

1 **Genetic architecture and evolution of the *S* locus supergene in *Primula vulgaris***

2 Jinhong Li^{1,2,*}, Jonathan M. Cocker^{1,2,*}, Jonathan Wright³, Margaret A. Webster^{1,2}, Mark McMullan³,
3 Sarah Dyer^{3,†}, David Swarbreck³, Mario Caccamo^{3,†}, Cock van Oosterhout⁴ & Philip M. Gilmartin^{1,2,‡}

4 ¹ School of Biological Sciences, University of East Anglia, Norwich Research Park, Norwich,
5 NR4 7TJ, UK.

6 ² John Innes Centre, Norwich Research Park, Norwich NR4, 7UH, United Kingdom.

7 ³ The Earlham Institute, Norwich Research Park, Norwich, NR4 7UH, United Kingdom.

8 ⁴ School of Environmental Sciences, University of East Anglia, Norwich Research Park, Norwich,
9 NR4 7TJ, UK.

10

11 * These authors contributed equally.

12 † Current address: National Institute for Agricultural Botany, Huntingdon Road, Cambridge,
13 CB3 0LE UK.

14 ‡ Corresponding author

15

16 **Summary**

17 Darwin's studies on heterostyly in *Primula* described two floral morphs, pin and thrum, with reciprocal
18 anther and stigma heights that promote insect-mediated cross-pollination. This key innovation evolved
19 independently in several angiosperm families. Subsequent studies on heterostyly in *Primula*
20 contributed to the foundation of modern genetic theory and the neo-Darwinian synthesis. The
21 established genetic model for *Primula* heterostyly involves a diallelic *S* locus comprising several genes,
22 with rare recombination events that result in self-fertile homostyle flowers with anthers and stigma at
23 the same height. Here we reveal the *S* locus supergene as a tightly-linked cluster of thrum-specific
24 genes that are absent in pins. We show that thrums are hemizygous not heterozygous for the *S* locus,
25 which suggests that homostyles do not arise by recombination between *S* locus haplotypes as previously
26 proposed. Duplication of a floral homeotic gene 51.7 MYA, followed by its neofunctionalisation,
27 created the current *S* locus assemblage which led to floral heteromorphy in *Primula*. Our findings
28 provide new insights into the structure, function and evolution of this archetypal supergene.

29 **Introduction**

30 Heterostyly evolved independently in at least 28 families of animal-pollinated angiosperms¹. In the
31 Primulaceae the majority of species² produce dimorphic flowers³, a characteristic inherited as a simple
32 Mendelian trait; alleles are defined as *S* (*Short style*) and *s* (*long style*)⁴. The two floral forms are known
33 as pin and thrum; thrums behave as heterozygous *S/s* and pins homozygous *s/s*⁵. Classical genetic
34 studies on mutation, linkage and recombination by Bateson⁵, Bridges⁶, Ernst^{7,8}, Haldane⁹, Darlington¹⁰,
35 and others, established *Primula* as an early genetic model, and led to the definition of a co-adapted
36 linkage group of three genes at the *S* locus, *G* (*Griffel* (*style length*)), *P* (*Pollen*) and *A* (*Antheren* (*anther*)
37 *position*)⁷, which control distinct aspects of heteromorphic flower development; this locus defined the
38 archetypal supergene¹¹. Studies of heterostyly in *Primula* contributed significantly to the foundation of
39 modern genetic theory and the neo-Darwinian synthesis. Supergenes have subsequently been shown to
40 control other multi-trait complex phenotypes in plants, animals and fungi¹².

41 Pin flowers have a long style and low anthers, thrum flowers have a short style and high anthers (Fig.1a).
42 This reciprocal herkogamy promotes insect-mediated cross-pollination between floral morphs, which
43 actively enhances efficiency of reciprocal pollen transfer¹³; such biotic pollination is associated with an
44 elevated speciation rate in angiosperms¹⁴. Differences in stigma shape, papillae length, pollen size and
45 corolla mouth diameter characterise dimorphic *Primula* flowers^{3,15}; a sporophytic self-incompatibility
46 system¹⁶ inhibits intra-morph pollination¹, with different efficacy in different *Primula* species¹⁷. Self-
47 fertile homostyle flowers (Fig.1a) occasionally occur¹³; although originally considered mutants⁸, later
48 studies led to the widely-accepted view that self-fertile homostyles arise by recombination in
49 heterozygous thrums between dominant (*GPA*) and recessive (*gpa*) haplotypes, with associated
50 disruption of coupling between male and female self-incompatibility functions (e.g. *gPA* and *Gpa*)^{18,19}.
51 This interpretation defined the order of genes at the *Primula S* locus, and has formed the backdrop to
52 the last 60 years of research into the *S* locus supergene, including models on the evolution of
53 heterostyly^{20,21} and population genetic analyses²²⁻²⁴ in natural homostyle populations^{25,26}.

54 More recent studies aimed at identifying *S* locus genes involved examination of flower development²⁷,
55 analysis of differentially-expressed floral genes²⁸, characterisation of *S* locus-linked sequences^{29,30},
56 molecular genetic analysis of *S*-linked mutant phenotypes³¹⁻³³, creation of genetic and physical
57 maps^{34,35}, assembly of a partial genome sequence³⁶, and construction of BAC contigs spanning the *S*
58 locus³⁵. Despite these extensive investigations, the genetic architecture of the *S* locus has, until now,
59 been an unresolved enigma. Here we compare the *S* haplotype sequences from pin, thrum, long
60 homostyle and short homostyle plants. The *s* haplotype lacks a 278 kb sequence containing five thrum-
61 specific genes present in thrum and homostyles; thrums are therefore hemizygous not heterozygous for
62 the *S* locus. We demonstrate that this 278 kb region is the only thrum-specific genomic region
63 transcribed in flowers, and by genetic and natural population analyses demonstrate complete linkage to
64 the *S* locus; our data indicate that homostyles cannot occur by recombination as proposed. We also
65 provide an estimate of the evolutionary-age of assembly for the *S* locus supergene.

