

A peer-reviewed version of this preprint was published in PeerJ on 20 May 2014.

[View the peer-reviewed version](https://doi.org/10.7717/peerj.391) (peerj.com/articles/391), which is the preferred citable publication unless you specifically need to cite this preprint.

Bombarely A, Coate JE, Doyle JJ. 2014. Mining transcriptomic data to study the origins and evolution of a plant allopolyploid complex. PeerJ 2:e391 <https://doi.org/10.7717/peerj.391>

Mining transcriptomic data to study the origins and evolution of a plant allopolyploid complex

Allopolyploidy combines two progenitor genomes in the same nucleus, and is a common mechanism for producing new species, especially in plants. Deciphering the origins of polyploid species is a complex problem, due to, among other things, extinct progenitors, multiple origins, gene flow between different polyploid populations, and loss of parental contributions through gene or chromosome loss. In this work, we studied three allopolyploid species in the genus *Glycine*, which includes the cultivated soybean (*G. max*). Previous work based on two nuclear sequences showed that these allopolyploids combine the genomes of extant diploid species in the *G. tomentella* complex. We use several phylogenetic and population genomics approaches to clarify the origin of these species using single nucleotide polymorphism data and a guided transcriptome assembly. The results support the hypothesis that each of the three polyploid species are fixed hybrids combining the homoeologous genomes of its two putative parents. Based on mapping to the soybean reference genome, there appear to be no large regions for which one homoeologous contribution is missing. Phylogenetic analyses of 27 selected transcripts using a coalescent approach also indicates multiple origins for *G. tomentella* polyploid species, and suggest that origins occurred within the last several hundred thousands years.

3 1 - Department of Plant Biology, Cornell University, Ithaca, New York 14853 USA; 2 - Department of Biology, Reed
4 College, Portland, Oregon 97202-8199 USA;

5 INTRODUCTION

6 Polyploidy (whole genome duplication, WGD) is a key process in plant evolution. All seed plants are
7 fundamentally polyploidy, with a second WGD event characterizing all flowering plants (Soltis et al., 2009; Jiao et
8 al., 2011), and additional events found in many lineages (Jiao et al., 2011). It has been estimated that 15% of all
9 flowering plant speciation events involve polyploidy (Wood et al., 2009). Systematists generally recognize
10 autopolyploidy and allopolyploidy as distinct types of polyploidy events, based on the level of divergence of the
11 diploid genomes that formed the polyploid. The terms are best thought of as describing elements of a continuum that
12 ranges from the doubling of a single genome (autopolyploidy), to the incorporation of differentiated genomes in a
13 single nucleus by hybridization (allopolyploidy). From a genetic perspective, allopolyploids are characterized by
14 diploid-like meiotic behavior and limited interaction between the two homoeologous genomes. The duplicated
15 chromosomes of an autopolyploid initially can interact randomly, leading to polysomic segregation, but it is
16 generally assumed that this is a transient state; diploidization leads to the eventual presence of homoeologous
17 genomes. It is difficult, if not impossible to determine from the genomes of older polyploids (paleopolyploids,
18 mesopolyploids) how differentiated their progenitor genomes were.

19 The initial “fixed hybrid” condition of an allopolyploid erodes over time as homoeologous loci are lost
20 (Lynch and Conery, 2000; Maere et al., 2005); this process of “fractionation” is thought to occur preferentially from
21 one subgenome, but the precise mechanisms remain unknown (Freeling et al., 2012). The earliest stages of polyploid
22 evolution may contribute disproportionately to gene loss and genomic rearrangement through genomic shock
23 (McClintock, 1984). For example, individuals of the ca. 100 year-old allopolyploid, *Tragopogon miscellus*, have lost
24 entire chromosomes of one parent (Chester et al., 2012). Diversity in polyploids can be due to mutational divergence
25 from parental diploids, but also due to multiple origin (Symonds et al., 2010).

26 The ready availability of genomic and transcriptomic data has opened new opportunities for studying the
27 evolution of polyploids (Grover et al., 2012; Ilut et al., 2012; Dufresne et al., 2014) at the scale of whole genomes,
28 though the application of such data to complex polyploid genomes is challenging. At the same time, the field of
29 systematics has seen what has been called a new paradigm for studying species relationships, involving genealogical
30 approaches (Edwards, 2009). Genealogical methods have lately begun to be applied to both autopolyploids (Arnold
31 et al., 2012; Hollister et al., 2012) and allopolyploids (e.g. Slotte et al., 2011; Jones et al., 2013; Slotte et al., 2013).

32 The genus *Glycine* includes the cultivated soybean (*G. max*) and its wild progenitor (*G. soja*), both annual
33 species native to northeastern Asia, as well as approximately 30 perennial species native to Australia classified as
34 subgenus *Glycine* (Ratnaparkhe et al., 2010). Like many plant species, *Glycine* has a complex history of polyploidy:
35 in addition to events shared with all angiosperms (Jiao et al., 2011) and eudicots (Jiao et al., 2012), the soybean
36 genome retains evidence from a whole genome duplication (WGD) around 50 million years ago (MYA) shared with a
37 large subset of legumes (Blanc and Wolfe, 2004; Schlueter et al., 2004; Cannon et al., 2010), and particularly from a
38 more recent polyploidy event that increased the chromosome number from $2n = 20$ to $2n = 40$ (Shoemaker et al.,
39 2006; Doyle and Egan, 2010; Schmutz et al., 2010; Doyle, 2012). This *Glycine*-specific WGD occurred around 10
40 MYA, which is the estimated time of homoeologous gene divergence (e.g., Egan and Doyle, 2010; Schmutz et al.,
41 2010), and 5 MYA, when the annual and perennial species diverged from a common already-polyploid ancestor
42 (Doyle and Egan, 2010).

43 In addition to these older events, eight perennial *Glycine* species are allopolyploids with $2n = 78$ or 80
44 hypothesized to have arisen by hybridization involving various combinations of eight extant diploid species, several
45 multiple times and involving both progenitors as chloroplast genome donors (Doyle et al., 2004). Their origins were
46 hypothesized initially from crossing data and more recently from gene phylogenies, but inferences have been made
47 from only two nuclear genes—histone H3-D (Doyle et al., 1999; Doyle et al., 2002) and the 18S-26S nuclear
48 ribosomal cistron internal transcribed spacer region (nrDNA ITS: (Singh et al., 2001; Rauscher et al., 2004)). Both of
49 these markers confirmed the fixed hybridity of *Glycine* allopolyploid species, but it is not known to what extent the

entire genomes of these plants retains contributions from both parental diploid species. Furthermore, an estimate of the date of origin has been made for only one of the eight allopolyploids (Doyle et al., 1999).

A better understanding of the origin and evolution of the *Glycine* allopolyploid complex will complement its exploitation in studying the impact of allopolyploidy on a range of morphological and physiological characters (Coate and Doyle, 2010; Coate et al., 2012; Ilut et al., 2012; Coate et al., 2013; Hegarty et al., 2013). Here we apply phylogenetic and coalescent methods to a transcriptomic dataset from three of these allopolyploid species and their diploid progenitors that was originally generated to study the effects of polyploidy on their ability to cope with stress from excess light (Coate et al., 2013).

MATERIAL AND METHODS

Taxon Sampling and Transcriptome Sequencing

Three *Glycine* (Fig. 1) allopolyploid “triads,” defined as an allopolyploid species and its two putative diploid progenitors, were sampled : 1) the allopolyploid, *G. tomentella* T1 ($2n = 78$) and the diploid species, *G. tomentella* D1 (E-genome of (Hymowitz et al., 2010); $2n = 38$) and *G. tomentella* D3 (D-genome; $2n = 40$); 2) *G. dolichocarpa* (= *G. tomentella* T2; $2n = 80$) and its putative progenitors *G. tomentella* D3 and *G. syndetika* (= *G. tomentella* D4; $2n = 40$); and 3) *G. tomentella* T5 ($2n = 78$) and its hypothesized progenitors, *G. tomentella* D1 and *G. clandestina* ($2n = 40$). Each species was represented by 2-5 accessions sampled from the CSIRO Division of Plant Industry Perennial *Glycine* Germplasm Collection (Table 1).

Plants were grown in a common growth chamber with a 12 h/12 h light/dark cycle, 22 °C/18 °C day/night temperature regime, and a light intensity of either 125 mmol m⁻² s⁻¹ (LL) or 800 mmol m⁻² s⁻¹ (EL). Different light intensities were used for the purposes of a separate study examining light stress responses (Coate et al., 2013). Single leaflets were pooled from six individuals per accession, and RNA-Seq libraries were constructed from the pooled tissue. All samples were taken from approximately 1-week-old, fully expanded leaves, and were collected 0.5–2.0 h into the light period. For each light treatment, all tissue was collected in a single morning and immediately frozen in liquid nitrogen. Total RNA was isolated from pooled leaf tissue using the Plant RNeasy Kit with on-column DNase treatment (Qiagen, Valencia, CA, USA). Single-end RNA-Seq libraries were constructed following the Illumina mRNA-seq Sample Preparation Kit protocol (Illumina, San Diego, CA, USA), with the following modifications: (1) two rounds of polyA selection were performed using the Dynabeads mRNA DIRECT Kit (Life Technologies, Carlsbad, CA, USA); (2) RNA was fragmented for 2 min at 70 °C using the RNA fragmentation reagents kit (Life Technologies; and (3) Illumina PE adapters were replaced with custom-made adapters containing 3nt barcodes in order to facilitate multiplexing of samples (see (Coate et al., 2013) for adapters and Supplementary Table S1 for the barcodes sequences). Sequencing was performed on either the GAIIx or HiSeq 2000 platform (Illumina), generating 88 nt or 100 nt reads, respectively. Equimolar amounts of three (GAIIx) or four (HiSeq 2000) barcoded libraries were combined and sequenced per channel.

Read Processing and Single Nucleotide Polymorphism (SNP) Calling

Perennial species reads were processed with Fastq-mcf (Aronesty, 2013) to trim low quality extremes (min. quality 30) and remove short reads (min. read length 50 bp). They were aligned to the soybean genome (version 1.0, downloaded from www.phytozome.net/soybean) using Bowtie2 (Langmead and Salzberg, 2012) with the default

parameters. Mapping files from the same accession were merged. Reads without preferential mapping (same score for two or more mapping hits) and with a mapping score below 20 were removed. SNP calling was performed using Samtools (Li et al., 2009). SNPs supported with read coverages below 5 were removed. VCF files were combined and formatted to Structure and Hapmap formats using the Perl script MultiVcfTool (<https://github.com/aubombarely/GenoToolBox/blob/master/SeqTools/MultiVcfTool>)

Homoeologue read identification and transcript-guided assembly

For homoeologous SNP identification, a consensus diploid transcriptome was rebuilt for each of the species groups (A, with *G. clandestina* and *G. canescens* accessions; D1, with *G. tomentella* D1 accessions; D3 with *G. tomentella* D3 accessions and D4 with *G. syndetika* accessions) using Samtools (Li et al., 2009) and Gffread from the Cufflinks software package (Trapnell et al., 2010). A progenitor reference set was created for each of the polyploid species joining the diploid transcriptome sets (T1=D1+D3, T2=D3+D4 and T5=A+D1). Reads from the polyploid species were mapped with these references using Bowtie2. Sam mapping files were processed to identify reads according the preferential mapping with each of the progenitors using the Perl script, SeparateHomeolog2Sam (<https://github.com/aubombarely/GenoToolBox/blob/master/SeqTools/SeparateHomeolog2Sam>). Reads with mapping score AS and XS = 0 (No SNPs) were kept and used to rebuilt the polyploid transcriptomes using Samtools (Li et al., 2009) and Gffread (from the Cufflinks package, (Trapnell et al., 2012)). Once the reads were separated according its preferential mapping, they were mapped back to the soybean genome. SNPs were called as described above.

Population structure analysis

The programs Structure (Pritchard et al., 2000) and fineStructure (Lawson et al., 2012) were used to analyze population structure of the two SNP datasets, with and without polyploid SNPs separated by homoeologue, described above. For Structure, each of the datasets was divided into three subsets of 20,000 SNPs selected with a random function incorporated in the MultiVcfTool. 5 replicates were run for each of the subsets with a burn-in of 10,000, from K=1 to K=15 using the default parameters. The optimal number of clusters was identified based on the rate of change in the log probability of data between successive K values (Evanno et al., 2005). Results were visualized using R (barplot function).

