their *C. rubella* parents.

We tested whether heterozygous TE insertions within a range of different window sizes close to genes (200 bp, 1 kbp, 2 kbp, 5 kbp, and 10 kbp) were associated with ASE by performing Fisher exact tests in R 3.0.2. We tested whether the expression of the allele on the same chromosome as a nearby (within 1 kbp) TE insertion was reduced compared to ASE at against genes without nearby TE insertions using a Wilcoxon rank sum test. Similar tests were conducted to test for an effect on relative ASE of TE insertions with uniquely mapping siRNAs.

Acknowledgements

The authors thank Michael Nowak, Stockholm University, for valuable comments on the manuscript and Daniel Skelly, Duke University, for helpful advice on ASE analyses. Sequencing was performed by the SNP&SEQ Technology Platform in Uppsala. The facility is part of the National Genomics Infrastructure (NGI) Sweden and Science for Life Laboratory. The SNP&SEQ Platform is also supported by the Swedish Research Council and the Knut and Alice Wallenberg Foundation. The computations were performed on resources provided by SNIC through Uppsala

Multidisciplinary Center for Advanced Computational Science (UPPMAX) under Project b2012122.

References

- 1 Jacob F, Monod J. Genetic regulatory mechanisms in synthesis of proteins. *J Mol Biol.* 1961; 3:318–356.
- 2 King MC, Wilson AC. Evolution at two levels in humans and chimpanzees. *Science* .1975; 188(4184):107-16
- 3 Wray GA. The evolutionary significance of cis-regulatory mutations. *Nat Rev Genet.* 2007; 8:206–216
- 4 Carroll SB. Evo-devo and an expanding evolutionary synthesis: a genetic theory of morphological evolution. *Cell.* 2008; 134:25–36.
- 5 Stern DL, Orgogozo V. The loci of evolution: how predictable is genetic evolution? *Evolution.* 2008; 62(9):2155–2177.
- 6 Wittkopp PJ, Kalay G. Cis-regulatory elements: molecular mechanisms and evolutionary processes underlying divergence. *Nat Rev Genet.* 2012; 13:59–69.
- 7 Albert FW, Kruglyak L. The role of regulatory variation in complex traits and disease. *Nat Rev Genet.* 2015; 16:197–212.
- 8 Doebley J, Lukens L. Transcriptional regulators and the evolution of plant form. *Plant Cell.* 1998; 10:1075–1082.
- 9 Carroll SB. Endless forms: the evolution of gene regulation and morphological diversity. *Cell.* 2000; 101:577–580.
- 10 Hoekstra HE, Coyne JA. The locus of evolution: evo devo and the genetics of adaptation. *Evolution.* 2007; 61:995–1016.
- 11 Prud'homme B, Gompel N, Carroll SB. Emerging principles of regulatory evolution. *Proceedings of the National Academy of Sciences.* 2007; 1:8605–8612.
- 12 Shapiro MD, Bell MA, Kingsley DM. Parallel genetic origins of pelvic reduction in vertebrates. *Proceedings of the National Academy of Sciences.* 2006; 103:13753–13758.
- 13 Miller CT, Beleza S, Pollen AA, Schluter D, Kittles RA, Shriver MD, Kingsley DM. Cis-Regulatory changes in Kit ligand expression and parallel evolution of pigmentation in sticklebacks and humans. *Cell.* 2007; 131:1179–1189.
- 14 Cleves PA, Ellis NA, Jimenez MT, Nunez SM, Schluter D, Kingsley DM et al. Evolved tooth gain in sticklebacks is associated with a cis-regulatory allele of