66 **Identification and assembly of the *S* locus**

67 We previously used four *S*-linked probes to assemble two BAC contigs flanking the *S* locus; these were
68 integrated into a genetic map with the gap between contigs predicted to contain some, or all, of the *S*
69 locus genes³⁵. A fifth *S*-linked probe³³, *GLO^T*, also identified a BAC clone which we could not position
70 relative to our *S* locus map³⁵. In parallel, we initiated the *de novo* assembly of a *P. vulgaris* reference
71 genome using a self-fertile homozygous long homostyle (*S^{LH1}/S^{LH1}*) from the Somerset population
72 identified by Crosby²⁵. We also generated genome sequence data from individual pin (*s/s*) and thrum
73 (*S/s*) plants, pools of their pin and thrum progeny, and a short homostyle (*S^{SH1}/s*)³⁵ (Supplementary
74 Table 1a). Fig. 1a shows relevant floral phenotypes and genotypes.

75 Using the *GLO^T* BAC (BAC70F11) we searched a long homostyle genome assembly (Supplementary
76 Table 1b) to identify and link two genome sequence contigs. This step initiated the assembly of a
77 contiguous 455,881 bp sequence encompassing the entire *S^{LH1}* haplotype from this highly homozygous
78 inbred line (Supplementary Fig. 1a). This assembly contains a 278,470 bp sequence which is absent
79 from pins and flanked by a ~3 kb tandem repeat that is present only as a single copy in the *s* haplotype

80 (Fig. 1); each repeat contains a *Cyclin-like F box (CFB)* gene. We therefore focused on this region as
81 the presumptive *S* locus. Sequences flanking the S^{LHI} 278,470 bp region on the left (75,084 bp), and
82 right (96,327 bp) share extensive similarity to the *s* haplotype (Fig. 1b, Supplementary Fig. 1b).

83 Next, we designed PCR primers for left- and right-border regions of S^{LHI} , and separate *s* haplotype-
84 specific primers (Supplementary Table 2). Analyses with pin, thrum, long and short homostyle genomic
85 DNA confirmed pin as *s/s*, and long homostyle as homozygous S^{LHI}/S^{LHI} (Fig. 1c). Supported by
86 sequence alignment (Supplementary Sequence Analysis 1, 2), these data also show that thrum and the
87 short homostyle share the same left- and right-border sequences as the long homostyle, and that they
88 are both heterozygous for the *s* and S^{LHI} flanking markers (Fig. 1c). The established model defines
89 homostyles as recombinants between *S* and *s* haplotypes; if this is the case, long and short homostyles
90 should possess reciprocal combinations of *s* and *S* haplotype left and right border sequences, but they
91 do not (Fig. 1c).

92 **Comparative analysis of *S* haplotypes**

93 We then focused on the 278 kb region from S^{LHI} that is absent from the *s* haplotype. This region contains
94 five predicted gene models, CCM^T , GLO^T , CYP^T , PUM^T and KFB^T which were manually curated and
95 are supported by RNA-Seq data as thrum-specific in expression; four other models identify transposon
96 sequences which were discounted as functional *S* locus genes and excluded from further analysis
97 (Supplementary Fig. 2a, b). CCM^T (Conserved Cysteine Motif) encodes a protein with a C-terminal
98 domain that is conserved in monocots and dicots (Supplementary Sequence Analysis 3); proteins
99 containing this novel domain are rich in either proline or negatively charged amino acids. One of these,
100 $PIG93$ from *Petunia x hybrida*, is a partner of $PSK8$, a protein involved in brassinolide signalling³⁷. A
101 second CCM -like gene with 90% sequence similarity is found in both pin and thrum genomes. GLO^T
102 was originally defined as a thrum-specific allele³³ of *P. vulgaris* GLO , a floral homeotic gene
103 responsible for the *S* locus-linked mutant phenotype *Hose in Hose*^{32,35}. These data show GLO and GLO^T
104 as distinct loci; the encoded proteins share 82% sequence identity but the *Hose in Hose* mutation, in
105 which GLO is dominantly up-regulated, does not affect heterostyly³². CYP^T encodes a cytochrome P450
106 similar to *Arabidopsis* $CYP72B1$, a brassinolide 26-hydroxylase³⁸. CYP^T is one of four $CYP72$ class

107 genes in the *P. vulgaris* thrum genome, the other three are present in both pin and thrum; the closest
108 encodes a protein with 65% sequence identity to *CYP^T*. *PUM^T* encodes a Pumilio-like³⁹ RNA-binding
109 protein, and *KFB^T* encodes a protein with similarity to the *Arabidopsis* Kiss-Me-Deadly Kelch repeat F
110 Box protein involved in regulating cytokinin activity⁴⁰; Both *PUM^T* and *KFB^T* are unique to the 278 kb
111 region with no homologues found in our pin genome sequence. The tandemly duplicated sequences
112 flanking the *S* locus contain Cyclin-like F Box genes, *CFB^{TL}* and *CFB^{TR}* (Supplementary Fig. 2a); in pin,
113 a single *CFB^P* exists. Gene model predictions also identified seven genes in the 75 kb to the left of
114 *CFB^{TL}*, and eight genes in the 96 kb to the right of *CFB^{TR}*, designated S Flanking Gene Left (*SFG^L*) and
115 Right (*SFG^R*) (Supplementary Fig. 2a); these genes are present in both pin and thrum.