For fineStructure each of the two SNP datasets were divided into 20 different subsets each mapping to one soybean reference chromosome. Analyses were performed following the instructions from the fineStructure web for the unlinked model (http://www.maths.bris.ac.uk/~madjl/finestructure/data_example.html). Results were presented as a heatmap of distances between each of the accessions. A principal component analysis (PCA) was performed over the same distance matrix using fineStructure software. The PCA figure was created using R.

Reconstruction of phylogenies using concatenated SNPs

SNPs from the dataset in which SNPs from allopolyploids were partitioned into their two homoeologues ("homoeologue data set") were concatenated to create a supermatrix with 36 operational taxonomic units (OTUs). The two homoeologous gene copies from each allopolyploid were treated as individual OTUs; for example the D1 and D3 homoeologues of T1 individuals were treated as D1T1 and D3T1, respectively. *G. max*, accession William82 was used as outgroup. The alignment files were produced changing the SNPs Hapmap format to fasta using a Perl script. The resulting matrix was used in two analyses. First, maximum likelihood (ML) was used, implemented in

127 PhyML (Guindon and Gascuel, 2003) the substitution model with GTR as. 100 bootstrap replicates were conducted).
128 Second, in order to visualize reticulations in the dataset, a network method, NeighborNet, was implemented in the
129 SplitsTree package (Huson and Bryant, 2006) with the default parameters. Trees were visualized and drawn using
130 FigTree (Rambaut, 2012).

131 *Gene-based analyses*

132 A subset of transcripts was selected for phylogenetic and network analyses based on the following criteria:
133 No more than 10% of Ns for the guided assembly consensus sequence in any of the accessions after the
134 homoeologue read identification; alignments with at least 1000 bp; and genes with their corresponding *G. max*
135 homolog identified as an existing pair retained from the most recent (5-10 million years; (Doyle and Egan, 2010))
136 *Glycine* whole genome duplication (WGD) event, as compiled by Du et al. (Du et al., 2012). Sequence alignments
137 were based on the transcriptome-guided assembly. Sequence for each of the genes was collected with a Perl script
138 (FastaSeqExtract, GenoToolBox script package), joined with the Cat Linux command and changed to the required
139 sequence alignment format using a BioPerl script (bp_sreformat.pl). The 95 alignments selected were used in an
140 exploratory phylogenetic analysis using the Bayesian MCMC method, BEAST (Drummond et al., 2012) (HKY
141 substitution model, 10,000,000 MCMC). Alignments that produced trees in which *G. max* was not sister to perennial
142 *Glycine* species in the consensus tree were removed. Generally the removed alignments showed tree topologies with
143 two large clades with long branches, indicating the possibility of inclusion of paralogous genes from the older whole
144 genome duplication (ca. 50 MY, common to the Leguminosae; reviewed in Doyle 2012) instead the orthologue.

145 27 genes selected after this filtering were analyzed using three different methods: 1) Phylogenies were
146 reconstructed using ML using PhyML (Guindon and Gascuel, 2003) with 1,000 bootstraps. jModelTest2 was used to
147 choose the best substitution model (Darriba et al., 2012). According to the Bayesian Information Criterion (BIC)
148 HKY was the preferred model (40% of the genes), followed by K80 (26% of the genes) (Supplementary Table S2). 2)
149 Networks were constructed using NeighborNet in SplitsTree4 with the default parameters (Huson and Bryant, 2006).
150 3) Bayesian analysis was performed using BEAST v2.0 (Drummond et al., 2012). The two homoeologous gene
151 copies from each allopolyploid were treated as individual OTUs as in the concatenated analysis, and *G. max*,
152 accession William82 was again used as outgroup. Based on the jModelTest2 results, HKY was used as the
153 substitution model. The MCMC chain was set to 100,000,000 MCMC generations, taking samples every 1000
154 generations. Divergence ages were estimated by scaling the tree root (divergence between *G. max* and perennials) to 5
155 Myr (Egan and Doyle, 2010). All trees were drawn using FigTree (Rambaut, 2012).

156 *Species tree reconstruction*

157 Species tree reconstruction under the coalescent was performed using the 27 selected genes in *BEAST
158 (Drummond et al., 2012). The 24 accessions, including two homoeologues for each allopolyploid accession, were
159 grouped in 11 operational taxonomic units (OTUs) for this analysis: *G. canescens*, *G. clandestina*, *G. tomentella* D1,
160 *G. tomentella* D3, *G. syndetika* (D4), *G. tomentella* T1-D1, *G. tomentella* T1-D3, *G. dolichocarpa* T2-D3, *G.*
161 *dolichocarpa* T2-D4, *G. tomentella* T5-A and *G. tomentella* T5-D1. *G. max* was used as outgroup. Based on
162 jModelTest2 results, HKY was used as substitution model. The MCMC chain was set to 100,000,000 MCMC
163 generations, taking samples every 1000 generations. Divergence dates were estimated as described above. All the
164 trees were drawn using FigTree (Rambaut, 2012).

165 RESULTS

166 *Phylogenomics dataset generation*

167 Between 7-60 million reads from leaf transcriptomes of 24 accessions representing 8 *Glycine* perennial
168 species were mapped to the *Glycine max* genome (v1.0) (Schmutz et al., 2010). Reads mapped to 22,500-25,000
169 genes (~40% of soybean gene models; Table 1); this represents between 4.5 and 11.6% of the genome. 200,000-
170 965,000 single nucleotide polymorphisms (SNPs) were identified relative to *G. max*; 6.3-12.6% of SNP positions
171 were polymorphic in diploid species (*G. clandestina*, *G. canescens*, *G. tomentella* D1 (referred as D1 hereafter), *G.*
172 *tomentella* D3 (referred as D3) and *G. syndetika* (referred as D4)), and 18.4-28.8% in polyploid species (*G.*
173 *tomentella* T1 (referred as T1), *G. tomentella* T5 (referred as T5) and *G. dolichocarpa* T2 (referred as T2); Table 2).
174 The interpretation of these positions as standard heterozygosity is complicated by the recent (5-10 MYA: (Doyle and
175 Egan, 2010)) WGD in the ancestral *Glycine* genome. In a gene for which soybean has lost one of the homoeologous
176 copies from this event, but the perennial species for which it is serving as reference has retained both copies,
177 polymorphic SNPs may be due to reads from two different homoeologous loci in the perennial, rather than two
178 alleles at a single locus. Low levels of conventional heterozygosity are expected in *Glycine* species, because of their
179 strongly selfing reproductive biology, with much reproduction occurring through cleistogamous (closed, selfing)
180 flowers.

181 The much higher percentage of polymorphic positions in polyploid individuals (T1, T2, T5) likely is also
182 due to the mapping of reads from two homoeologous copies to a single target, in this case due to recent polyploidy:
183 for example, mapping reads from tetraploid ($2n = 80$) T2 to a single locus in the diploid ($2n = 40$) *G. max* reference
184 genome will result in reads from both its D3 and D4 homoeologous subgenomes mapping to the same target,
185 increasing the chance of observing a polymorphism at a given site. Separating reads from T1, T2, and T5 polyploid
186 individuals was possible where the read has at least a SNP that could be related to one homoeologous genome
187 contributor (e.g., D3 and D4 differed by a SNP and this difference was retained in the D3 and D4 homoeologous
188 genomes of T2; diploid-distinguishing polymorphism (DDP); see (Ilut et al., 2012)). Between 11.4 and 20.8% of
189 reads were assigned to one of the progenitors (Table 3).

190 Between 124,984 and 399,884 SNPs were produced for each accession. The filtering of the missing data
191 produced 237,243 and 75,958 polymorphic positions for all the accessions before and after the homoeologous read
192 assignment respectively. SNPs per chromosome ranged from 7,455 (chromosome 14) to 16,494 (chromosome 8) and
193 from 2,288 (chromosome 14) to 5,300 (chromosome 8) before and after the homoeologous read assignment
194 respectively.

195 Transcriptome-guided assemblies produced between ~1,800 and ~6,600 full-length sequences (as mapped to
196 the *G. max* gene models) for each diploid accession. For polyploid subtranscriptomes this number was much lower
197 because only reads that mapped preferentially to one of the diploid consensus species and reads that mapped equally
198 but with no polymorphism (perfect match) were used during the transcriptome-guided assembly. Any read that
199 mapped equally to two or more positions with one or more polymorphisms was discarded because it was impossible
200 to assign it to any of the diploid progenitors, reducing the mapping coverage of the reference gene models. Between
201 ~350 and ~1,350 full length sequences were assembled for the polyploid homoeologous subtranscriptomes of which
202 between 4 to 19% were duplicated genes from the last *Glycine* WGD event (Schmutz et al., 2010). For phylogenetic
203 analysis, full length sequences are not needed so a phylogenetic analysis dataset was created with 27 genes (see
204 Material and Methods for the criteria used to generate this dataset)(Table 5).

205 *Genome-wide distribution of homoeologous SNPs.*

206 For each allopolyploid accession, the ca. 120,000-400,000 SNPs (Table 3) that could be identified to
 207 homoeologous subgenome were mapped to the soybean reference genome (Schmutz et al., 2010). This produced a
 208 map that is analogous to chromosome painting (genomic in situ hybridization, GISH) experiments using the reads
 209 from which the SNPs were derived, which we term “electronic chromosome painting” (e-Chromopainting). Similar
 210 patterns were seen for all accessions, with high densities of SNPs at the ends of each soybean chromosome and far
 211 lower densities in pericentromeric regions (Fig. 2). This pattern is expected using reads from transcriptome data,
 212 because of the sparse distribution of genes in pericentromeric regions of the soybean genome (Schmutz et al., 2010).
 213 Notably, in all allopolyploid accessions, SNPs from both homoeologues were distributed across the entire genome,
 214 and no regions were identified in which SNPs from only one homoeologue were mapped (Fig 2; Supplementary Fig.
 215 1-10).

216 **Population structure analyses.**

217 Structure (Pritchard et al., 2000) was first run using all available SNPs, without separating SNPs from
 218 polyploids into homoeologous groups. Structure was run from $K = 1-15$; $K = 6$ was identified as one of the optimal
 219 preferred values of K using the delta K method of Evanno et al. (2005; Supplementary Fig. 11). Five of these six
 220 groups corresponded to diploid taxa: D1, D3, D4, *G. canescens*, and *G. clandestina* (Fig. 3a). The sixth group was
 221 represented only as a minor component in D4 accession 2073. Diploid accessions showed little or no evidence of
 222 admixture, with the exception of D4 accession 2073 (Fig. 3). In contrast, all polyploid accessions were admixed,
 223 each with approximately 50% contributions from two different diploid groups. The genomic makeup of each
 224 accession was as expected from previous hypotheses (e.g., Doyle et al., 2002; Fig. 1): T1 accessions showed
 225 admixture from D1 and D3, T2 accessions from D3 and D4, and natural T5 accessions from D1 and *G. clandestina*.
 226 The synthetic T5 accession (A58-1) was also admixed, with contributions from D1 and *G. canescens*, as expected
 227 (Joly et al., 2004).

228 A second Structure analysis was conducted with each polyploid accession treated as two separate OTUs,
 229 using the homoeologue dataset (Table 2). As with the previous analysis, the analysis was run for $K = 1-15$. The
 230 Evanno method (Evanno et al., 2005) identified $K = 6$ and 9 as the preferred values (Supplementary Fig. 11). In the
 231 case of $K = 9$ the group representation shows the same structure than the $K = 6$ (Supplementary Fig. 12). Results for
 232 diploids were similar to those obtained in the previous analysis (Fig. 3b). Subgenomes from natural allopolyploids
 233 and the synthetic T5 allopolyploid (A58-1) were shown to belong exclusively to diploid groups, with little or no
 234 evidence of admixture, indicating that the SNP filtering into homoeologous contributions was successful.