- Bmp6. Proceedings of the National Academy of Sciences. 2014; 111:13912–13917.
- 15 Chang J, Zhou Y, Hu X, Lam L, Henry C, Green EM et al. The molecular mechanism of a cis-regulatory adaptation in yeast. PLoS Genet 2013; 9:e1003813.
- 16 Metzger BPH, Yuan DC, Gruber JD, Dubeau F, Wittkopp PJ. Selection on noise constrains variation in a eukaryotic promoter. Nature. 2015; doi:10.1038
- 17 Doebley J, Stec A, Hubbard L. The evolution of apical dominance in maize. Nature. 1997; 386:485–488.
- 18 Studer A, Zhao Q, Ross-Ibarra J, Doebley J. Identification of a functional transposon insertion in the maize domestication gene *tb1*. Nat Genet. 2011; 43:1160–1163.
- 19 Hanikenne M, Talke IN, Haydon MJ, Lanz C, Nolte A, Motte P et al. Evolution of metal hyperaccumulation required *cis*-regulatory changes and triplication of *HMA4*. Nature. 2008; 453: 391-395.
- 20 Sicard A, Thamm A, Marona C, Lee YW, Wahl V, Stinchcombe JR et al. Repeated evolutionary changes of leaf morphology caused by mutations to a homeobox gene. Curr Biol. 2014; 24:1880–1886.
- 21 Li P, Filiault D, Box MS, Kerdafrrec E, van Oosterhout C, Wilczek AM et al. Multiple *FLC* haplotypes defined by independent *cis*-regulatory variation underpin life history diversity in *Arabidopsis thaliana*. Genes & Dev. 2014; 28: 1635-1640.
- 22 Fraser HB. Genome-wide approaches to the study of adaptive gene expression evolution: systematic studies of evolutionary adaptations involving gene expression will allow many fundamental questions in evolutionary biology to be addressed. Bioessays. 2011; 33:469–477.
- 23 Pastinen T. Genome-wide allele-specific analysis: insights into regulatory variation. Nat Rev Genet. 2010; 11:533–538.
- 24 Wittkopp PJ, Haerum BK, Clark AG. Regulatory changes underlying expression differences within and between *Drosophila* species. Nat Genet. 2008; 40:346–350.
- 25 Bullard JH, Purdom E, Hansen KD, Dudoit S. Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. BMC Bioinformatics. 2010; 11:94
- 26 Fraser HB, Moses AM, Schadt EE. Evidence for widespread adaptive evolution of gene expression in budding yeast. Proceedings of the National Academy of Sciences. 2010; 107:2977–2982.
- 27 House MA, Griswold CK, Lukens LN. Evidence for selection on gene expression in cultivated rice (*Oryza sativa*). Mol Biol Evol. 2014; 31:1514–1525.

- 28 Lemmon ZH, Bukowski R, Sun Q, Doebley JF. The Role of cis Regulatory Evolution in Maize Domestication. *PLoS Genet.* 2014; 10:e1004745.
- 29 Hollister JD, Gaut BS. Epigenetic silencing of transposable elements: a trade-off between reduced transposition and deleterious effects on neighboring gene expression. *Genome Res.* 2009; 19:1419–1428.
- 30 Hollister JD, Smith LM, Guo Y-L, Ott F, Weigel D, Gaut BS. Transposable elements and small RNAs contribute to gene expression divergence between *Arabidopsis thaliana* and *Arabidopsis lyrata*. *Proceedings of the National Academy of Sciences.* 2011; 108:2322–2327.
- 31 Wang X, Weigel D, Smith LM. Transposon variants and their effects on gene expression in *Arabidopsis*. *PLoS Genet.* 2013; 9:e1003255.
- 32 He F, Zhang X, Hu J, Turck F, Dong X, Goebel U, Borevitz J, de Meaux J. Genome-wide analysis of cis-regulatory divergence between species in the *Arabidopsis* genus. *Mol Biol Evol.* 2012; 29:3385–3395.
- 33 He F, Zhang X, Hu J-Y, Turck F, Dong X, Goebel U, Borevitz JO, de Meaux J. Widespread interspecific divergence in cis-regulation of transposable elements in the *Arabidopsis* genus. *Mol Biol Evol.* 2012; 29:1081–1091.
- 34 Stapley J, Santure AW, Dennis SR. Transposable elements as agents of rapid adaptation may explain the genetic paradox of invasive species. *Mol Ecol.* 2015; doi: 10.1111/mec.13089
- 35 Sun W, Shen Y-H, Han M-J, Cao Y-F, Zhang Z. An adaptive transposable element insertion in the regulatory region of the *EO* gene in the domesticated silkworm, *Bombyx mori*. *Mol Biol Evol.* 2014; 31:3302–3313.
- 36 González J, Karasov TL, Messer PW, Petrov DA. Genome-wide patterns of adaptation to temperate environments associated with transposable elements in *Drosophila*. *PLoS Genet.* 2010; 6:e1000905.
- 37 Ito H, Gaubert H, Bucher E, Mirouze M, Vaillant I, Paszkowski J. An siRNA pathway prevents transgenerational retrotransposition in plants subjected to stress. *Nature.* 2011; 472:115–119.
- 38 Makarevitch I, Waters AJ, West PT, Stitzer M, Hirsch CN, Ross-Ibarra J et al. Transposable elements contribute to activation of maize genes in response to abiotic stress. *PLoS Genet.* 2015; 11:e1004915.
- 39 Naito K, Zhang F, Tsukiyama T, Saito H, Hancock CN, Richardson AO et al. Unexpected consequences of a sudden and massive transposon amplification on rice gene expression. *Nature.* 2009; 461:1130–1134.
- 40 Slotte T, Hazzouri KM, Agren JA, Koenig D, Maumus F, Guo YL et al. 2013. The *Capsella rubella* genome and the genomic consequences of rapid mating system evolution. *Nat Genet.* 2013; 45(7):831-835
- 41 Douglas GM, Gos G, Steige KA, Salcedo A, Holm K, Josephs EB et al. 2015.