116 To further investigate *S* haplotype differences, we aligned thrum and the short homostyle genome
117 contigs to the 455 kb *S^{LHI}* region (Supplementary Fig. 3a, b); although *S* and *S^{SHI}* assemblies are not
118 contiguous, they show homology across the 278 kb region. We also aligned genome sequence reads
119 from pin, thrum, long²⁵ and short³⁵ homostyle to the *S^{LHI}* assembly and plotted sequence read depth
120 across the 455 kb region (Fig. 2a). Sequences flanking the 278 kb insertion show a read depth of ~60
121 in all four genomes. However, between *CFB^{TL}* and *CFB^{TR}* we see differences; the long homostyle
122 (*S^{LHI}/S^{LHI}*) behaves as a homozygote, but both thrum (*S/s*) and the short homostyle (*S^{SHI}/s*) have half
123 this read depth, and pin (*s/s*) lacks this region (Fig. 2a); they behave genetically as heterozygotes but
124 our data show they are hemizygous for a region that is absent in pin. Alignment of all four genomes
125 over this region further show that thrum, long and short homostyles share the same boundary regions
126 (Figs. 1c, 2a, Supplementary Figs. 2, 3). This detail, coupled to the presence of all five *S* locus genes
127 in thrum, long and short homostyle (Fig. 2a) show that these homostyles did not arise by recombination
128 as proposed^{18,19}, and that *S*, *S^{LHI}* and *S^{SHI}* haplotypes all reside within an equivalent region that is absent
129 from pins.

130 To determine whether the 278 kb region is the only thrum-specific region in the genome, we searched
131 for additional thrum-specific genome sequences encoding genes. In two parallel analyses we identified
132 transcripts that were only expressed in thrums, and also mapped pin genomic sequencing reads to a
133 thrum genome assembly. We then examined the depth and breadth of pin genome reads mapped to the

134 thrum genome in the regions defined by the thrum-specific transcripts; *k*-means clustering analysis
135 resolved the transcribed regions into two clusters (Fig. 2b); deep and broad read coverage defined
136 presence of the region in both pin and thrum genomes, low read depth or low coverage identified a
137 region as thrum-specific, with pin sequence alignments representing erroneously mapped sequence
138 reads (Supplementary Table 3). Nine thrum-specific regions were thus identified; these define four of
139 the five thrum-specific genes from the 278 kb region; *GLO^T*, *CYP^T*, *PUM^T* and *KFB^T* (Fig. 2b). *CCM^T*
140 is expressed at a low level (see below) and is the only gene from the cluster not represented.
141 Identification of three contigs for *GLO^T* and *KFB^T*, and two for *CYP^T*, is due to the use of a non-
142 scaffolded thrum genome assembly (Supplementary Table 1b), and the length of *GLO^T* and *CYP^T* (see
143 below). We conclude that there are no other flower-expressed genes unique to thrums and that the 278
144 kb sequence is the only thrum-specific genomic region. Significantly, these data show that the thrum *S*
145 haplotype does not contain any additional genes compared to the long homostyle *S^{LHI}* haplotype. These
146 analyses revealed 391 gene models that are uniquely expressed in thrums, and 270 gene models that are
147 uniquely expressed in pins, but present in both pin and thrum genomes; these are candidates for direct
148 or indirect targets of the *S* locus genes that control pin and thrum flower development.

149 Next we investigated whether the sequences flanking the thrum-specific 278 kb region could also be
150 part of the *S* locus that contained pin- and thrum-specific alleles of genes involved in the control of
151 heterostyly. If sequences flanking the thrum-specific region contain genes that also contribute to *S*
152 locus function, restriction of recombination between pin and thrum alleles would be required to
153 maintain integrity and functionality of the locus. However, if these flanking regions are freely
154 recombining this would indicate that the thrum-specific region alone contains the entire *S* locus gene
155 cluster. We therefore undertook a recombination analysis investigating the pattern of nucleotide
156 polymorphisms (SNPs) across the flanking sequences, comparing the alleles present in a pin and a thrum
157 plant. These data (Supplementary Fig. 4) reveal that sequences flanking the thrum-specific region
158 contain blocks of significantly reduced polymorphism, which is consistent with recent recombination
159 events. The sequences flanking the thrum-specific 278 kb region thus seem to be homogenised by

160 recombination between pin and thrum alleles and suggest that they are not involved in the control of
161 the heterostyly phenotype.

162 **Linkage of *GLO^T* and the *S* locus**

163 *GLO^T* was initially identified as thrum-specific in a small segregating population³³. To demonstrate
164 unequivocal linkage of *GLO^T*, and therefore the *S^{LHI}* assembly, to the *S* locus, we revisited a three-point-
165 cross with 2075 progeny³⁵ used previously to place *Oakleaf*³¹ (<1.7 cM) and *Hose in Hose*⁴¹ (<1.6 cM)
166 on either side of the *S* locus³⁵ (Fig. 3a, Supplementary Table 4). This cross also yielded the short
167 homostyle *Hose in Hose* plant³⁵ used here (Fig. 1a). We analysed DNA from pin and thrum parents,
168 pools of pin and thrum non-recombinant progeny, and two double-recombinant (*Oakleaf-S-Hose in*
169 *Hose*) thrum progeny by PCR analysis with *GLO^T* and *GLO* specific primers (Fig. 3, Supplementary
170 Table 2); *GLO* is present in both pin and thrum³².

171 The parent plants show the original linkage profiles (Fig. 3b); we found no linkage disruption between
172 *GLO^T* and thrum phenotype using pools of 100 non-recombinant progeny. Furthermore, double-
173 recombinants show that recombination between *Oakleaf* and *S*, or *S* and *Hose in Hose*, does not disrupt
174 linkage between *GLO^T* and thrum phenotype (Fig. 3b). These data place the 455 kb assembly between
175 *Oakleaf* and *Hose in Hose*, within the *S* locus BAC assembly³⁵ (Fig. 3a); previous studies did not
176 identify any BACs that link the 455kb region to BAC contigs S-left and S-right³⁵. To increase mapping
177 resolution, we analysed natural populations of *P. vulgaris* and *P. veris*. Pooled genomic DNA from
178 200 pin plants of each species was analysed by PCR using *GLO^T* and *GLO* specific primers (Fig. 3c).
179 A single thrum plant was used as control because loss of a dominant marker in one individual would
180 not be detected in a thrum pool. In total, 500 pin plants were analysed (Figs. 3b,c), none showed
181 recombination; these data demonstrate that *GLO^T* and the surrounding region is in tight thrum-specific
182 linkage (<0.2 cM) with the *S* locus in both *P. vulgaris* and *P. veris* (Fig. 3a).