235 Complementary to the second Structure analysis, the data were analyzed using ChromoPainter and
 236 FineStructure (Lawson et al., 2012). ChromoPainter produces a co-ancestry matrix (as a measure of the ancestry
 237 sharing between individuals) based on the haplotype information provided by shared chunks (regions) of biallelic
 238 markers between individuals (Lawson et al., 2012). The two SNP datasets were filtered by selecting only the biallelic
 239 markers, producing a subset with 220,952 and 71,610 SNPs (before and after homoeologous read assignment,
 240 respectively) distributed along all 20 soybean chromosomes. Regions identified by ChromoPainter for each
 241 accession ranged from 516 (D4 2321) to 567 (*G. clandestina* 1253) and from 202 (D4 1300 and 2321) to 221 (D4
 242 2073) (before and after homoeologous read assignment respectively). Principal component analysis (PCA) and
 243 population relationship analysis using a Bayesian approach were performed over the co-ancestry matrix using
 244 FineStructure (Lawson et al., 2012). PCA before homoeologous read assignment (Fig. 4a) shows seven well-
 245 differentiated groups, one per species with the exception that *G. canescens* and *G. clandestina* clustered together.
 246 Diploid species formed the vertices of a trapezoid. A-genome species (*G. canescens*, *G. clandestina* and D4) formed
 247 a more dispersed group than either D1 or D3. Each polyploid species fell between its putative diploid progenitors,
 248 consistent with each being an admixture (fixed hybrid). After the homoeologous read assignment (Fig. 4b), each of
 249 the polyploid subgenomes clustered with its diploid progenitors, producing three clear clusters: D1, D3, and A-
 250 genome (comprising *G. canescens*, *G. clandestina* and D4, as expected). Heatmaps were used to visualize the

population relationships produced by FineStructure, complementing the information shown by the PCA figures. The heatmap before homoeologous read assignment (Fig. 4c), showed four intense regions (red, magenta and blue colors) corresponding to the four species groups of the PCA (Fig. 4a). Each polyploid showed the expected similarity to its progenitors; similarly, as expected the two *G. clandestina* accessions were more similar to one another than either was to *G. canescens*. Also, T5 A58, the artificial polyploid produced from a cross between *G. canescens* 1232 and D1 1316, showed the expected relationships with these accessions. Other T5 polyploids also showed a stronger signal from D1 1316 than from other D1 accessions. T2 accessions did not show any stronger signal with any particular D3 accession than with others, but they did with the D4 accessions 1300 and 2321, relative to 2073. T1 accessions 1288 and 1763 also showed a stronger signal with particular D1 and D3 accessions, whereas T1 accession 1361 showed a weaker signal with the D1 and D3 accessions included here. After the homoeologous read assignment (Fig. 4d), some of these signals were intensified, such as the relationship between T5-D1 subgenomes and particular D1 accessions, but other relationships that were suggested when all SNPs were considered were not observed (for example there is not a stronger signal of D1 1316 with the T5 accessions). These differences may be due to the methodology used for the homoeologous read assignment.

Phylogeny and network analysis of concatenated SNPs.

Phylogenetic and network analyses were conducted using the homoeologue dataset, with SNPs concatenated to create a single supermatrix. The maximum likelihood (ML) tree, rooted with *G. max*, identified four subclades comprising two major clades: 1) the A-genome, with subclades of D4 vs. *G. clandestina* and *G. canescens*; and 2) the D-genome (D3) and E-genome (D1) (Fig. 5a). Each of the subclades showed a different pattern with respect to diploid and tetraploid subgenome relationships. In the *canescens/clandestina* clade, the A-subgenome of the synthetic allopolyploid (A58) was sister to the accession from which it was created (*G. canescens* 1232), as expected; the two natural T5 allopolyploids were sister to *G. clandestina* 1126, as expected from other data (e.g., Doyle et al. 2002). In the D4 clade, diploid accession 2073 was sister to all remaining accessions, a unique placement consistent with its apparently admixed nature (Fig. 3a). The polyploid subgenomes formed a paraphyletic group, with the two diploid accessions sister to the D4 subgenome of one T2 accession (1134). A similar pattern was seen in the D3 subclade, where T2 accessions formed a paraphyletic group, and all four diploid accessions formed a clade sister to T2 accession 1134. Also embedded within the T2 accessions was a clade consisting solely of T1 accessions. T1 accessions also formed a monophyletic group within the D1 clade, where natural T5 accessions and D1 accessions also formed monophyletic groups. Surprisingly, there was not a sister relationship between the D1-subgenome of synthetic allopolyploid A58 and the D1 accession from which it was formed (1316). Similar topologies were produced by neighbor-joining analysis (data not shown).

NeighborNet was used to analyze the full homoeologue dataset to identify minority patterns of relationships in the data. When rooted with *G. max*, the topology (Fig. 5b) was very similar to the ML tree (Fig. 5a), even having such features as the sister relationship of D4 2073 to other D4 accessions, and the monophyly of T1 homoeologues in both the D1 and D3 clades. There was clear evidence of character support for alternative relationships, but those relationships were minor in comparison with the major phylogenetic signal.

Gene-based phylogenetic and network analyses

Gene trees were constructed for the 27 genes (described in the Material and Method) using several different phylogenetic and network methods. Similar topologies for trees from individual genes were obtained with BEAST and PhyML. All 27 trees showed the split between the A-genome clade and the D1/D3 clade seen in the ML tree reconstructed from concatenated SNPs (Fig. 5a). However, many individual gene trees showed unexpected groupings of one or more accessions, particularly within the A-genome clade, where several trees grouped accessions from *G.*

canescens with *G. syndetika*-D4 instead of with *G. clandestina* (for example ML and BEAST trees for the gene Glyma04g39670, Supplementary Figure 17 and 45). Relationships within the major subclades varied among the 27 gene trees. For example, nine of the 27 trees showed separate clades for *G. canescens* (plus the A58 sequence) and *G. clandestina* (e.g., Fig. 6, a and c), but in only three of them did diploid species form monophyletic groups (Supplementary Figures 13 to 67). Overall, there were far more departures from expectations in the A-genome clade than in the D1/D3 clade.

There were numerous cases where alleles from diploid accessions formed monophyletic groups (e.g., 12 of 27 BEAST topologies had alleles from all four D3 accessions in a clade, often with high posterior probability). At some loci, alleles from one or more polyploids formed monophyletic clades; for example, at Glyma06g18640 (Supplementary Figure 50), all taxa, including both homoeologous subgenomes of each polyploid, formed separate clades, with the exception of *G. clandestina*. However, this was unusual, and paraphyletic groupings of alleles were common, particularly in polyploids. For example, at 26 of 27 loci, T2-D3 alleles were not monophyletic, at least some having closer relationships to D3 or T1-D3 alleles, and in gene Glyma01g35620, T5-D1_1969 was most closely related to D1_1156 whereas T5-D1_1487 was most closely related to D1_1157 and D1_1316 (Supplementary Figure 41). On the assumption that alleles in tetraploids all originated from diploid progenitors, such paraphyletic relationships suggest multiple origins or perhaps subsequent gene flow from diploids after polyploid formation. The BEAST trees, calibrated with the 5 MYA divergence of *G. max* and the perennial subgenus (Innes et al., 2008), allowed dates of allele divergence to be estimated. Among comparisons of interest are the minimum distances between alleles from a tetraploid and alleles from its diploid progenitor (e.g., T2-D3 vs. D3) or alleles from the same progenitor in a second tetraploid (e.g., T2-D3 vs. T1-D3); the latter represent “diploid” alleles as well, under the assumption that there has been no gene flow between the two tetraploids, something that is reasonable for *G. tomentella* tetraploids (e.g., Doyle et al. 1986). Minimum distances between polyploid and diploid alleles (over)estimate the time of entry of that allele into the polyploid, which is typically assumed to be an origin of the polyploid (Doyle and Egan 2010). Minimum dates (Supplemental Table 3) were 0.31 MY for T1 (measured at the D1 locus), 0.29 MY for T5 (measured at the D1 locus), and 0.38 MY for T2 (measured at the D3 locus). Error bars on these estimates, however, were substantial.

NeighborNet (implemented in SplitsTree 4; (Huson and Bryant, 2006)) was used to construct networks for each of the 27 genes. This method showed patterns consistent with intragenic recombination; the Pairwise Homoplasy Index (PHI) of Bruen et al. (2006), also implemented in SplitsTree, was significant for 11 of the 27 genes (data not shown). The dominant patterns in NeighborNet topologies were similar to the overall pattern shown in phylogenetic analyses of the 27 genes, and thus to results for the full homoeologous SNP dataset. As with other methods, NeighborNet networks suggested multiple inputs of alleles from diploid progenitors into polyploids (e.g., gene Glyma02g11580, Fig. 6c).

Species tree reconstruction under the coalescent

Species trees were reconstructed using the coalescent approach implemented in *BEAST (Heled and Drummond, 2010), which used information contained in the individual gene trees from the 27 genes described above. The overall *BEAST tree (Fig. 7a) topology was similar to that of trees from concatenated SNPs. By definition, each of the allopolyploid homoeologous genomes was a single OTU despite the possibility of independent origins; each of these was grouped with its putative progenitor species. Within the D1 genome clade, the T1 and T5 polyploids were sisters to one another; similarly, T1 and T2 were sisters in the D3 clade. The DensiTree output (Supplementary Figure 40) indicated considerable uncertainty only within the D3 clade, where both other possible topologies (T2 sister to D3, T1 sister to D3) appeared in a substantial number of trees. As expected, divergence dates of polyploids from their diploid progenitors estimated by *BEAST were higher than minimum estimates from the 27

individual loci, all being greater than 300,000 years (Fig. 7a).

DISCUSSION

The *Glycine* subgenus *Glycine* polyploid complex appears ideally suited as a model for studying allopolyploid evolution, because it comprises eight independently formed but closely related allopolyploid species triads (an allotetraploid and its two diploid progenitors; Fig. 1) that overlap in their composition. We are exploiting this model system to study the effect of allopolyploidy on a wide range of phenotypes, including transcriptome size, morphology, anatomy, climate niche, photosynthesis, and photoprotection (Coate and Doyle, 2010; Coate et al., 2012; Ilut et al., 2012; Coate et al., 2013; Coate and Doyle, 2013; Hegarty et al., 2013; Coate et al., 2014; Habert et al., 2014)

Various lines of evidence culminated in the hypotheses of reticulate relationships within the complex (Doyle et al. 2002) some of which are shown in Fig. 1. Chromosome number polymorphism ($2n = 38, 40, 78, 80$) was observed in what was initially considered the single taxon, *Glycine tomentella* (Newell and Hymowitz 1978). Patterns of sterility and partial chromosome pairing in artificial crosses among polyploid *G. tomentella* plants were consistent with the presence of shared homoeologous diploid genomes among polyploids (Grant et al. 1984; Singh et al. 1988). Isozyme studies of allopolyploids and candidate diploid species led to the characterization of numerous “races” of *G. tomentella* designated “D” for diploid, or “T” for tetraploid (Doyle et al., 1986; Singh et al., 1998). Eventually, molecular phylogenetic studies assumed a dominant role in refining hypotheses of relationships; (Hsing et al., 2001; Brown et al., 2002; Doyle et al., 2002; Rauscher et al., 2002). However, these DNA-based studies are based on only two nuclear markers: the internal transcribed spacers of the 18S-5.8S-26S nuclear ribosomal gene cistron (nrDNA ITS) and the low copy gene, histone H3-D. Relationships of chloroplast genomes are broadly consistent with these results (Hsing et al., 2001), but are complicated by incongruence with nuclear markers, likely due to a combination of incomplete lineage sorting and introgression (Doyle et al., 2004).

It is known that the earliest stages of polyploid evolution can produce genomic shock (McClintock, 1984), in which whole chromosomes, chromosomal segments, or individual genes may be lost (e.g., Chester et al., 2012). Later stages in the evolutionary cycle of a polyploid can involve the progressive loss of DNA in the phenomenon of “genomic downsizing” (Leitch and Bennett, 2004) and the loss of individual genes from one homoeologue in the process of “fractionation” (Schnable and Freeling, 2011; Freeling et al., 2012). In addition to the loss of genes, the process of concerted evolution, notably through gene conversion, can result in the replacement of a gene from one genome by its homoeologue, (e.g., Wang et al., 2007). Because of these processes, although it is expected that an allopolyploid will be a fixed hybrid, combining the homoeologous genomes of its two parents, this may never have been completely true across the entire genome since the earliest stages of its evolution, and becomes progressively less likely to be true as time passes. Thus, we were interested in knowing whether the hypotheses of relationships formulated for species of the perennial *Glycine* allopolyploid complex using single genes were true of the whole genomes, or whether there was evidence of loss of homoeologous genes or chromosomal segments from one or more allopolyploid species. Given the results from *Tragopogon miscellus*, which was found to be polymorphic for unequal contributions from its two diploid progenitors at the chromosome and chromosome segment level (Chester et al., 2012), we also wished to sample multiple individuals of *Glycine* allopolyploids.