- Hybrid origins and the earliest stages of diploidization in the highly successful recent polyploid *Capsella bursa-pastoris*. *Proceedings of the National Academy of Sciences*. 2015; 112(9):2806-2811
- 42 Foxe JP, Slotte T, Stahl EA, Neuffer B, Hurka H, Wright SI. Recent speciation associated with the evolution of selfing in *Capsella*. *Proceedings of the National Academy of Sciences*. 2009; 106:5241–5245.
- 43 Guo Y-L, Bechsgaard JS, Slotte T, Neuffer B, Lascoux M, Weigel D et al. Recent speciation of *Capsella rubella* from *Capsella grandiflora*, associated with loss of self-incompatibility and an extreme bottleneck. *Proceedings of the National Academy of Sciences*. 2009; 106:5246–5251.
- 44 St Onge KR, Källman T, Slotte T, Lascoux M, Palmé AE. Contrasting demographic history and population structure in *Capsella rubella* and *Capsella grandiflora*, two closely related species with different mating systems. *Mol Ecol*. 2011; 20:3306–3320.
- 45 Brandvain Y, Slotte T, Hazzouri KM, Wright SI, Coop G. Genomic Identification of Founding Haplotypes Reveals the History of the Selfing Species *Capsella rubella*. *PLoS Genet*. 2013; 9:e1003754.
- 46 Sicard A, Stacey N, Hermann K, Dessoly J, Neuffer B, Bäurle I et al. Genetics, evolution, and adaptive significance of the selfing syndrome in the genus *Capsella*. *Plant Cell*. 2011; 23:3156–3171.
- 47 Slotte T, Hazzouri KM, Stern D, Andolfatto P, Wright SI. Genetic architecture and adaptive significance of the selfing syndrome in *Capsella*. *Evolution*. 2012; 66:1360–1374.
- 48 Agren JA, Wang W, Koenig D, Neuffer B, Weigel D, Wright SI. Mating system shifts and transposable element evolution in the plant genus *Capsella*. *BMC Genomics*. 2014; 15:602.
- 49 Skelly DA, Johansson M, Madeoy J, Wakefield J, Akey JM. A powerful and flexible statistical framework for testing hypotheses of allele-specific gene expression from RNA-seq data. *Genome Res*. 2011; 21:1728–1737.
- 50 Williamson RJ, Josephs EB, Platts AE, Hazzouri KM, Haudry A, Blanchette M et al. Evidence for widespread positive and negative selection in coding and conserved noncoding regions of *Capsella grandiflora*. *PLoS Genet*. 2014; 10:e1004622.
- 51 Cubillos FA, Stegle O, Grondin C, Canut M, Tisné S, Gy I et al. Extensive cis-regulatory variation robust to environmental perturbation in *Arabidopsis*. *Plant Cell*. 2014; 26:4298–4310.
- 52 Sauret-Güeto S, Schiessl K, Bangham A, Sablowski R, Coen E. JAGGED controls *Arabidopsis* petal growth and shape by interacting with a divergent polarity field. *Plos Biol*. 2013; 11:e1001550.
- 53 Schiessl K, Muiño JM, Sablowski R. *Arabidopsis* JAGGED links floral organ

- patterning to tissue growth by repressing Kip-related cell cycle inhibitors. *Proceedings of the National Academy of Sciences*. 2014; 111:2830–2835.
- 54 Scott RJ, Spielman M, Bailey J, Dickinson HG. Parent-of-origin effects on seed development in *Arabidopsis thaliana*. *Development*. 1998; 125:3329–3341.
- 55 Wolff P, Weinhofer I, Seguin J, Roszak P, Beisel C, Donoghue MT et al. High-resolution analysis of parent-of-origin allelic expression in the *Arabidopsis* Endosperm. *PLoS Genet*. 2011; 7:e1002126.
- 56 McManus CJ, Coolon JD, O'Duff M, Eipper-Mains J, Graveley BR, Wittkopp PJ. Regulatory divergence in *Drosophila* revealed by mRNA-seq. *Genome Res*. 2010; 20: 816-825
- 57 Arunkumar R, Ness RW, Wright SI, Barrett SCH. The Evolution of Selfing Is Accompanied by Reduced Efficacy of Selection and Purging of Deleterious Mutations. *Genetics*. 2014; 199(3):817-829
- 58 Agren JA, Wright SI. Selfish genetic elements and plant genome size evolution. *Trends Plant Sci*. 2015; doi: 10.1016/j.tplants.2015.03.007
- 59 Doyle JJ, Doyle JL. A rapid DNA isolation procedure for small quantities of fresh leaf tissue. *Phytochem bull*. 1987; 19: 11-15.
- 60 Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal*. 2011; 17:10–12.
- 61 Wang L, Wang S, Li W. RSeQC: quality control of RNA-seq experiments. *Bioinformatics*. 2012; 28:2184–2185.
- 62 Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*. 2010; 26:841–842.
- 63 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013; 29:15–21.
- 64 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res*. 2010; 20:1297–1303.
- 65 DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet*. 2011; 43:491–498.
- 66 Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A et al. From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr Protoc Bioinformatics*. 2013; 11:11.10.1–11.10.33.
- 67 Slotte T, Foxe JP, Hazzouri KM, Wright SI. Genome-wide evidence for efficient positive and purifying selection in *Capsella grandiflora*, a plant species with a large effective population size. *Mol Biol Evol*. 2010; 27:1813-1821.