183

184

185 ***S* locus gene expression and function**

186 Having shown that homostyles did not occur by recombination we sought to determine their molecular
187 basis by comparing gene expression across the four haplotypes. Expression analysis of genes within,
188 and flanking, the *S* locus was undertaken by mapping four replicate RNA-Seq datasets from pin and
189 thrum flowers to the *S^{LH1}* assembly (Fig. 4a and Supplementary Table 5). *GLO^T*, *CYP^T*, *PUM^T*, *KFB^T*
190 and *CCM^T* all show thrum-specific expression (Figs. 2b, 4a). *CFB^{TL}* is expressed at a low level in both
191 pin and thrum flowers; *CFB^{TR}* is not expressed. Genes flanking the *S* locus are expressed in both pin
192 and thrum flowers, except *SFG^{R6}* which has low expression in thrum and is not detected in pin; *SFG^{L1}*
193 is expressed at a low level in both pin and thrum (Fig. 4a, Supplementary Table 5). These analyses
194 reveal the 278 kb region as an island of thrum-specific gene expression.

195 Gene model predictions (Supplementary Fig. 2a) for *CCM^T*, *GLO^T*, *CYP^T*, *PUM^T*, *KFB^T*, and three *CFB*
196 alleles were confirmed by alignment to RNA-Seq data to define intron-exon boundaries. Two *S* locus
197 genes are surprisingly large, *GLO^T* spans 25 kb with two introns over 10 kb; *CYP^T* spans 68 kb with 10,
198 20 and 30 kb introns (Supplementary Fig. 2c). Interestingly, the *GLO^T S^{SH1}* allele contains a 2.5 kb
199 retro-transposon in exon 2 which disrupts and severely truncates the encoded protein; mutation of *GLO^T*
200 in the short homostyle is associated with loss of anther elevation; style length and pollen size are
201 unaffected. The long homostyle (*S^{LH1}*) *CYP^T* allele has a single base insertion in exon 3 that introduces
202 a disruptive premature stop codon, and is associated with loss of style length suppression; anther height
203 and pollen size are unaffected. We also sequenced an independent long homostyle (*S^{LH2}*) from the
204 Chiltern Hills²⁵ which represents a second *CYP^T* mutant allele with a G-C transversion in exon 2 that
205 results in an Asp126His substitution. *CFB^{TR}* has an 11 bp deletion compared to *CFB^{TL}* and *CFB^P* that
206 introduces a premature stop codon. The architecture of *S* and *s* haplotypes is summarised in Fig. 4b.
207 Comparison of alleles is presented in Supplementary Sequence Analysis 3.

208 **Date of the *GLO^T* duplication**

209 The first indication that *GLO^T* was a discrete locus from *GLO* came when we identified distinct BAC
210 clones for each gene, together with insight from other studies of B function MADS box genes which
211 suggested duplication could underpin diversification of novel floral morphologies⁴². The short
212 homostyle *GLO^T* mutation is not complemented by *GLO*, or by ectopic expression of *GLO*; the short

213 homostyle is in the *Hose in Hose* background³². The recent report of a partial *P. veris* genome
214 sequence³⁶ noted the duplication of *GLO* and referred to the genes as *GLO1* and *GLO2* but could not
215 show linkage of *GLO^T* (*GLO2*) to the *S* locus. Demonstration that these genes represent distinct loci
216 with *GLO^T* at the *S* locus provides the opportunity to date the duplication event associated with assembly
217 of the *S* locus supergene. To determine the age of duplication we isolated *GLO* and *GLO^T* sequences
218 from six *Primula* species, and used these with sequences from other species to conduct a Bayesian
219 relaxed-clock phylogenetic analysis with a combination of secondary calibrations (Fig. 5,
220 Supplementary Tables 6a,b). The index of substitution saturation value⁴³ for *GLO* and *GLO^T* sequences
221 (0.1187) was significantly lower than the Iss critical value (0.7318, $p < 0.0001$) indicating low saturation
222 between these sequences. These analyses yielded a mean (5-95% Highest Posterior Density) age
223 estimate of 51.7 (33.1-72.1) MYA for the duplication leading to the divergence of *GLO* and *GLO^T*
224 lineages.

225 The duplication and neofunctionalisation of *GLO^T* represents a landmark evolutionary event at the *S*
226 locus, and precedes estimates for the *Primula*-*Androsace* divergence; estimates for this node are 32 (20-
227 51) MYA⁴⁴, and 44 (33-54) MYA⁴⁴ with fossil priors being set with a log normal distribution, and 40
228 (30-51) MYA with fossils modelled as exponential priors⁴⁵. The *Androsace* were predicted to be the
229 first taxon within the *Primulaceae* to exhibit heteromorphy⁴⁶, our data indicate that the *GLO-GLO^T*
230 duplication predates this divergence, which implies heterostyly evolved following a single duplication
231 event in the *Primulaceae*. Two models have been proposed for the evolution of *Primula* heterostyly,
232 the first postulates a long homostyle²¹, and the other an approach herkogamous pin-form flower²⁰, as
233 the original floral form. The duplication and neofunctionalization of *GLO^T* would be consistent with
234 both models if this was the first gene at the (Fig. 5). The *S* locus sequence, structure and timing of the
235 *GLO^T* duplication, and analysis of other genes at the *S* locus genes, will inform further evolutionary
236 genetic analysis of primary and secondary homostyly in *Primula* and help to determine the sequence of
237 events leading to the establishment of the *S* locus gene cluster.