High-throughput sequencing produces massive amounts of genome-wide data, and thus has great potential for, systematic and evolutionary studies in general (Gilad et al., 2009), and for addressing our questions in particular. However, it is not trivial to extract relevant information from short read sequencing data, particularly for allopolyploids, where the interest is often in deconvoluting the complex genome into its two homoeologous subgenomes (Grover et al., 2012; Ilut et al., 2012). Here we used a leaf transcriptome dataset from an experiment that was originally designed to explore the effects of allopolyploidy on photoprotection (Coate et al., 2013) to elucidate evolutionary patterns in multiple individuals from three overlapping *Glycine* allopolyploid species triads

381 (Fig. 1). Also included in the study was a synthetic allopolyploid, A-58, which mimics the nuclear genome
382 composition of the natural T5 allopolyploid (Joly et al., 2004).

383 ***Glycine* allopolyploids are fixed hybrids throughout their genomes**

384 Analyses using all SNPs identified from the full dataset showed that all three of these allopolyploids are
385 indeed fixed hybrids, combining diploid genomes as depicted in Fig. 1. Structure results indicated an essentially
386 equal contribution from both parental diploids in all three cases (Fig. 3a); PCA analysis also was consistent with this
387 hypothesis, placing each polyploid approximately midway between its putative progenitors, as expected for an F1
388 hybrid (Fig. 4a).

389 In order to determine whether or not the polyploids have contributions from their parents across their entire
390 genomes, reads were partitioned by homoeologous genome and mapped to the soybean reference genome (Schmutz
391 et al., 2010). As portrayed by e-chromosome painting (Fig. 2), it is clear that no individual sampled from any of the
392 three allopolyploid species has any major regions represented by only one homoeologue. Coverage is sparse in
393 pericentromeric and centromeric regions, as expected due to the low density of genes in these regions of the soybean
394 genome (Schmutz et al., 2010). The degree of shared synteny between soybean and perennial *Glycine* species is as
395 yet unknown, but regardless of the order of chromosomal segments, it is clear that there has not been significant loss
396 of homoeologous genes. We mapped reads to over 22,000 of the approximately 46,000 genes of the soybean genome
397 (Schmutz et al., 2010). These numbers include both homoeologous copies from the 5-10 MYA polyploidy event that
398 shaped the modern “diploid” ($2n = 38,40$) *Glycine* genome. We were able to deconvolute between 4 and 19% of
399 these genes into their homoeologous contributions in each of the three recent allopolyploids. Using genomic in situ
400 hybridization (GISH), (Chester et al., 2012) showed examples of allopolyploid *T. miscellus* plants that had all four
401 chromosomes or chromosome segments of one diploid parent (4:0), but also examples of plants with 3:1 ratios of
402 homoeologous chromosomes or chromosomal segments. Our e-chromosome painting method cannot distinguish the
403 3:1 condition from an equal contribution from both parents segments, so it is possible that such plants exist in our
404 sample.

405 Structure analysis using the partitioned homoeologous SNPs corroborated results with the full, unpartitioned
406 dataset, in placing each polyploid homoeologous genome with its putative progenitor (Fig. 3b). The FineStructure
407 PCA included three major groupings, each of which included diploids and the expected polyploid homoeologous
408 subgenomes derived from them (Fig. 4b). The grouping of D4 accessions and A-genome (*G. canescens* and *G.*
409 *clandestina*) into a single cluster is not surprising, because *G. syndetika* (D4) is also a member of the A-genome
410 (Ratnaparkhe et al., 2011). Genome groups were originally defined on the basis of reproductive compatibility in
411 artificial crosses (Ratnaparkhe et al. 2011), and indeed *G. syndetika* (D4) 2073 shows evidence of admixture with *G.*
412 *canescens* and *G. clandestina* (Fig. 3). In contrast, D1 and D3, though both classified as “*G. tomentella*”, belong to
413 two different genome groups (E and D, respectively; Ratnaparkhe et al. 2011). This greater genetic similarity of the
414 three A-genome species is not reflected in relative divergence dates; for example, the *BEAST analysis dates the
415 divergence between *G. syndetika* and the two other A-genome species at slightly earlier than the divergence between
416 D1 and D3 (Fig. 7a). Thus, reproductive barriers likely arose earlier after divergence in the D1/D3 lineage than
417 within the A-genome.

418 Allopolyploid evolution in *Glycine* fits “Darlington’s Rule” (Darlington, 1937)—that allopolyploids should
419 form between species that are reproductively isolated, often due to chromosomal differences, whereas reproductively
420 compatible diploids tend to form homoploid hybrids. No allopolyploids are known to have formed among A-genome
421 species, and only one of the eight known *Glycine* allopolyploids involves hybridization within a genome group
422 (tetraploid *G. tabacina* is the product of the most divergent species cross possible within the B-genome; Fig. 1;
423 (Doyle et al., 2004)). In contrast, D1 and D3 have different chromosome numbers ($2n = 38$ vs. 40, respectively),
424 which contribute to their ability to form fertile diploid hybrids. D1 has also formed allopolyploids with D5A, another

425 $2n = 40$ “*G. tomentella*”; however, reproductive incompatibility occurs between $2n = 40$ *G. tomentella* taxa, and
426 other allopolyploids in the complex combine genomes of $2n = 40$ taxa (Fig. 1).

427 ***Gene histories, allele divergence times, and sources of genetic diversity in polyploids.***

428 Gene trees from the 27 loci that met criteria designed to provide orthologous, highly transcribed genes with
429 sufficient characters for phylogeny reconstruction and inferences of polyploid origins mostly conformed to
430 expectations based on previous work using the low copy nuclear locus, histone H3D (Brown et al., 2002; Doyle et
431 al., 2002; González-Orozco et al., 2012), the nrDNA ITS (Singh et al., 2001; Rauscher et al., 2004), and chloroplast
432 noncoding sequences (Hsing et al. 2001). The use of BEAST and *BEAST (Heled and Drummond, 2010) allowed us
433 to estimate divergence times of alleles and species for the first time for some of these taxa. Dating polyploid origins
434 is complicated by numerous factors (Doyle and Egan, 2010). For one thing, if the polyploid has arisen recurrently,
435 then there is no single date that marks “the” origin. The relevant date for testing the anthropogenic disturbance
436 hypothesis would be the oldest origin. However, because it is unlikely that a polyploid allele and any of a set of
437 diploid progenitor alleles will coalesce at exactly the time of polyploid origin, distances for any given polyploid
438 event will be overestimates of the actual time of origin. Further complicating matters, the error bars on our BEAST
439 divergence estimates were large relative to the estimates themselves. Nevertheless, it appears likely that these *G.*
440 *tomentella* allopolyploids are hundreds of thousands rather than tens of thousands of years old. *BEAST estimates
441 should be averages of all origins of a polyploid taxon, and these, too are several hundred thousand years for each
442 allopolyploid. Thus, it appears likely that these polyploid species were present in Australia long before humans
443 arrived in Australia (Hudjashov et al., 2007). The fact that these three species, and possibly other allopolyploid
444 members of the complex, may have evolved at roughly the same time is intriguing. In the ca. 5 MY since the
445 perennial members of *Glycine* diverged from the annual lineage (Egan and Doyle, 2010), there is no evidence of
446 polyploidy until these species were formed, probably within the last 1 MY. Perhaps the onset of severe aridity in
447 Australia around 3 MYA, heralding the change to the present extreme wet-dry glacial cycles (Crisp et al., 2004)
448 could have provided ecological opportunities for polyploids. It will be interesting to refine our estimates through
449 increased sampling of these three triads, and to obtain estimates for the other five allopolyploid species.

450 REFERENCES

- 451 **Arnold B, Bomblies K, Wakeley J** (2012) Extending Coalescent Theory to Autotetraploids. *Genetics* **192**: 195–204
- 452 **Aronesty E** (2013) Comparison of sequencing utility programs. *Open Bioinform J*
- 453 **Blanc G, Wolfe K** (2004) Widespread paleopolyploidy in model plant species inferred from age distributions of
454 duplicate genes. *Plant Cell* **16**: 1667–1678
- 455 **Brown AHD, Doyle JL, Grace JP, Doyle JJ** (2002) Molecular phylogenetic relationships within and among
456 diploid races of *Glycine tomentella* (Leguminosae). *Australian Systematic Botany* **15**: 37–47
- 457 **Cannon SB, Ilut D, Farmer AD, Maki SL, May GD, Singer SR, Doyle JJ** (2010) Polyploidy did not predate the
458 evolution of nodulation in all legumes. *PLoS ONE* **5**: e11630
- 459 **Chester M, Gallagher JP, Symonds VV, Cruz da Silva AV, Mavrodiev EV, Leitch AR, Soltis PS, Soltis DE**
460 (2012) Extensive chromosomal variation in a recently formed natural allopolyploid species, *Tragopogon miscellus*
461 (Asteraceae). *P Natl Acad Sci Usa* **109**: 1176–1181
- 462 **Coate JE, Doyle JJ** (2010) Quantifying whole transcriptome size, a prerequisite for understanding transcriptome
463 evolution across species: an example from a plant allopolyploid. *Genome Biol Evol* **2**: 534–546
- 464 **Coate JE, Doyle JJ** (2013) Genomics and transcriptomics of photosynthesis in polyploids.
- 465 **Coate JE, Luciano AK, Seralathan V, Minchew KJ, Owens TG, Doyle JJ** (2012) Anatomical, biochemical, and
466 photosynthetic responses to recent allopolyploidy in *Glycine dolichocarpa* (Fabaceae). *Am J Bot* **99**: 55–67
- 467 **Coate JE, Powell AF, Owens TG, Doyle JJ** (2013) Transgressive physiological and transcriptomic responses to
468 light stress in allopolyploid *Glycine dolichocarpa* (Leguminosae). *Heredity* **110**: 160–170
- 469 **Coate JE, Bar H, Doyle JJ** (2014) Extensive translational regulation of gene expression in an allopolyploid
470 correlates with long term retention of duplicated genes. *The Plant Cell*. doi: <http://dx.doi.org/10.1105/tpc.113.119966>
- 471 **Crisp M, Cook L, Steane D** (2004) Radiation of the Australian flora: what can comparisons of molecular
472 phylogenies across multiple taxa tell us about the evolution of diversity in present-day communities? *Philos Trans R*
473 *Soc Lond, B, Biol Sci* **359**: 1551–1571
- 474 **Darlington CD** (1937) Recent advances in cytology. *Recent advances in cytology*
- 475 **Darriba D, Taboada GL, Doallo R, Posada D** (2012) jModelTest 2: more models, new heuristics and parallel
476 computing. *Nat Meth* **9**: 772
- 477 **Doyle JJ** (2012) Polyploidy in Legumes. In PS Soltis, DE Soltis, eds, *Polyploidy and Genome Evolution*. Springer
478 Berlin Heidelberg, Berlin, Heidelberg, pp 147–180
- 479 **Doyle JJ, Doyle JL, Brown A, Palmer RG** (2002) Genomes, Multiple Origins, and Lineage Recombination in the
480 *Glycine Tomentella* (Leguminosae) Polyploid Complex: Histone H3-D Gene Sequences. *Evolution*
- 481 **Doyle JJ, Doyle JL, Brown AH** (1999) Origins, colonization, and lineage recombination in a widespread perennial
482 soybean polyploid complex. *P Natl Acad Sci Usa* **96**: 10741–10745
- 483 **Doyle JJ, Doyle JL, Rauscher J, Brown A** (2004) Diploid and Polyploid Reticulate Evolution Throughout the
484 History of the Perennial Soybeans (*Glycine* Subgenus *Glycine*). *New Phytol* **161**: 121–132
- 485 **Doyle JJ, Egan AN** (2010) Dating the origins of polyploidy events. *New Phytol* **186**: 73–85