- 68 Degner JF, Marioni JC, Pai AA, Pickrell JK, Nkadori E, Gilad Y et al. Effect of read-mapping biases on detecting allele-specific expression from RNA-sequencing data. *Bioinformatics*. 2009; 25:3207–3212.
- 69 Boeva V, Zinovyev A, Bleakley K, Vert J-P, Janoueix-Lerosey I, Delattre O et al. Control-free calling of copy number alterations in deep-sequencing data using GC-content normalization. *Bioinformatics*. 2011; 27:268–269.
- 70 Alexa A, Rahnenführer J, Lengauer T. Improved scoring of functional groups from gene expression data by decorrelating GO graph structure. *Bioinformatics*. 2006; 22:1600–1607.
- 71 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA et al. The variant call format and VCFtools. *Bioinformatics*. 2011; 27:2156–2158.
- 72 Kofler R, Betancourt AJ, Schlötterer C. Sequencing of pooled DNA samples (Pool-Seq) uncovers complex dynamics of transposable element insertions in *Drosophila melanogaster*. *PLoS Genet*. 2012; 8:e1002487.

Figure Legends

Fig. 1. ASE results for flower buds. Distributions of ASE ratios (*C. rubella*/Total) for all assayed genes (A, B, C), and for genes with at least 0.95 posterior probability of ASE (D, E, F). Ratio of *C. rubella* to total for genomic reads, for genes with significant ASE (G, H, I), and the distribution of the dispersion parameter that quantifies variability in ASE across genes (J, K, L). All distributions are shown for each of the three interspecific F1s inter 3.1 (left), inter4.1 (middle) and inter5.1 (right).

Fig. 2. ASE results for leaves. Distributions of ASE ratios (*C. rubella*/Total) for all assayed genes (A, B, C), and for genes with at least 0.95 posterior probability of ASE (D, E, F). Ratio of *C. rubella* to total for genomic reads, for genes with significant ASE (G, H, I), and the distribution of the dispersion parameter that quantifies variability in ASE across genes (J, K, L). All distributions are shown for each of the three interspecific F1s inter 3.1 (left), inter4.1 (middle) and inter5.1 (right).

Fig. 3. Many cases of ASE are specific to individuals or samples. Venn diagrams showing intersections of genes with ASE in flower buds (A) and leaves (B) of the three F1 individuals, and (C) in all leaf and flower samples, for the set of genes assayed in all F1s.

Fig. 4. Enrichment of genes with ASE in narrow QTL regions. There is an excess of genes with ASE in narrow QTL regions for flower buds (A) but not for leaves (B). Histograms show the distribution of numbers of genes with ASE that fall within

narrow QTL regions, based on 1000 random permutations of the observed number of genes with ASE among all genes where we could assess ASE. Arrows indicate the observed number of genes with ASE that are located in narrow QTL regions.

Fig. 5. Enrichment of TEs near genes with ASE. The Figure shows odds ratios (OR) of the association between genes with ASE and TEs, with TE insertions scored in four different window sizes. Odds ratios for flower buds are shown for all three F1s studies, with values for flower buds in black and leaves in grey.

Fig. 6. The effect of TE insertions on relative allelic expression. Boxplots show the relative allelic expression (expression of the allele on same haplotype as TE insertion relative to expression of both alleles) for genes near heterozygous TE insertions, scored in a range of window sizes ranging from 0 bp (within the gene) to 10 kbp from the gene. A. The relative allelic expression is reduced for genes with nearby TE insertions. B. The degree of reduction of relative allelic expression is stronger for genes near TE insertions targeted by uniquely mapping siRNA. C. There is no reduction of relative allelic expression for genes near TE insertions that are not targeted by uniquely mapping siRNA.

Tables

Table 1. Genes amenable to analysis of ASE in flower bud and leaf samples from the three *C. grandiflora* x *C. rubella* F1s, counts of genes with evidence for ASE and the estimated false discovery rate (FDR) and proportion of genes with ASE.

F1 designation	Sample	Genes amenable to ASE analysis ^a	Analyzed genes ^b	Heterozygous SNPs in analyzed genes	Number of genes with ASE PP $\geq 0.95^c$	FDR	ASE proportion ^d
Inter3.1	Flower buds	18299	16857	262120	4728	0.0013	0.38
Inter4.1		18270	17837	272126	5744	0.0022	0.42
Inter5.1		18144	17448	262696	5176	0.0020	0.40
Inter3.1	Leaves	18299	14877	238786	5105	0.0012	0.44
Inter4.1		18270	15784	249181	8129	0.0024	0.62
Inter5.1		18144	15478	240653	4795	0.0018	0.41

^aTotal number of genes with heterozygous SNPs in coding regions remaining after filtering.