238 **Conclusions**

239 We show that the *S* locus supergene is a tightly-linked cluster of five thrum-specific genes, spanning a
240 278 kb sequence that is absent in pins (Fig. 2a), this finding defines the basis for Bateson and Gregory's
241 *S* haplotype dominance⁵. The annotation *S/s* and *s/s* for thrum and pin could be represented by *S/-* and
242 *-/-*, but we suggest retention of the traditional nomenclature with recognition of *s* as a null haplotype.
243 Floral heteromorphy in *Primula* has evolved after duplication of a floral homeotic gene 51.7 MYA,
244 followed by its neofunctionalisation, creating the current *S* locus assemblage. This insight has profound
245 implications for our understanding of a key evolutionary innovation of flowering plants. The molecular
246 basis of the *Primula S* locus supergene appears to be different from those proposed for the control of
247 butterfly mimicry, and avian and insect social behaviour⁴⁷⁻⁴⁹. It is also unlike the mating-type locus in
248 ascomycete fungi which comprises two distinct idiomorphs⁵⁰. Ernst originally proposed that *Primula*
249 homostyles arose by mutation⁸, he was correct, and mutations in *CYP^T* and *GLO^T* homostyle alleles
250 earmark these genes as candidates for the style length suppression (*G*), and anther elevation (*A*),
251 functions⁷ respectively. Darwin suggested the primary function of heterostyly evolved to promote out-
252 crossing¹³, generating novel variation that is the substrate of natural selection. The parallel evolution
253 of heterostyly in diverse angiosperm families¹ has exploited insect-mediated pollination, which in turn
254 is associated with an accelerated rate of speciation in angiosperms¹⁴. Deciphering the genetic
255 architecture of the *Primula S* locus as the first heterostyly supergene provides a blueprint for the
256 comparative evolutionary genetic analysis of this key adaptation in other angiosperm families, as well
257 as for the molecular characterisation of other pollination syndromes underpinning both biodiversity and
258 food security.

259

260 **Methods**

261 **Plant Material**

262 The long homostyle plant (*S^{LHI}*) used for DNA sequencing was a homozygote derived from a population
263 originally described by Crosby in 1940²⁵ at Wyke Champflower, Somerset, UK, which had undergone
264 several generations of selfing to generate a homozygous line which greatly facilitated assembly of the

265 genome sequence. The independent long homostyle population in the Chiltern Hills discovered by
266 Crosby in 1944²⁶ provided our second long homostyle (S^{LH2}) from Hawridge, Buckinghamshire, UK.
267 Pin and thrum *P. vulgaris* were grown from seed (<http://www.wildseed.co.uk>) as described previously²⁷.
268 Pin and thrum plants selected for genome sequencing were crossed to generate an F1 population. The
269 short homostyle was originally identified in a mapping population of *P. vulgaris* plants³⁵ *P. veris* for
270 genome sequencing were grown from seed collected at the Durham University Mountjoy site. *P. elatior*
271 leaf material was collected from Bull's Wood (<http://www.suffolkwildlifetrust.org>) with permission of
272 Suffolk Wildlife Trust. *P. farinosa* were obtained from Kevock Garden Plants
273 (<http://www.kevockgarden.co.uk/>), *P. vialii* and *P. denticulata* were from the laboratory collection. The
274 population of *P. veris* used for *S* locus linkage analysis was sampled from Lolly Moor
275 (<http://www.norfolkwildlifetrust.org.uk>) with permission of Norfolk Wildlife Trust. *P. vulgaris* used
276 for *S* locus linkage analysis was sampled with permission of Norfolk County Council from the B1135
277 roadside verge between Ketteringham and Browick, near Wymondham, Norfolk. Plants from the three-
278 point mapping cross have been described previously³⁵.

279 **Preparation of sequencing libraries**

280 DNA and RNA preparation was as described previously^{30,33}. All genomic DNA and RNA-Seq libraries
281 for sequencing were prepared at The Genome Analysis Centre using standard Illumina protocols.
282 Genomic paired-end libraries: An Illumina TruSeq library was prepared using a protocol optimized for
283 1µg of input genomic DNA (Illumina 15026486 Rev. C). Mate-pair libraries: The protocol was
284 optimized for 4-10 µg of high molecular weight DNA; following fragmentation, samples were size
285 fractionated to enable generation of mate-pair libraries of 5, 7 and 9 kb (Illumina 15035209 Rev. D).
286 RNA paired-end libraries: Libraries were constructed using the Illumina TruSeq RNA protocol
287 (Illumina 15026495 Rev.B). The *S* locus region was assembled as outlined in Supplementary Methods
288 and Supplementary Fig. 1. The assembly was validated by comparison of independent assemblies from
289 different Illumina paired-end and and mate-paired sequencing libraries of thrum, long and short
290 homostyle individuals, gaps between contigs and regions of Ns were resolved by PCR amplification
291 and Sanger sequencing of the products.

292 **Genomic DNA PCR analysis**

293 Genomic DNA was isolated as described previously³⁰. Primers are shown in Supplementary Table 2.
294 PCR was performed in 50 µl reactions with 100-200 ng genomic DNA using Promega GoTaq G2 Green
295 Master mix (cat# M7822). *Pfu*Turbo DNA polymerase (Agilent Technologies, cat#600250) was used
296 when PCR products were to be sequenced. Amplification conditions for primers TRB-F and -R and
297 CFB-F and -R: 95°C 3 min; 95°C 30 sec, 61°C 30 sec, 72°C 1 min x35 cycles; 72°C 5 min.
298 Amplification cycles for primers PRB-F and -R, TLB-F and -R, GLOT-F and -R: 95°C 3 min; 95°C
299 30 sec, 58°C 30 sec, 72°C 1 min, x35 cycles : 72°C 5 min. Amplification cycles for primers GLO-F
300 and -R: 95°C 3 min; 95°C 30 sec, 65°C 30 sec, 72°C 1 min, x35 cycles; 72°C 5 min.