- 486 **Doyle JJ, Schuler MA, Godette WD, Zenger V, Beachy RN, Slightom JL** (1986) The glycosylated seed storage
487 proteins of *Glycine max* and *Phaseolus vulgaris*. Structural homologies of genes and proteins. *Journal of Biological*
488 ...
- 489 **Drummond AJ, Suchard MA, Xie D, Rambaut A** (2012) Bayesian phylogenetics with BEAUti and the BEAST
490 1.7. *Mol Biol Evol* **29**: 1969–1973
- 491 **Du J, Tian Z, Sui Y, Zhao M, Song Q, Cannon SB, Cregan P, Ma J** (2012) Pericentromeric effects shape the
492 patterns of divergence, retention, and expression of duplicated genes in the paleopolyploid soybean. *Plant Cell* **24**:
493 21–32
- 494 **Dufresne F, Stift M, Vergilino R, Mable BK** (2014) Recent progress and challenges in population genetics of
495 polyploid organisms: an overview of current state-of-the-art molecular and statistical tools. *Mol Ecol* **23**: 40–69
- 496 **Edwards SV** (2009) Is a new and general theory of molecular systematics emerging? *Evolution* **63**: 1–19
- 497 **Egan AN, Doyle JJ** (2010) A comparison of global, gene-specific, and relaxed clock methods in a comparative
498 genomics framework: dating the polyploid history of soybean (*Glycine max*). *Syst Biol* **59**: 534–547
- 499 **Evanno G, Regnaut S, Goudet J** (2005) Detecting the number of clusters of individuals using the software
500 STRUCTURE: a simulation study. *Mol Ecol* **14**: 2611–2620
- 501 **Freeling M, Woodhouse MR, Subramaniam S, Turco G, Lisch D, Schnable JC** (2012) Fractionation mutagenesis
502 and similar consequences of mechanisms removing dispensable or less-expressed DNA in plants. *Curr Opin Plant*
503 *Biol* **15**: 131–139
- 504 **Gilad Y, Pritchard JK, Thornton K** (2009) Characterizing natural variation using next-generation sequencing
505 technologies. *Trends Genet* **25**: 463–471
- 506 **González-Orozco CE, Brown AHD, Knerr N, Miller JT, Doyle JJ** (2012) Hotspots of diversity of wild Australian
507 soybean relatives and their conservation in situ. *Conserv Genet* **13**: 1269–1281
- 508 **Grover CE, Salmon A, Wendel JF** (2012) Targeted sequence capture as a powerful tool for evolutionary analysis.
509 *Am J Bot* **99**: 312–319
- 510 **Guindon S, Gascuel O** (2003) A simple, fast, and accurate algorithm to estimate large phylogenies by maximum
511 likelihood. *Syst Biol* **52**: 696–704
- 512 **Habert RS, Brown AHD, Doyle JJ** (2014) Allopolyploidy, climate niche modeling, and evolutionary “success” in
513 *Glycine* (Leguminosae). *American Journal of ...*
- 514 **Hegarty M, Coate J, Sherman-Broyles S, Abbott R, Hiscock S, Doyle J** (2013) Lessons from natural and artificial
515 polyploids in higher plants. *Cytogenet Genome Res* **140**: 204–225
- 516 **Heled J, Drummond AJ** (2010) Bayesian inference of species trees from multilocus data. *Mol Biol Evol* **27**: 570–
517 580
- 518 **Hollister JD, Arnold BJ, Svedin E, Xue KS, Dilkes BP, Bomblies K** (2012) Genetic adaptation associated with
519 genome-doubling in autotetraploid *Arabidopsis arenosa*. *PLoS Genet* **8**: e1003093
- 520 **Hsing Y-LC, Hsieh J-S, Peng C-L, Chou C-H, Chiang T-Y** (2001) Systematic Status of the *Glycine tomentella* and
521 *G. tabacina* Species Complexes (Fabaceae) Based on ITS Sequences of Nuclear Ribosomal DNA. *J Plant Res* **114**:
522 435–442
- 523 **Hudjashov G, Kivisild T, Underhill PA, Endicott P, Sanchez JJ, Lin AA, Shen P, Oefner P, Renfrew C, Villems**
524 **R, et al** (2007) Revealing the prehistoric settlement of Australia by Y chromosome and mtDNA analysis. *P Natl Acad*
525 *Sci Usa* **104**: 8726–8730

- 526 **Huson DH, Bryant D** (2006) Application of phylogenetic networks in evolutionary studies. *Mol Biol Evol* **23**: 254–
527 267
- 528 **Hymowitz T, Singh RJ, Kollipara KP** (2010) The Genomes of the Glycine. *In Plant Breeding Reviews*. John Wiley
529 & Sons, Inc, Oxford, UK, pp 289–317
- 530 **Ilut DC, Coate JE, Luciano AK, Owens TG, May GD, Farmer A, Doyle JJ** (2012) A comparative transcriptomic
531 study of an allotetraploid and its diploid progenitors illustrates the unique advantages and challenges of RNA-seq in
532 plant species. *Am J Bot* **99**: 383–396
- 533 **Innes RW, Ameline-Torregrosa C, Ashfield T, Cannon E, Cannon SB, Chacko B, Chen NWG, Couloux A,**
534 **Dalwani A, Denny R, et al** (2008) Differential accumulation of retroelements and diversification of NB-LRR
535 disease resistance genes in duplicated regions following polyploidy in the ancestor of soybean. *Plant Physiol* **148**:
536 1740–1759
- 537 **Jiao Y, Leebens-Mack J, Ayyampalayam S, Bowers JE, McKain MR, McNeal J, Rolf M, Ruzicka DR, Wafula**
538 **E, Wickett NJ, et al** (2012) A genome triplication associated with early diversification of the core eudicots. *Genome*
539 *Biol* **13**: R3
- 540 **Jiao Y, Wickett NJ, Ayyampalayam S, Chanderbali AS, Landherr L, Ralph PE, Tomsho LP, Hu Y, Liang H,**
541 **Soltis PS, et al** (2011) Ancestral polyploidy in seed plants and angiosperms. *Nature* **473**: 97–U113
- 542 **Joly S, Rauscher JT, Sherman-Broyles SL, Brown AHD, Doyle JJ** (2004) Evolutionary dynamics and preferential
543 expression of homeologous 18S-5.8S-26S nuclear ribosomal genes in natural and artificial glycine allopolyploids.
544 *Mol Biol Evol* **21**: 1409–1421
- 545 **Jones G, Sagitov S, Oxelman B** (2013) Statistical inference of allopolyploid species networks in the presence of
546 incomplete lineage sorting. *Syst Biol* **62**: 467–478
- 547 **Langmead B, Salzberg SL** (2012) Fast gapped-read alignment with Bowtie 2. *Nat Meth* **9**: 357–359
- 548 **Lawson DJ, Hellenthal G, Myers S, Falush D** (2012) Inference of population structure using dense haplotype data.
549 *PLoS Genet* **8**: e1002453
- 550 **Leitch IJ, Bennett MD** (2004) Genome downsizing in polyploid plants. *Biological Journal of the Linnean* ...
- 551 **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome**
552 **Project Data Processing Subgroup** (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**:
553 2078–2079
- 554 **Lynch M, Conery JS** (2000) The evolutionary fate and consequences of duplicate genes. *Science* **290**: 1151–1155
- 555 **Maere S, De Bodt S, Raes J, Casneuf T, Van Montagu M, Kuiper M, Van de Peer Y** (2005) Modeling gene and
556 genome duplications in eukaryotes. *P Natl Acad Sci Usa* **102**: 5454–5459
- 557 **McClintock B** (1984) The significance of responses of the genome to challenge. *Science* **226**: 792–801
- 558 **Pritchard J, Stephens M, Donnelly P** (2000) Inference of population structure using multilocus genotype data.
559 *Genetics* **155**: 945–959
- 560 **Rambaut A** (2012) FigTree version 1.4.0. <http://tree.bio.ed.ac.uk/software/figtree/>
- 561 **Ratnaparkhe MB, Singh RJ, Doyle JJ** (2010) Glycine. *In Wild Crop Relatives: Genomic and Springer Berlin*
562 *Heidelberg, Berlin, Heidelberg*, pp 83–116
- 563 **Rauscher JT, Doyle JJ, Brown AHD** (2004) Multiple origins and nrDNA internal transcribed spacer homeologue
564 evolution in the Glycine tomentella (Leguminosae) allopolyploid complex. *Genetics* **166**: 987–998

565 **Rauscher JT, Doyle JJ, Brown AHD** (2002) Internal transcribed spacer repeat-specific primers and the analysis of
566 hybridization in the *Glycine tomentella* (Leguminosae) polyploid complex. *Mol Ecol* **11**: 2691–2702

567 **Schlueter J, Dixon P, Granger C, Grant D, Clark L, Doyle JJ, Shoemaker R** (2004) Mining EST databases to
568 resolve evolutionary events in major crop species. *Genome* **47**: 868–876

569 **Schmutz J, Cannon SB, Schlueter J, Ma J, Mitros T, Nelson W, Hyten DL, Song Q, Thelen JJ, Cheng J, et al**
570 (2010) Genome sequence of the palaeopolyploid soybean. *Nature* **463**: 178–183

571 **Schnable JC, Freeling M** (2011) Genes identified by visible mutant phenotypes show increased bias toward one of
572 two subgenomes of maize. *PLoS ONE* **6**: e17855

573 **Shoemaker RC, Schlueter J, Doyle JJ** (2006) Paleopolyploidy and gene duplication in soybean and other legumes.
574 *Curr Opin Plant Biol* **9**: 104–109

575 **Singh RJ, Kim HH, Hymowitz T** (2001) Distribution of rDNA loci in the genus *Glycine* Willd. *Theor Appl Genet*
576 **103**: 212–218

577 **Singh RJ, Kollipara KP, Hymowitz T** (1998) The genomes of *Glycine canescens* FJ Herm., and *G. tomentella*
578 Hayata of Western Australia and their phylogenetic relationships in the genus *Glycine* Willd. *Genome*

579 **Slotte T, Bataillon T, Hansen TT, St Onge K, Wright SI, Schierup MH** (2011) Genomic determinants of protein
580 evolution and polymorphism in *Arabidopsis*. *Genome Biol Evol* **3**: 1210–1219

581 **Slotte T, Hazzouri KM, Ågren JA, Koenig D, Maumus F, Guo Y-L, Steige K, Platts AE, Escobar JS, Newman**
582 **LK, et al** (2013) The *Capsella rubella* genome and the genomic consequences of rapid mating system evolution. *Nat*
583 *Genet* **45**: 831–835

584 **Soltis DE, Albert VA, Leebens-Mack J, Bell CD, Paterson AH, Zheng C, Sankoff D, dePamphilis CW, Wall**
585 **PK, Soltis PS** (2009) Polyploidy and angiosperm diversification. *Am J Bot* **96**: 336–348

586 **Symonds VV, Soltis PS, Soltis DE** (2010) Dynamics of polyploid formation in *Tragopogon* (Asteraceae): recurrent
587 formation, gene flow, and population structure. *Evolution* **64**: 1984–2003

588 **Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn JL, Pachter L**
589 (2012) Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat*
590 *Protoc* **7**: 562–578

591 **Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ, Pachter L**
592 (2010) Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching
593 during cell differentiation. *Nat Biotechnol* **28**: 511–U174

594 **Wang X, Tang H, Bowers JE, Feltus FA, Paterson AH** (2007) Extensive concerted evolution of rice paralogs and
595 the road to regaining independence. *Genetics* **177**: 1753–1763

596 **Wood TE, Takebayashi N, Barker MS, Mayrose I, Greenspoon PB, Rieseberg LH** (2009) The frequency of
597 polyploid speciation in vascular plants. *P Natl Acad Sci Usa* **106**: 13875–13879

Figure 1

Schema of the *Glycine* perennial polyploid complex.

Schema of the *Glycine* perennial polyploid complex. Diploid progenitors are represented by circles and allotetraploid species by circles. Chromosome number are represented by $2n = XX$. Species used in this study (*G. tomentella* D1, *G. tomentella* D3, *G. syndetika*, *G. canescens*, *G. clandestina*, *G. dolichocarpa*, *G. tomentella* T1 and *G. tomentella* T5) are green colored.

