

1 **Running head:**

2 **PRINCIPAL TREE INFERENCE FROM GENOMIC BLOCKS**

3

4 **Title:**

5 **Using genomic location and coalescent simulation to investigate gene tree discordance in**

6 ***Medicago* L.**

7

8 Sousa, F.^{1,‡}, Bertrand, Y.J.K.^{1,‡}, Doyle, J.J.², Oxelman, B.¹ and Pfeil, B.E.^{1,*}

9

10 ¹Department of Biological and Environmental Sciences, University of Gothenburg, Box 461, 40530
11 Gothenburg, Sweden

12 ²Department of Plant Biology, Cornell University, 404 Mann Library Building, Ithaca, NY 14853,
13 USA

14

15 *Corresponding author: bernard.pfeil@gu.se

16 ‡Equally shared contributions

17

18 **ABSTRACT**

19 Several well-documented evolutionary processes are known to cause conflict between
20 species level phylogenies and gene level phylogenies. Three of the most challenging processes for
21 species tree inference are incomplete lineage sorting, hybridisation and gene duplication, which
22 may result in unwarranted comparisons of paralogous genes. Several existing methods have dealt
23 with these processes but none has yet been able to untangle all three at once. Here, we propose a
24 step-wise method by which these processes can be discerned using information on genomic location
25 coupled with coalescent simulations. In the first step, highly discordant genes within genomic
26 blocks (putative paralogues) are identified and excluded from the dataset and, in the second step,
27 blocks of linked genes are grouped according to their hybrid history. Existing multispecies
28 coalescent software can then be applied to recover the principal tree(s) that make up the species
29 tree/network without violating the underlying model. The potential of the approach is evaluated on
30 simulated data derived from a species network composed of nine species, of which one is of hybrid
31 origin, and displaying a single gene duplication that leads to paralogous comparisons. We apply our
32 method to an empirical set of 12 genes from seven species sampled in the plant genus *Medicago*
33 that display phylogenetic discordance. We identify the causes of the discordance and demonstrate
34 that the *Medicago orbicularis* lineage experienced an episode of ancient hybridisation. Our results

35 show promise as a new way to explore phylogenetic sequence data that can significantly improve
36 species tree inference in presence of hybridisation and undetected paralogy or other causes leading
37 to extremely discordant gene trees.

38

39 **KEYWORDS:** phylogenetic incongruence, species tree, gene tree, coalescent simulation, genomic
40 location, paralogy, hybridisation, incomplete lineage sorting, principal tree

41

42 Species tree inference from molecular data has, in recent years, substantially shifted away
43 from the concatenation of sequences from several genes to approaches based on the multispecies
44 coalescent (Rannala and Yang, 2003; Degnan and Rosenberg, 2009). The multispecies coalescent
45 assumes that each gene has its own phylogenetic history (Kingman, 1982) conditioned on the
46 common underlying species phylogeny. However, several evolutionary processes may induce
47 significant differences in phylogenies from unlinked genes, of which three of the best understood
48 are: incomplete lineage sorting (ILS), homoploid hybridisation and paralogy (Maddison 1997).

49 Incomplete lineage sorting occurs when alleles coalesce more deeply than species
50 divergence, as a consequence of stochastic differential assortment of ancestral polymorphism (e.g.,
51 Takahata and Nei 1985; Pamilo and Nei 1988; Maddison 1997). Although the effect of ILS on
52 phylogenies can be profound (Degnan and Salter 2005), especially in sexually reproducing
53 organisms with large populations and when the time between speciation events is short (Pamilo and
54 Nei 1988; Nei and Kumar 2000), ILS itself is assumed to be ubiquitous, as an unavoidable
55 consequence of neutral population processes. Gene tree incongruence has been explained by ILS in
56 several studies that have applied coalescent-based methodologies (Pollard et al. 2006; Carstens and
57 Knowles 2007; Syring et al. 2007; Blanco-Pastor et al. 2012).