301 **Bioinformatic and evolutionary analyses**

302 Bioinformatic and evolutionary analyses are described in Supplementary Methods due to space
303 constraints.

304 **Data availability**

305 Sequence data are available through Genbank accessions KT257663-KT257681, under Bioproject
306 PRJEB9683 <http://www.ebi.ac.uk/ena/data/view/PRJEB9683>, and <http://opendata.earlham.ac.uk/primula/>

307

308 Correspondence and requests for materials should be addressed to P.M.G. (p.gilmartin@uea.ac.uk).

309

310 **Acknowledgements** We thank Martin Lappage, Mike Hughes and Pam Wells for horticultural support;
311 colleagues at TGAC for Illumina sequencing; Anil Thanki for TGAC Browser support; Olivia Kent for
312 *P. elatior* *GLO* and *GLO^T* sequences; Norfolk Wildlife Trust, Suffolk Wildlife Trust and Norfolk
313 County Council for permission to sample *P. veris*, *P. elatior* and *P. vulgaris* respectively; Matt Gage,
314 Brendan Davies and Dianna Bowles for comments on the manuscript; Wenjia Wang for advice on *k*-
315 means analysis; BBSRC for funding via grant BB/H019278/2, and prior awards G11027 and P11021;
316 The Gatsby Foundation for early stage funding; University of Leeds, Durham University and University

317 of East Anglia for support over several years of the project to PMG. PMG's lab is hosted at the John
318 Innes Centre under the UEA-JIC Norwich Research Park collaboration.

319 **Author Contributions** J.L. contributed to project design, performed all molecular analyses, generated
320 the *S* locus assembly, manually annotated *S* locus gene structures, and undertook data analysis. J.M.C.
321 contributed to bioinformatic analyses, including automated annotation of the *S* locus region, undertook
322 *in silico* gene expression and *k*-means clustering analyses, assembled genome sequences and library
323 scaffolds, generated the molecular phylogeny, undertook recombination analysis of the *S* locus flanking
324 regions and contributed to project design. J.W. assembled genome sequences and library scaffolds,
325 contributed to genome annotation and generated the automated gene model predictions across the *S*
326 locus, aligned sequence reads to the *S* locus assembly and contributed to project design. M.A.W.
327 contributed the inbred long homostyle line, other genetic resources and classical genetics, identified the
328 short homostyle mutant, and generated the three-point cross used to demonstrate linkage. M.M. and
329 C.v.O. contributed to the molecular phylogeny construction, evolutionary data analysis and
330 recombination analysis. S.A., D.S. and M.C. contributed to the genome sequencing strategy, assembly
331 and annotation that underpins this project. P.M.G. conceived, designed and directed the project,
332 contributed to data analysis, prepared the figures and drafted the manuscript, with revision input from
333 C.v.O.; all authors contributed to editing the manuscript.

334 **Figure Legends**

335 **Figure 1 *P. vulgaris* floral phenotypes and genotypes.**

336 **a**, Heterostyly phenotypes and genotypes with respect to *s*, *S*, *S^{LHI}* and *S^{SHI}* haplotypes; the short
337 homostyle carries the *Hose in Hose* mutation³². Anther (A) and stigma (S) **b**, Comparison of *s*, *S*, *S^{LHI}*
338 and *S^{SHI}* haplotypes, sequence present in the *S^{LHI}*, *S^{SHI}* and *S*, but absent from *s* (red); duplicated flanking
339 sequence present as a single copy in the *s* haplotype (yellow); flanking sequences common to all
340 haplotypes left (blue), and right (green); not to scale. The *GLO^T* BAC location is shown. PCR primers
341 used for amplification of flanking regions (→←) (Supplementary Table 2). **c**, PCR analysis of genomic
342 DNA from pin (P), thrum (T), long (LH) and short homostyle (SH) plants (shown in **a**), using primers

343 (as in **b**), that distinguish left (LB) and right (RB) borders of S^{LH1} , S^{SH1} and S haplotypes from the s
344 haplotype; sizes as indicated. See also Supplementary Sequence Analysis 1.

345 **Figure 2 Organisation of S locus haplotypes.** **a**, The S^{LH1} haplotype showing S locus genes (red)
346 duplicated flanking CFB loci (yellow), left and right flanking genes SFG^L 1-7 and SFG^R 1-8 (black)
347 (see also Supplementary Fig. 2a). Illumina sequence read depth from pin (black), thrum (blue), short
348 (red) and long homostyle (yellow) genomic DNA. **b**, Scatter plot analysis showing breadth of read
349 coverage (%) and \log_{10} depth of read coverage for pin progeny pool genome sequence reads mapped to
350 thrum genome contigs encoding genes with thrum-specific expression. Two clusters, defined by k -
351 means analysis are shown; transcript regions in contigs where reads map with low depth and breadth
352 (red): 1, KFB^T ; 2, GLO^T ; 3, CYP^T ; 4, GLO^T ; 5, KFB^T ; 7, PUM^T ; 8, GLO^T ; 8, CYP^T ; 9, KFB^T
353 (Supplementary Table 3). Transcript regions in contigs to which pin progeny pool genome sequence
354 reads map with high depth and breadth (blue).

355 **Figure 3 Linkage of the S haplotype to the thrum phenotype.** **a**, Map of the S locus region, distances
356 in cM³⁵. The S haplotype 278 kb region (red) between duplicated ~3 kb CFB loci (yellow) is shown
357 relative to sequenced BAC contigs³⁵ with 75 kb left flanking (blue) and 96 kb right flanking (green)
358 sequences (Fig. 1b). **b**, PCR analysis of GLO^T linkage using Pin (P) and thrum (T) plants and 100
359 pooled non-recombinant (no x-over) progeny from a three-point cross³⁵ (Supplementary Table 4)
360 compared to non- S locus GLO as control. Two thrum plants (T1 and T2) from double-recombination
361 events (xx-over)³⁵, *Oakleaf* to S and S to *Hose in Hose*, are also shown. **c**, PCR analyses of GLO^T
362 linkage to thrum in natural populations using 200 pooled pin (P) *P. vulgaris*, and 200 pooled pin *P.*
363 *veris* plants, compared to individual thrum (T) plants, with GLO as control; sizes in kb.