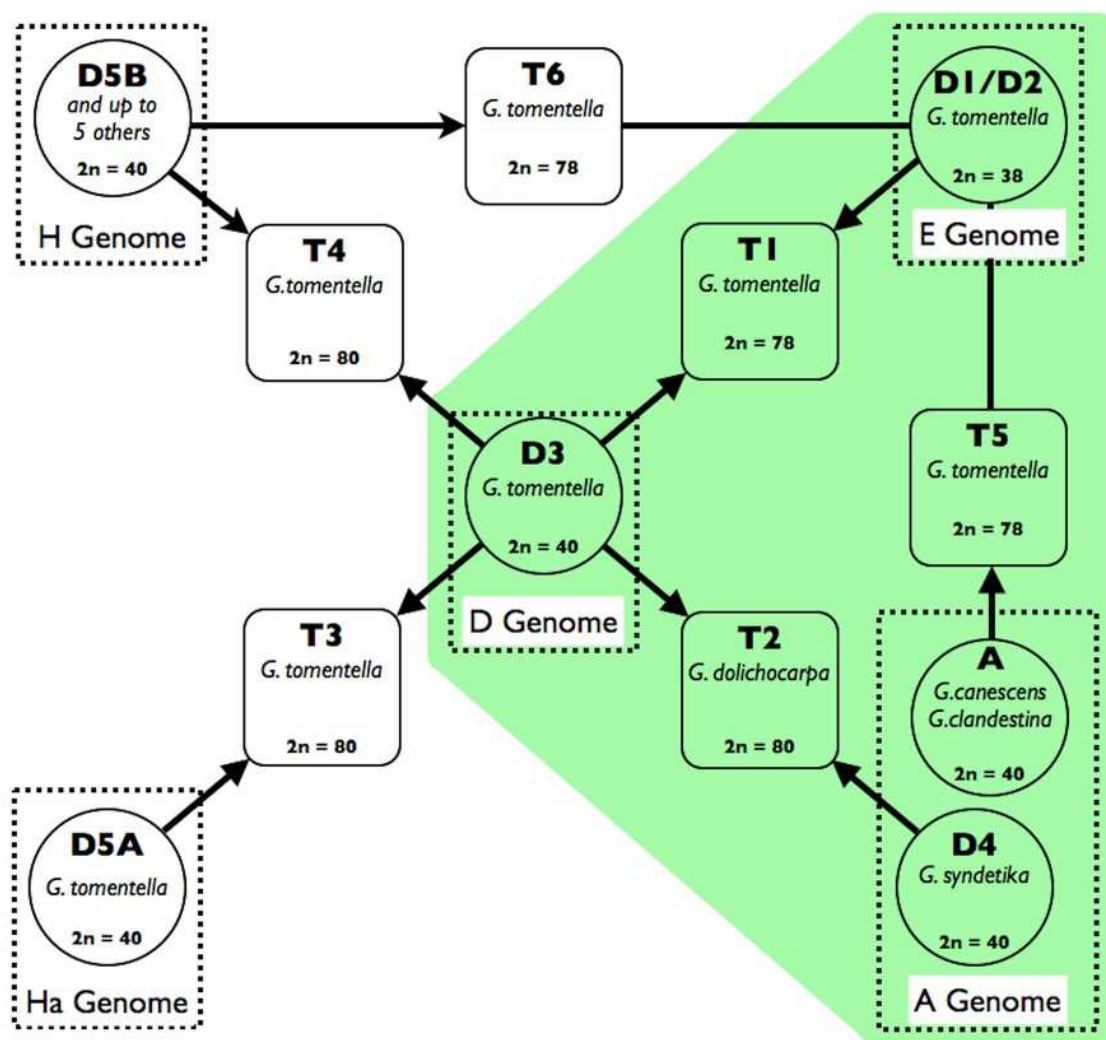


Figure 2

Electronic chromosome painting for *G. dolichocarpa* T2 accession 1134.

Electronic chromosome painting for *G. dolichocarpa* T2 accession 1134. SNP positions on the 20 soybean chromosomes are represented by blue lines (D3 progenitor) or red lines (D4 progenitor).

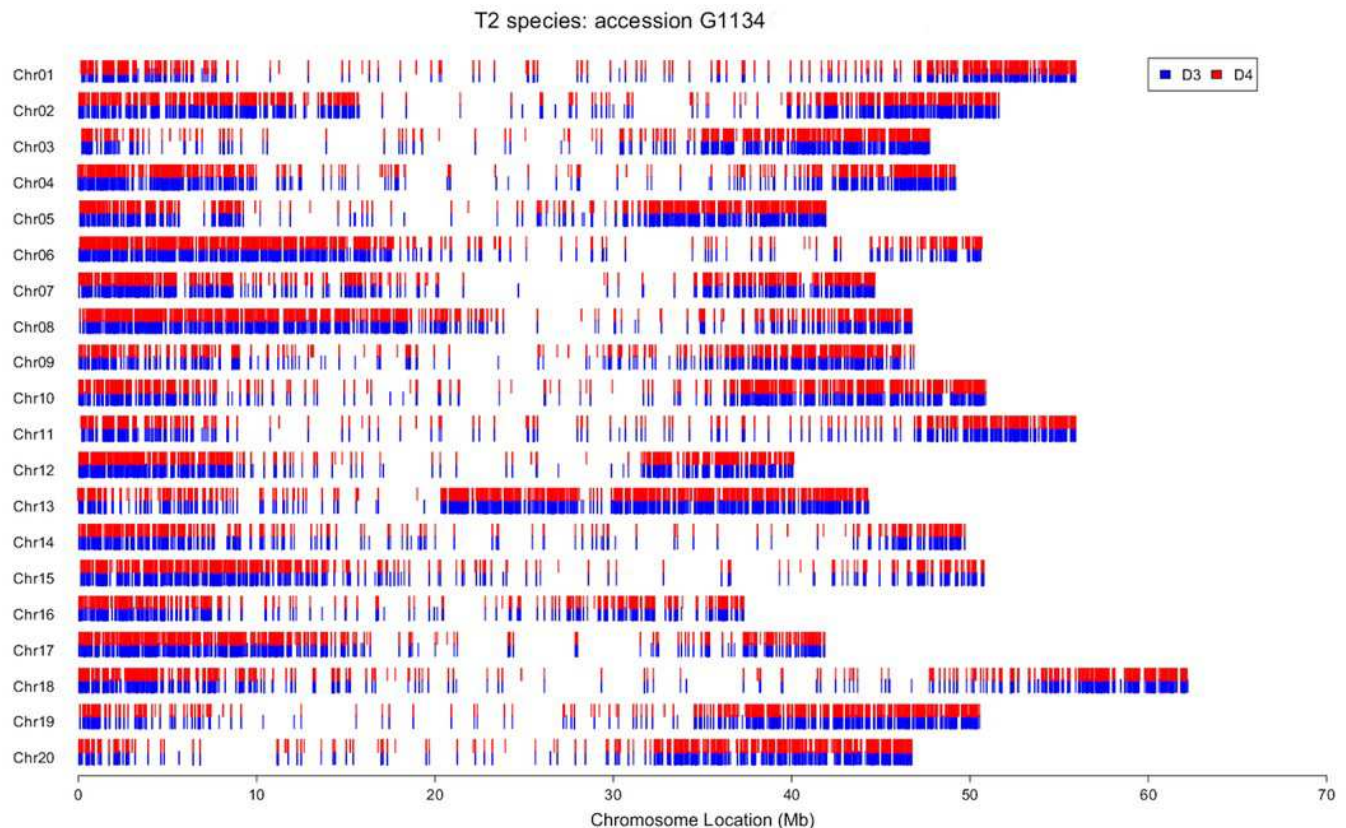


Figure 3

Structure analysis for *Glycine* perennial polyploid accessions

SNP analysis using Structure for a set of 20,000 random SNPs for *Glycine* polyploid complex accessions (A) without homoeolog separation and (B) with homoeolog separation for a K = 6. The five progenitor diploid species are represented by red (*G. clandestina*), dark red (*G. canescens*), yellow (*G. tomentella* D1), blue (*G. syndetika* D4) and green (*G. tomentella* D3).

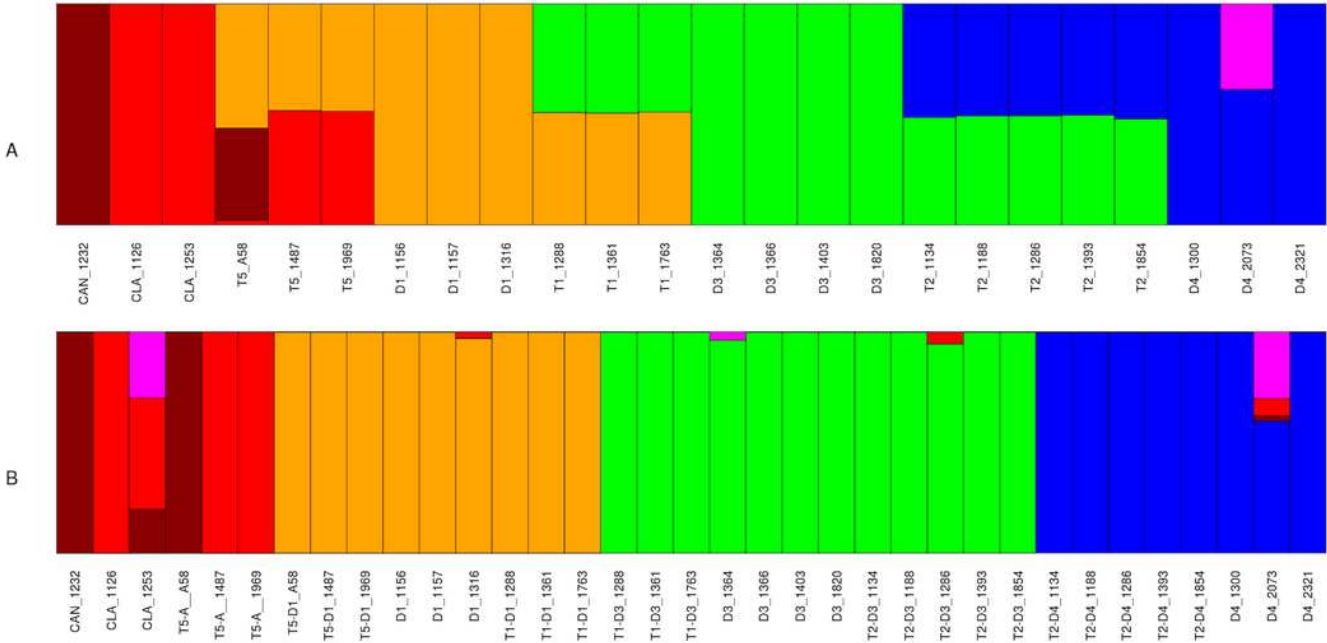


Figure 4

FineStructure Analysis for *Glycine* perennial polyploid species

Analysis using FineStructure for (a, c) 220,952 SNPs for the *Glycine* perennial polyploid complex (species groups A, D1, D3, D4, T1, T2 and T5) without homoeologue separation. 7 clusters can be distinguished (one per species group) in the PCA analysis where polyploids are admixtures of the diploid progenitor groups (a). The heatmap (b) shows diploid hybrid signal for polyploids, for example T5_A58 shows a stronger signal with its progenitors: CAN_1232 (blue) and D1_1316 (intense orange). (b, d) 70,910 SNPs for the *Glycine* perennials polyploid complex after homoeologue separation. 3 clusters can be distinguished in the PCA analysis (b): right cluster, species A and D4; bottom-left cluster, species D1; and top-left cluster, species D3. (d) The heatmap signal is divided into the same three major clusters.