^bNumber of genes amenable to ASE analyses with expression data in at least one of the replicates of the sample.

^cGenes with evidence for ASE (posterior probability ≥ 0.95).

^dDirect estimate of the ASE proportion independent of significance cutoffs.

Table 2. Selfing syndrome candidate genes identified based on ASE, QTL information, and Arabidopsis annotation.

<i>C. rubella</i> ortholog	Arabidopsis ortholog	Arabidopsis annotation	GO biological process terms related to floral and reproductive development
Carubv10012851m ^{a,b}	AT3G24340	CHR40	regulation of flower development
Carubv10016094m ^{a,b}	AT3G24650	ATABI3, ABI3, SIS10	embryo development, cotyledon development
Carubv10007602m ^{a,b}	AT4G21600	ENDO5	brassinosteroid biosynthetic process
Carubv10000655m ^{b,d}	AT5G08130	BIM1	brassinosteroid mediated signaling pathway, primary shoot apical meristem specification
Carubv10006681m ^{b,d}	AT4G28720	YUC8	brassinosteroid mediated signaling pathway
Carubv10021883m ^{a,c}	AT1G68480	JAG	sepal formation, flower development, abaxial cell fate

			specification, anther development, carpel development, stamen development, petal formation, specification of floral organ identity
Carubv10021345m ^{a,c}	AT1G68640	PAN, TGA8	petal formation, sepal formation, regulation of flower development
Carubv10013321m ^{a,c}	AT3G22420	ATWNK2, WNK2, ZIK3	photoperiodism, flowering
Carubv10016406m ^{a,c}	AT3G23270	-	pollen tube growth
Carubv10014951m ^{a,c}	AT3G23440	EDA6, MEE37	megagametogenesis
Carubv10014152m ^{a,c}	AT3G23630	ATIPT7, IPT7	pollen tube growth, reciprocal meiotic recombination
Carubv10010238m ^{a,c}	AT3G62210	EDA32	polar nucleus fusion
Carubv10004312m ^{a,c}	AT4G16760	ATACX1, ACX1	pollen development
Carubv10005585m ^{a,c}	AT4G17030	AT-EXPR, EXPR, ATEXLB1,	sexual reproduction

ATEXPR1, EXLB1			
Carubv10007441m ^{a,c}	AT4G20370	TSF	regulation of flower development, photoperiodism, flowering, positive regulation of flower development
Carubv10004229m ^{a,c}	AT4G20910	CRM2, HEN1	specification of floral organ identity, floral organ formation, petal formation, regulation of flower development, sepal formation, meristem initiation, meristem development, ovule development
Carubv10015623m ^{a,c}	AT4G21380	ARK3, RK3	recognition of pollen
Carubv10007227m ^{a,c}	AT4G21530	APC4	ovule development
Carubv10007633m ^{a,c}	AT4G21590	ENDO3	petal development, stamen development, pollen tube growth, ovule development

^alocated within narrow QTL regions

^bASE in all three F1s

^cASE in the F1 with data for three replicates, but not in all three F1s

^dlocated within QTL regions, but not narrow QTL regions

Table 3. Mean number of TE insertions in three interspecific F1s. The table shows the overall number, as well as heterozygous insertions with parent of origin information.

TE family	Mean copy number	Heterozygous insertions	Insertions specific to the <i>C. rubella</i> parental genome	Insertions specific to the <i>C. grandiflora</i> parental genome
CACTA	84	40	10	30
Copia	710	483	144	339
Gypsy	1124	602	153	449
Harbinger	176	109	26	83
hAT	83	55	16	40
Helitron	236	127	30	97
LINE	229	165	38	128
MuDR	203	109	28	81
SINE	113	92	9	83
Total	2958	1782	454	1330

Table 4. Enrichment of TEs near genes with ASE. The table shows mean counts over all three F1s, and Fisher exact test *P*-values. The four categories of counts correspond to numbers of genes with ASE (posterior probability of ASE ≥ 0.95) and TE insertions within a specific window size near the gene (+ASE,+TE), with ASE but without TEs (+ASE,-TE), without ASE but with TE insertions (-ASE,+TE), and with neither ASE nor TEs (-ASE,-TE). NS indicates not significant.