364 **Figure 4 Expression and genomic organisation of S locus genes.** **a**, Gene expression from the S
365 (red) and s (blue) haplotypes using pin and thrum RNA-Seq data represented as Log_{10} of the number of
366 fragments per kb of transcripts per million fragments mapped (FPKM) +1; gene models as defined in
367 Supplementary Fig. 2a. **b**, Pictorial representation of genes within the s , S , S^{LH1} and S^{SH1} haplotypes
368 shown alongside stylized flowers. The base insertion in S^{LH1} CYP^T (red +), G-C transversion in S^{LH2}

369 (red I) and transposon insertion in $S^{SHI} GLO^T$ (red Δ) are indicated. Sequences of mutant alleles are
370 compared in Supplementary Sequence Analysis 3.

371 **Figure 5 Phylogenetic analysis and the date of duplication of GLO^T from GLO .** Phylogram of B
372 function MADS box genes from *Antirrhinum* (Am.), *Petunia* (Pe.), *Arabidopsis* (A.) and *Primula* (P.)
373 species presented against an evolutionary time scale in millions of years (MYA); see Supplementary
374 Tables 6a,b. Thick blue lines represent the time scale range estimates at divergence branch points; the
375 thick red line defines the same for duplication of GLO^T from GLO .

376

377 References

- 378 1 Barrett, S. C. H. The evolution of plant sexual diversity. *Nature Reviews, Genetics* **3**, 274-284
379 (2002).
- 380 2 Richards, A. J. *Primula 2nd edition*. (2002).
- 381 3 Darwin, C. R. On the two forms or dimorphic condition in the species of *Primula*, and on their
382 remarkable sexual relations. *Journal of the Proceedings of the Linnean Society, Botany* **6**, 77-
383 96 (1862).
- 384 4 Gregory, R. P., De Winton, D. & Bateson, M. A. Genetics of *Primula sinensis*. *Journal of*
385 *Genetics* **13**, 219-253 (1923).
- 386 5 Bateson, W. & Gregory, R. P. On the inheritance of heterostylism in *Primula*. *Proceedings of*
387 *the Royal Society of London B Series* **76**, 581-586 (1905).
- 388 6 Bridges, C. B. The chromosome hypothesis of linkage applied to cases in sweetpeas and
389 *Primula*. *American Naturalist* **48**, 524-534 (1914).
- 390 7 Ernst, A. Weitere untersuchungen zur Phänanalyse zum Fertilitätsproblem und zur Genetik
391 heterostyler Primeln. II. *Primula hortensis*. *Archive der Julius Klaus Stiftung für*
392 *Vererbungsforschung Sozialanthropologie und Rassenhygiene* **11**, 1-280 (1936).
- 393 8 Ernst, A. Heterostylie-Forschung Versuche zur genetischen analyse eines organisations und
394 'Anpassungs' merkmals. *Zeitschrift für Induktive Abstammungs und Vererbungslehre* **71**, 156-
395 230 (1936).
- 396 9 De Winton, D. & Haldane, J. B. S. The genetics of *Primula sinensis*. III. Linkage in the diploid.
397 *Journal of Genetics* **31**, 67-100 (1935).
- 398 10 Darlington, C. D. Meiosis in diploid and tetraploid *Primula sinensis*. *Journal of genetics* **24**,
399 65-95 (1931).
- 400 11 Mather, K. The genetical architecture of heterostyly in *Primula sinensis*. *Evolution* **4**, 340-352
401 (1950).
- 402 12 Schwander, T., Libbrecht, R. & Keller, L. Supergenes and Complex Phenotypes. *Current*
403 *Biology* **24**, R288-R294 (2014).
- 404 13 Darwin, C. R. *The different forms of flowers on plants of the same species.*, (John Murray,
405 1877).
- 406 14 Dodd, M. E., Silvertown, J. & Chase, M. W. Phylogenetic analysis of trait evolution and species
407 diversity variation among angiosperm families. *Evolution* **53**, 732-744 (1999).
- 408 15 Webster, M. A. & Gilmartin, P. M. Analysis of late stage flower development in *Primula*
409 *vulgaris* reveals novel differences in cell morphology and temporal aspects of floral
410 heteromorphy. *New Phytologist* **171**, 591-603 (2006).