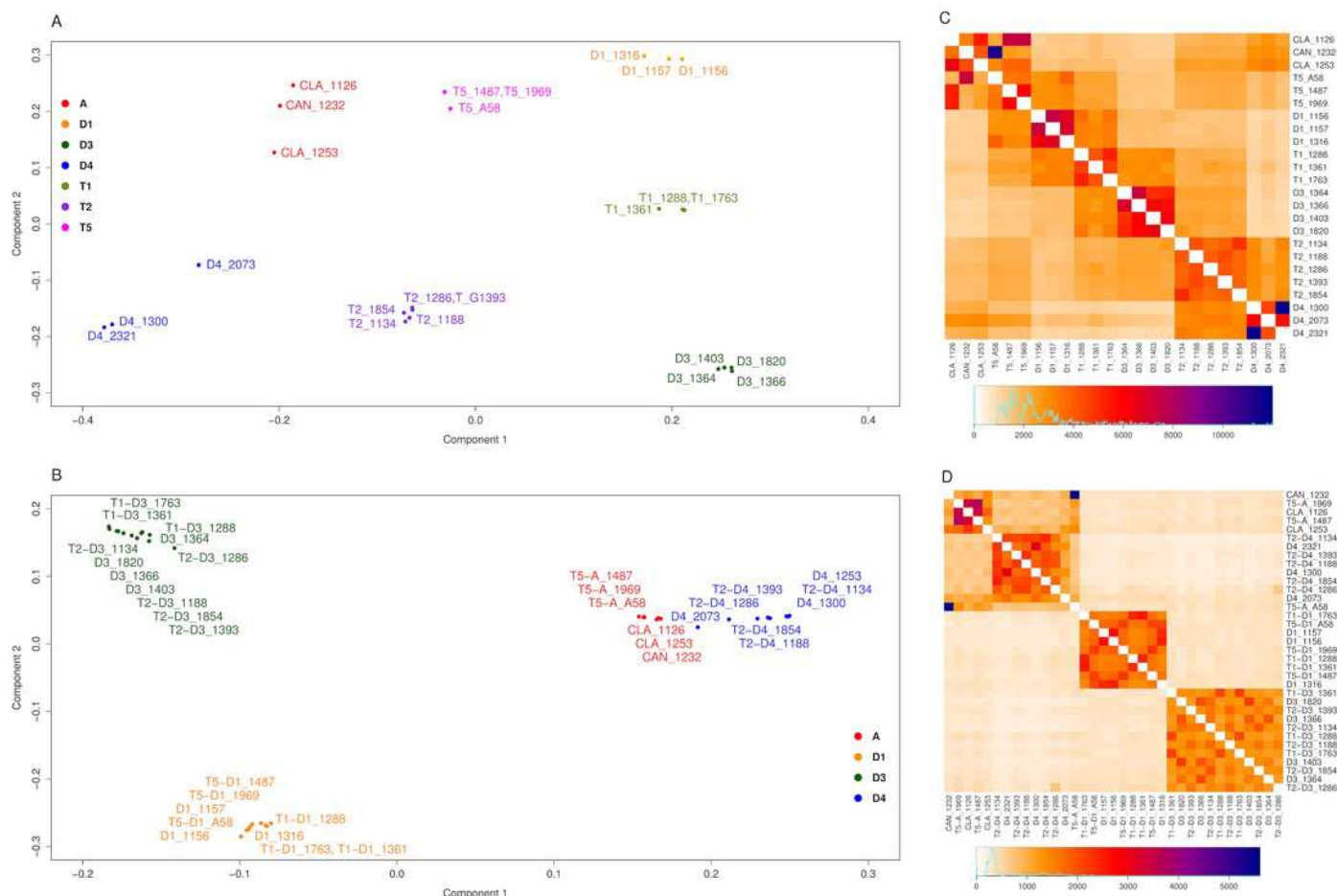


Figure 5

Phylogenetic relationship in the *Glycine* perennial polyploid complex

Phylogenetic relationships in the *Glycine* perennial polyploid complex after homoeologue separation. Branches are colored as in Fig. 4, based on the 5 different diploid species. In both ML (A) and NeighborNet (B), the same major species groups are visible (D1, D3, D4 and *G. canescens*/*G. clandestina*).

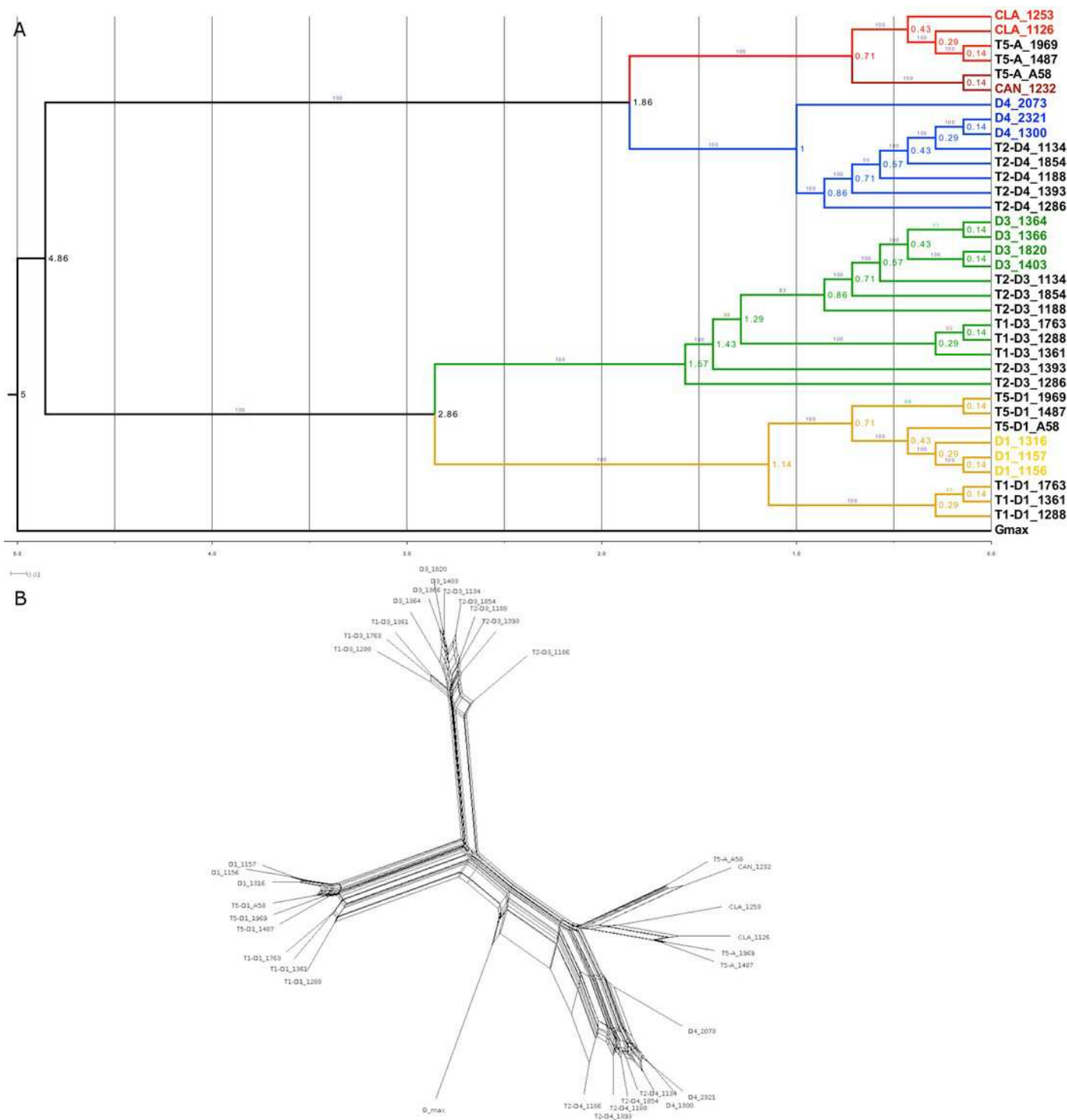


Figure 6

Phylogenetic analysis for the Glyma02g11580 locus

Phylogenetic analysis for the Glyma02g11580 locus using ML with bootstrapping values (A), NeighborNet (B) and BEAST (C) with posterior probabilities and showing node ages.

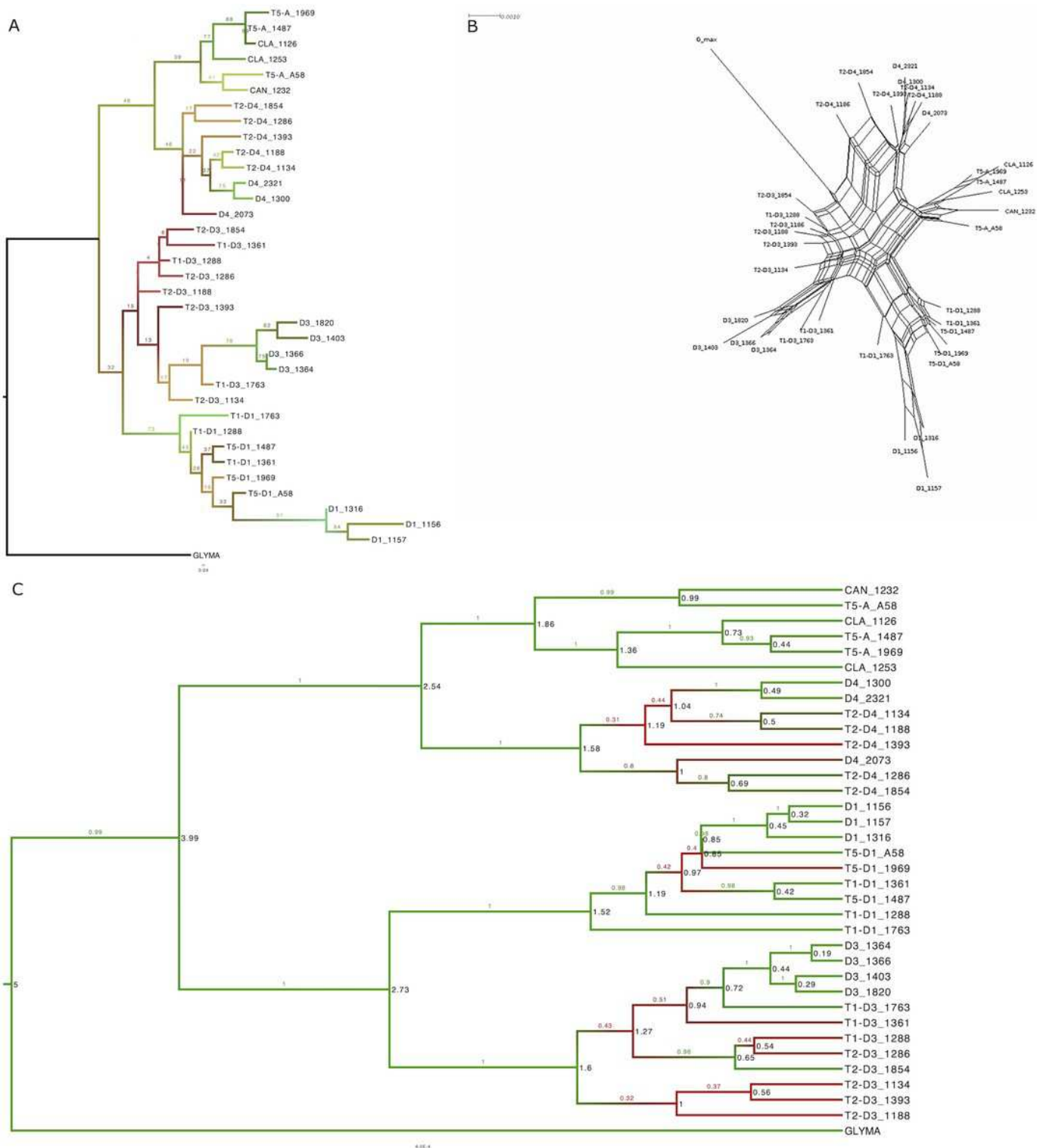


Figure 7

Phylogenetic tree with estimated divergence dates

*Beast tree with the estimated node ages and error bars representing the highest posterior density (HPD) interval for the 95% of the sampled values.