Sample	Window size (bp)	+ASE,+TE	+ASE,-TE	-ASE,+TE	-ASE,-TE	<i>P</i>
Flower buds	200	113	5103	136	12029	4.32×10^{-19}
	1000	218	4998	339	11826	5.07×10^{-16}
	2000	307	4909	540	11624	6.53×10^{-12}
	5000	566	4650	1108	11057	8.22×10^{-10}
	10000	958	4258	2006	10159	2.32×10^{-7}
Leaves	200	108	5902	115	9255	8.52×10^{-7}
	1000	216	5793	277	9093	1.49×10^{-4}
	2000	317	5693	435	8935	2.25×10^{-3}
	5000	595	5415	877	8493	NS
	10000	1027	4983	1576	7795	NS

Supporting Information

S1 Text. Detailed procedures for filtering genomic regions, qPCR details.

S2 Fig. Filtering of genomic regions on scaffolds 1 and 2. Genomic regions kept for analysis on scaffolds 1 and 2 after filtering for copy number variation, high repeat density, and elevated proportions of heterozygous SNP calls in an inbred *C. rubella* line. The top panels show the proportion kept for analysis after all filtering steps, and the lower panels show the proportion kept for analysis after each filtering stage, with copy number variation, fraction repeats and fraction heterozygous SNP calls indicated by purple lines.

S3 Fig. Filtering of genomic regions on scaffolds 3 and 4. Genomic regions kept for analysis on scaffolds 3 and 4 after filtering for copy number variation, high repeat density, and elevated proportions of heterozygous SNP calls in an inbred *C. rubella* line. The top panels show the proportion kept for analysis after all filtering steps, and the lower panels show the proportion kept for analysis after each filtering stage, with copy number variation, fraction repeats and fraction heterozygous SNP calls indicated by purple lines.

S4 Fig. Filtering of genomic regions on scaffolds 5 and 6. Genomic regions kept for analysis on scaffolds 5 and 6 after filtering for copy number variation, high repeat density, and elevated proportions of heterozygous SNP calls in an inbred *C. rubella* line. The top panels show the proportion kept for analysis after all filtering steps, and the lower panels show the proportion kept for analysis after each filtering stage, with

copy number variation, fraction repeats and fraction heterozygous SNP calls indicated by purple lines.

S5 Fig. Filtering of genomic regions on scaffolds 7 and 8. Genomic regions kept for analysis on scaffolds 7 and 8 after filtering for copy number variation, high repeat density, and elevated proportions of heterozygous SNP calls in an inbred *C. rubella* line. The top panels show the proportion kept for analysis after all filtering steps, and the lower panels show the proportion kept for analysis after each filtering stage, with copy number variation, fraction repeats and fraction heterozygous SNP calls indicated by purple lines.

S6 Fig. Distribution of allelic ratios in genomic DNA.

S7 Table. Parameter estimates and information on ASE analyses.

S8 Table. qPCR results.

S9 Table. Gene Ontology list enrichment analysis results. Gene Ontology (GO) terms significantly enriched (Weighted Fisher $P \leq 0.01$) among genes with ≥ 0.95 posterior probability of ASE in flowers or leaves of all three F1s. Only GO terms with more than two annotated genes are shown.