- 411 16 Shivanna, K. R., Heslop-Harrison, J. & Heslop-Harrison, Y. Heterostyly in *Primula* .2. Sites of
412 pollen inhibition, and effects of pistil constituents on compatible and incompatible pollen tube
413 growth. *Protoplasma* **107**, 319-337 (1981).
- 414 17 Richards, A. J. & Ibrahim, H. B. The Breeding System in *Primula veris* L .2. Pollen-Tube
415 Growth and Seed-Set. *New Phytologist* **90**, 305-314 (1982).
- 416 18 Lewis, D. Comparative incompatibility in angiosperms and fungi *Advances in Genetics*
417 *Incorporating Molecular Genetic Medicine* **6**, 235-285 (1954).
- 418 19 Dowrick, V. P. J. Heterostyly and homostyly in *Primula obconica*. *Heredity* **10**, 219-236
419 (1956).
- 420 20 Lloyd, D. G. & Webb, C. J. in *Evolution and Function of Heterostyly* (ed S.C.H. Barrett) 151-
421 175 (Springer Verlag, 1992).
- 422 21 Charlesworth, D. & Charlesworth, B. Model for the evolution of distyly. *American Naturalist*
423 **114**, 467-498 (1979).
- 424 22 Bodmer, W. F. The genetics of homostyly in populations of *Primula vulgaris*. *Philosophical*
425 *Transactions of the Royal Society of London Series B-Biological Sciences* **242**, 517-549 (1960).
- 426 23 Fisher, R. A. A Theoretical system of selection for homostyle *Primula*. *Sankhya* **9**, 325-342
427 (1949).
- 428 24 Piper, J. G., Charlesworth, B. & Charlesworth, D. A high-rate of self-fertilization and increased
429 seed fertility of homostyle primroses. *Nature* **310**, 50-51 (1984).
- 430 25 Crosby, J. L. High proportions of homostyle plants in populations of *Primula vulgaris*. *Nature*
431 **145**, 672-673 (1940).
- 432 26 Crosby, J. L. Selection of an unfavourable gene complex. *Evolutionary Ecology Research* **3**,
433 212-230 (1949).
- 434 27 Webster, M. A. & Gilmartin, P. M. A comparison of early floral ontogeny in wild-type and
435 floral homeotic mutant phenotypes of *Primula*. *Planta* **216**, 903-917 (2003).
- 436 28 McCubbin, A. G., Lee, C. & Hetrick, A. Identification of genes showing differential expression
437 between morphs in developing flowers of *Primula vulgaris*. *Sexual Plant Reproduction* **19**, 63-
438 72 (2006).
- 439 29 Li, J., Webster, M. A., Furuya, M. & Gilmartin, P. M. Identification and characterization of pin
440 and thrum alleles of two genes that co-segregate with the *Primula S* locus. *Plant Journal* **51**,
441 18-31 (2007).
- 442 30 Manfield, I. W. *et al.* Molecular characterization of DNA sequences from the *Primula vulgaris*
443 *S* locus. *Journal of Experimental Botany* **56**, 1177-1188 (2005).
- 444 31 Cocker, J. *et al.* *Oakleaf*: an *S* locus-linked mutation of *Primula vulgaris* that affects leaf and
445 flower development *New Phytologist* **10.1111/nph.13370** (2015).
- 446 32 Li, J. *et al.* *Hose in Hose*, an *S* locus-linked mutant of *Primula vulgaris* is caused by an unstable
447 mutation at the *Globosa* locus. *PNAS* **107**, 5664-5668 (2010).
- 448 33 Li, J. *et al.* The *S* locus-linked *Primula* homeotic mutant *sepaloid* shows characteristics of a B-
449 function mutant but does not result from mutation in a B-function gene. *Plant Journal* **56**, 1-12
450 (2008).
- 451 34 Yoshida, Y. *et al.* QTL analysis of heterostyly in *Primula sieboldii* and its application for morph
452 identification in wild populations. *Annals of Botany* **108**, 133-142 (2011).
- 453 35 Li, J. *et al.* Integration of genetic and physical maps of the *Primula vulgaris S* locus and
454 localization by chromosome *in situ* hybridisation *New Phytologist* **10.1111/nph.13373** (2015).
- 455 36 Nowak, M. D. *et al.* The draft genome of *Primula veris* yields insight into the molecular basis
456 of heterostyly. *Genome Biology* **16**, 16 (2015).
- 457 37 Verhoef, N. *et al.* Brassinosteroid biosynthesis and signalling in *Petunia hybrida*. *Journal of*
458 *Experimental Botany* **64**, 2435-2448 (2013).
- 459 38 Turk, E. M. *et al.* CYP72B1 inactivates brassinosteroid hormones: An intersection between
460 photomorphogenesis and plant steroid signal transduction. *Plant Physiology* **133**, 1643-1653
461 (2003).
- 462 39 Abbasi, N., Park, Y.-I. & Choi, S.-B. Pumilio Puf domain RNA-binding proteins in
463 *Arabidopsis*. *Plant signaling & behavior* **6**, 364-368 (2011).

464 40 Kim, H. J., Chiang, Y.-H., Kieber, J. J. & Schaller, G. E. SCFKMD controls cytokinin signaling
465 by regulating the degradation of type-B response regulators. *Proceedings of the National*
466 *Academy of Sciences of the United States of America* **110**, 10028-10033 (2013).

467 41 Webster, M. A. & Grant, C. J. The inheritance of calyx morph variants in *Primula vulgaris*
468 (Huds). *Heredity* **64**, 121-124 (1990).

469 42 Viaene, T. *et al.* Pistillata-Duplications as a Mode for Floral Diversification in (Basal) Asterids.
470 *Molecular Biology and Evolution* **26**, 2627-2645 (2009).

471 43 Xia, X. DAMBE5: A Comprehensive Software Package for Data Analysis in Molecular
472 Biology and Evolution. *Molecular Biology and Evolution* **30**, 1720-1728 (2013).

473 44 Magallón, S., Gómez-Acevedo, S., Sánchez-Reyes, L. L. & Hernández-Hernández, T. A
474 metacalibrated time-tree documents the early rise of flowering plant phylogenetic diversity.
475 *New Phytologist*, 10.1111/nph.13264 (2015).

476 45 Bell, C. D., Soltis, D. E. & Soltis, P. S. The age and diversification of the angiosperms re-
477 revisited. *American Journal of Botany* **97**, 1296-1303 (2010).

478 46 Mast, A. R. *et al.* Phylogenetic relationships in *Primula* L. and related genera (Primulaceae)
479 based on noncoding chloroplast DNA. *International Journal of Plant Sciences* **162**, 1381-1400
480 (2001).

481 47 Joron, M. *et al.* Chromosomal rearrangements maintain a polymorphic supergene controlling
482 butterfly mimicry. *Nature* **477**, 203-206 (2011).

483 48 Thomas, J. W. *et al.* The chromosomal polymorphism linked to variation in social behavior in
484 the white-throated sparrow (*Zonotrichia albicollis*) is a complex rearrangement and suppressor
485 of recombination. *Genetics* **179**, 1455-1468 (2008).

486 49 Wang, J. *et al.* A Y-like social chromosome causes alternative colony organization in fire ants.
487 *Nature* **493**, 664-668 (2013).

488 50 Turgeon, B. G. & Yoder, O. C. Proposed Nomenclature for Mating Type Genes of Filamentous
489 Ascomycetes. *Fungal Genetics and Biology* **31**, 1-5 (2000).

490