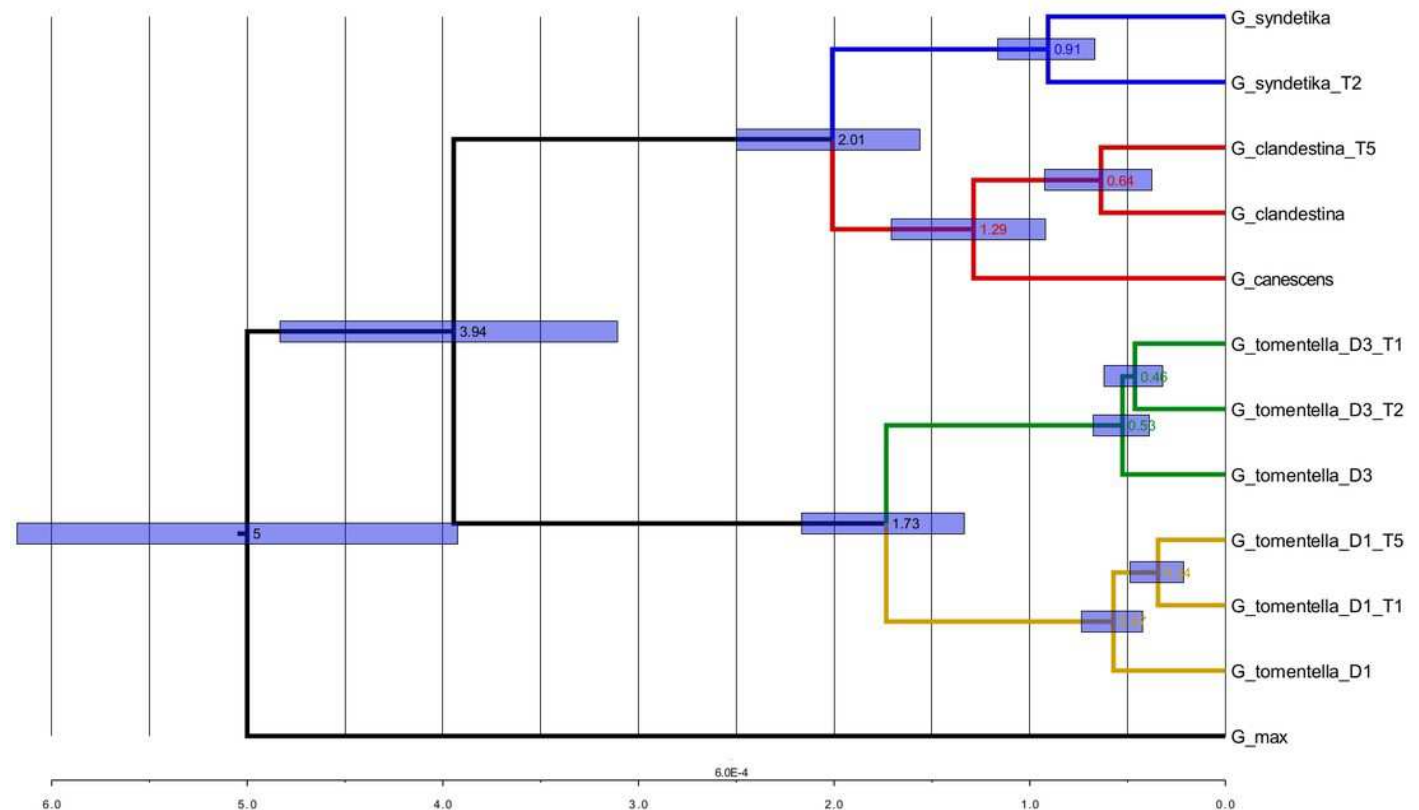


Table 1 (on next page)

Tables

Species	Accession	Samples	Raw Reads	Processed Reads	Mapped Reads	Represented Genes
<i>Glycine canescens</i>	1232	2	21332880	20696801	14381555	23833
<i>Glycine clandestina</i>	1126	2	19086864	18613018	11815996	23340
	1253	3	33546015	32326942	19117095	23723
<i>Glycine dolichocarpa</i>	1134	13	202427873	187120918	60712525	23643
	1188	2	19034633	18279858	11960713	22952
	1286	2	11814995	11216980	7422888	25278
	1393	2	21820163	21029983	13643602	23345
	1854	3	54748079	42826840	16032643	22718
<i>Glycine syndetika</i>	1300	3	25527322	23634740	14092961	24238
	2073	2	12132989	11072073	7087710	24438
	2321	2	32796391	30024544	13637368	22571
<i>Glycine tomentella</i> D1	1156	3	38218179	36988846	21905536	23041
	1157	2	16522541	15906072	9715890	23920
	1316	2	25207045	24375482	15417078	22749
<i>Glycine tomentella</i> D3	1364	1	10401944	9604350	6896983	22802
	1366	2	20631583	18098232	10766169	23364
	1403	3	31631369	28953234	17218424	23352
	1820	3	71185274	63055644	18625439	22871
<i>Glycine tomentella</i> T1	1288	2	14608219	14148847	9298349	23348
	1361	2	17964870	17627119	11217736	23758
	1763	2	21870236	20933661	14101838	23349
<i>Glycine tomentella</i> T5	A58_1	2	22447334	21996303	13389955	23042
	1487	2	21267274	20469069	13907305	23437
	1969	3	21324229	20847883	11136293	23522

Table 1: Sequencing, reads processing and mapping summary. Represented genes reflected the number of *Glycine max* reference genome genes where after the perennials reads mapping and expression measure have an expression > 0 (RPKM). Gray shading = allopolyploid species.

Species	Accession	% Gmax Coverage*	Raw SNPs	Processed SNPs**	Synonymous	Non-Synonymous
<i>Glycine canescens</i>	1232	7.2 [65.0]	589686	453,398 [7.7]	148321	123413
<i>Glycine clandestina</i>	1126	6.7 [61.4]	496746	375,943 [7.5]	115340	96562
	1253	7.5 [65.1]	617543	487,923 [8.3]	143952	124920
<i>Glycine dolichocarpa</i>	1134	11.6 [77.4]	1135676	965,643 [26.4]	242556	221326
	1188	7.4 [65.1]	550698	423,353 [28.9]	132471	113785
	1286	4.5 [45.5]	302661	224,653 [27.9]	67187	53595
	1393	6.7 [62.9]	470402	367,646 [28.8]	125339	104549
	1854	7.8 [65.3]	580531	471,020 [25.8]	140911	120274
<i>Glycine syndetika</i>	1300	7.5 [65.5]	605556	477,245 [7.8]	147362	125041
	2073	6.0 [57.6]	402798	282,215 [12.6]	91451	75333
	2321	8.0 [67.7]	670121	544,101 [6.3]	166409	143612
<i>Glycine tomentella</i> D1	1156	8.6 [69.7]	767614	621,043 [7.6]	190778	160781
	1157	6.2 [56.8]	455265	328,574 [7.1]	94377	77945
	1316	7.2 [62.3]	537439	412,518 [9.3]	120056	99666
<i>Glycine tomentella</i> D3	1364	5.0 [51.8]	335301	226,697 [7.8]	84917	65888
	1366	6.6 [59.7]	481258	360,327 [7.5]	111011	90015
	1403	6.4 [60.8]	476495	369,661 [6.6]	121526	99074
	1820	9.3 [69.6]	803774	641,145 [6.6]	188965	161826
<i>Glycine tomentella</i> T1	1288	6.9 [63.0]	498900	371,845 [19.6]	121418	102548
	1361	5.1 [54.2]	293339	200,738 [18.4]	75653	59378
	1763	7.1 [65.5]	533041	417,420 [19.4]	140203	116465
<i>Glycine tomentella</i> T5	A58_1	7.3 [64.6]	544331	430,552 [27.3]	135163	113781
	1487	7.0 [63.6]	516755	395,503 [26.8]	128199	105647
	1969	7.4 [65.9]	558920	444,468 [27.5]	146711	124933

Table 2: Summary of SNPs using *G. max* as reference genome. Gray shading = allopolyploid species.
 (* Between square brackets the coverage of the *G. max* transcriptome, including alternative splicings)
 (** Square brackets = percentage of heterozygous positions).

Species	Accession	Progenitor I	Mapped to Progenitor I (%)	SNPs for I*	Progenitor II	Mapped to Progenitor II (%)	SNPs for II*
<i>Glycine dolichocarpa</i>	1134	D3	11.4	399,884 [2.2]	D4	11.6	380,389 [2.1]
	1188	D3	20.8	227,610 [2.0]	D4	20.4	220,610 [2.1]
	1286	D3	20.3	124,984 [1.7]	D4	20.3	123,873 [1.8]
	1393	D3	19.6	197,132 [1.9]	D4	19.8	192,148 [1.9]
	1854	D3	17.9	245,354 [1.5]	D4	19.3	242,561 [1.7]
<i>Glycine tomentella</i> T1	1288	D1	14.9	143,232 [1.7]	D3	17.5	160,873 [1.9]
	1361	D1	15.0	155,360 [1.6]	D3	17.6	175,871 [2.0]
	1763	D1	14.8	158,777 [1.8]	D3	17.3	179,032 [2.0]
<i>Glycine tomentella</i> T5	A58_1	A	16.9	190,138 [2.1]	D1	20.5	222,134 [1.8]
	1487	A	17.1	174,051 [1.9]	D1	20.0	202,555 [1.7]
	1969	A	16.0	182,615 [2.4]	D1	18.6	214,799 [1.8]

Table 3: Summary of the mapped reads and SNPs produced after the homoeologous reads based in the selective mapping with its progenitors (* Square brackets = percentage of heterozygous positions).

Species Group	Gmax SNPs	A Group SNPs	D1 Group SNPs	D3 Group SNPs	D4 Group SNPs
A Species	9406	26,438 *	7096	6591	1465
D1 Species	11187	-	21,830 *	5933	7556
D3 Species	9299	-	-	25,157 *	7295
D4 Species	9314	-	-	-	23,324 *

Table 4: Summary of SNP count between species groups (polyploids are divided in two species according the progenitor origin). A Species includes *G. canescens*, *G. clandestina* and *G. tomentella* T5-A; D1 species includes *G. tomentella* D1, *G. tomentella* T1-D1 and *G. tomentella* T5-D1; D3 species includes *G. tomentella* D3, *G. tomentella* T1-D3 and *G. tomentella* T2-D3; D4 species includes *G. syndetika* and *G. tomentella* T2-D4. * The same species group contains the specific SNPs between accession of the same species.

GeneID	TreeLikelihood Mean	TreeLikelihood ESS	Gene Functional Annotation
Glyma01g35620	-4676.067	1361.926	Phytoene dehydrogenase
Glyma02g11580	-3986.812	1034.414	RNA binding protein
Glyma03g29330	-7611.686	894.813	Magnesium chelatase
Glyma03g36630	-2666.725	696.197	Rho GTPase activating protein
Glyma04g39670	-4142.157	2556.251	ATP-binding transport protein-related
Glyma05g05750	-3028.809	541.483	Beta-amylase
Glyma05g09310	-2578.198	305.191	Pyruvate kinase
Glyma05g26230	-3741.498	5156.337	Metalloprotease M41 FtsH
Glyma05g37840	-2138.404	3766.767	Haloacid dehalogenase-like hydrolase
Glyma06g18640	-3418.91	6613.752	Elongation factor Tu
Glyma07g03370	-2091.845	742.796	Palmytoil-monogalactosyldiacylglycerol delta-7 desaturase
Glyma07g17180	-2218.934	2833.162	Fructose-1,6-bisphosphatase
Glyma10g42100	-2903.849	1774.192	3-ketoacyl-CoA synthase
Glyma11g13880	-4644.419	7487.252	Lipoxygenase
Glyma11g33720	-3592.689	2273.389	DELLA protein
Glyma12g04150	-2061.527	4777.335	Fructose-bisphosphate aldolase
Glyma12g12230	-2177.72	1790.185	O-methyltransferase
Glyma13g17820	-2439.715	342.999	Polyubiquitin
Glyma14g03500	-2063.541	819.216	Phytoene synthase
Glyma16g00410	-4041.294	4475.536	heat shock protein 70
Glyma16g01980	-4985.489	387.993	Myb-like protein
Glyma16g04940	-2152.355	4586.42	Glyceraldehyde 3-phosphate dehydrogenase
Glyma18g04080	-2285.776	9535.754	26S proteasome regulatory complex, ATPase RPT4
Glyma19g03390	-2344.831	3190.5	Unknown
Glyma19g32940	-2176.029	2579.558	Fatty acid desaturase
Glyma20g24930	-2803.585	6535.602	3-ketoacyl-CoA synthase
Glyma20g32930	-2867.321	2078.549	Cytochrome P450 77A3

Table 5: Summary of the genes used in the BEAST and *BEAST analysis with the tree likelihood values and the functional annotation.