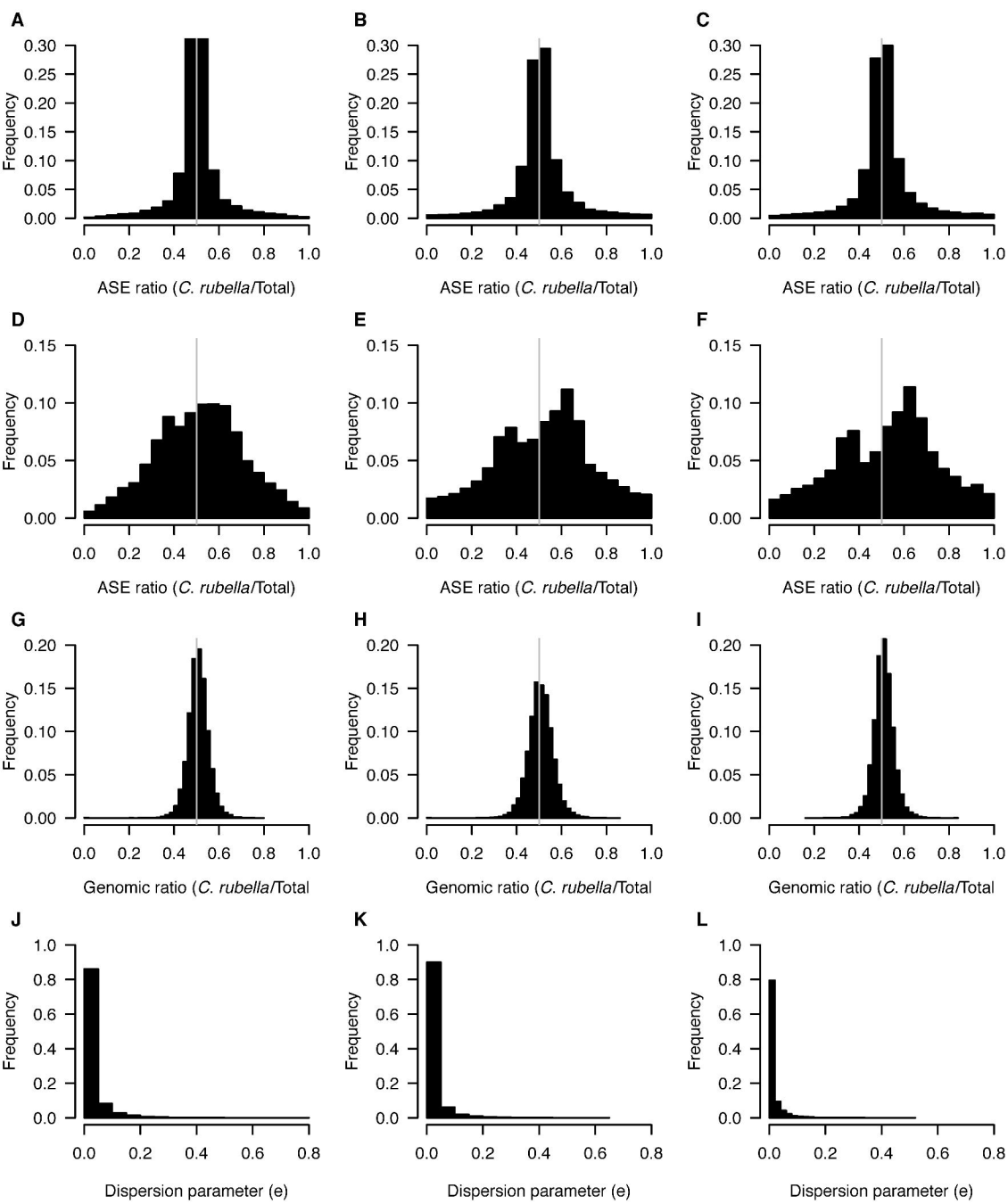
S10 Table. Heterozygosity in intergenic regions 1 kb upstream of genes with and without ASE. The table shows median heterozygosity for each F1, for genes with and

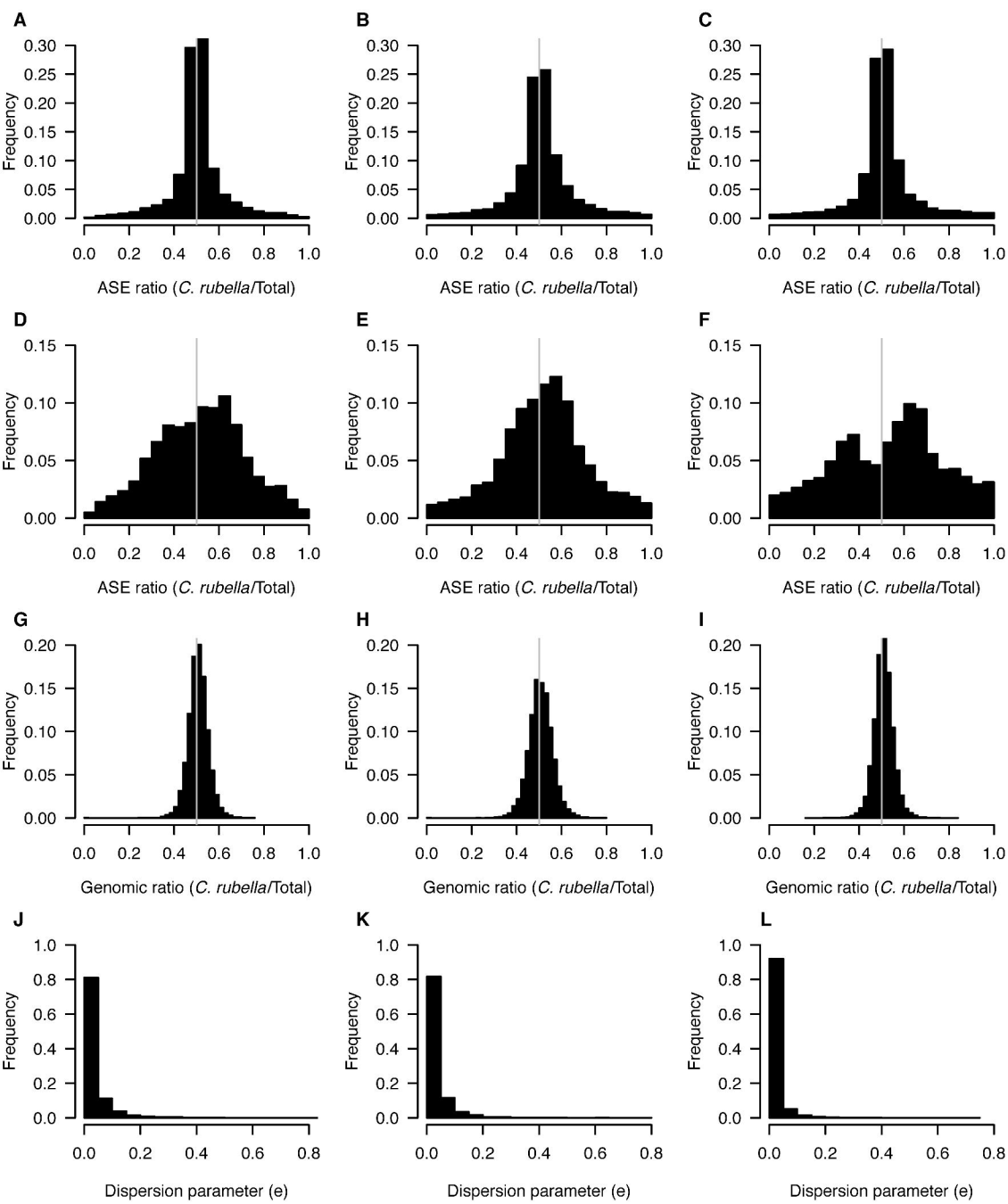
without ASE in flower buds and leaves, and results of Wilcoxon rank sum tests of a difference in median heterozygosity upstream of genes with and without ASE.

S11 Table. Enrichment of TEs near genes with ASE. The table shows odds ratios (OR), as well as the number of genes with ASE and heterozygous TE insertions (+ASE, +TE), with ASE but without heterozygous TE insertions (+ASE, -TE), without ASE and heterozygous TE insertions (-ASE, -TE) and P-values for Fisher exact tests. All values are shown for TEs scored within a range of window sizes within the gene, from 200 bp to 10000 bp.

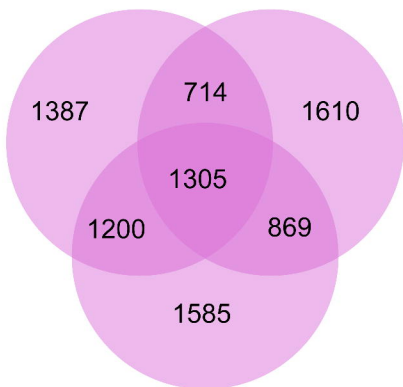
S12 Table. qPCR primers and probes.

S13 Table. Designations of interspecific F1s and geographical origins of parental accessions.

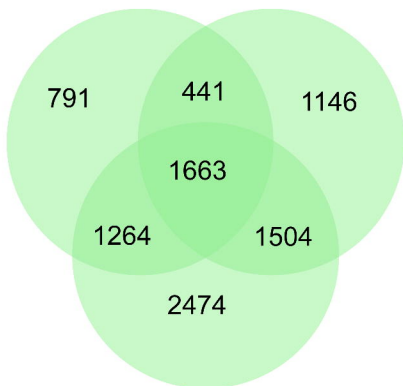




A



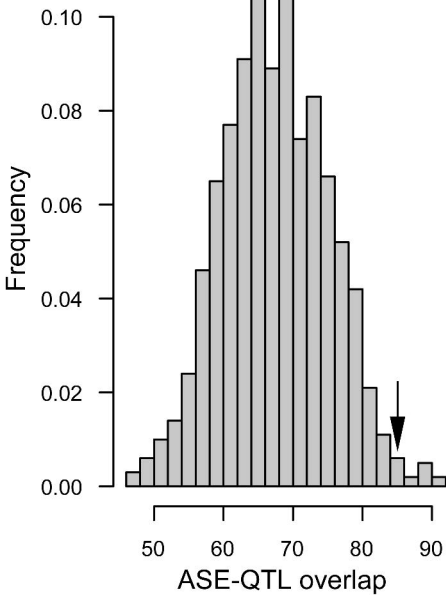
B



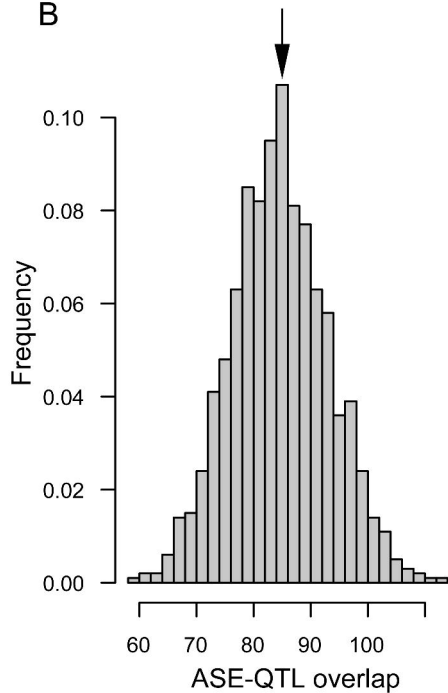
C



A



B



OR