58 Homoploid hybridisation, the formation of hybrids between different species at the same
59 ploidy level, can result in the introgression of alleles from one species into another, or even in the
60 formation of new species that contain alleles from both parents (e.g., Buerkle and Rieseberg 2008;
61 Hermansen et al. 2011; Bosse et al. 2014). Hybridisation has been widely documented in animals
62 and plants (Rieseberg 1997; Mallet 2007), with several studies demonstrating introgression of
63 alleles or single nucleotide polymorphisms (SNPs) while taking into account ILS (Buckley et al.
64 2006; Peters et al. 2007; Blanco-Pastor et al. 2012; Phillips et al. 2013; Ramadugu et al. 2013; Good
65 et al. 2015).

66 Genes are frequently subjected to duplication (the formation of paralogues) and loss. When
67 paralogues from different species are mistaken for orthologues, their comparison may result in
68 invalid inferences regarding species relationships (e.g., Doyle 1992; Oxelman et al. 2004).

69 Genealogical discordance due to paralogy can be further complicated by concerted evolution
70 (Sanderson and Doyle 1992; Nei and Rooney 2005). Many examples drawn from eukaryotic species
71 show that both tandem and whole genome duplications are very common (e.g., Lynch and Connery
72 2000; Cui et al. 2006; Chen et al. 2013; Van Zee et al. 2016) and can be followed by rapid gene loss
73 (Blanc and Wolfe 2004). Tandem duplications can be extensive in some diploid organisms, even
74 ones with relatively small genomes, such as *Arabidopsis thaliana* (Blanc and Wolfe 2004). The
75 frequent absence of a full set of duplicated genes, either due to gene loss or sampling bias, fosters
76 unintentional paralogous comparisons. We denote unrecognized instances of paralogy as “paralogy-
77 affected” genes throughout the paper.

78 The phylogenetic effect of each of these three confounding processes has been widely
79 studied (e.g., Maddison and Knowles 2006; Edwards et al. 2007; Holland et al. 2008; Maureira-
80 Butler et al. 2008; Joly et al. 2009; Meng and Kubatko 2009; Joly et al. 2010; Heled and
81 Drummond 2010; Jones et al. 2013; Yu et al. 2013). Several methods implement species tree
82 inference that accommodates ILS, by using the multispecies coalescent (e.g., BEST, Liu and Pearl
83 2007; STEM, Kubatko et al. 2009; *BEAST, Heled and Drummond 2010; MP-EST, Liu et al.
84 2010; PhyloNet, Liu et al. 2014; ASTRAL, Mirarab et al. 2014), with two of these incorporating
85 hybridisation in species tree inference (STEM-hy, Kubatko 2009, and PhyloNet). Another approach
86 that infers species trees despite incongruence, but without modelling any specific causal process, is
87 BUCKy (Ané et al. 2007). Other approaches determine whether ILS alone can explain the observed
88 gene tree incongruence, and thus infer hybridisation on a case-by-case basis rather than by co-
89 estimation (Buckley et al. 2006; Maureira-Butler et al. 2008; the JML method, Joly et al. 2009; the
90 ABBA-BABA test, Green et al. 2010; Bertrand et al. 2015).

91 Models that accommodate gene duplications in a coalescent framework have been proposed
92 (Rasmussen and Kellis 2012), but only for cases involving the presence of both duplicates at least in
93 some gene trees. However, no method, step-wise or simultaneous, has yet accommodated all three
94 sources of gene tree incongruence (ILS, hybridisation and paralogy). Instead, each method handles
95 only one or at most two of these sources of incongruence, and requires the assumption that the
96 other(s) are not operating (or do not affect the outcome). Given that each of these sources of
97 incongruence, when not handled appropriately, can mislead species tree inference (e.g., Oxelman et
98 al. 2004; Kubatko and Degnan 2007; Reid et al. 2012), it is critical that all three can be diagnosed
99 when acting in concert.

100 The legume genus *Medicago* L. (Fabaceae) exhibits severe phylogenetic incongruence
101 among nuclear and chloroplast markers (Bena 2001; Steele and Wojciechowski 2003; Maureira-
102 Butler et al. 2008; Steele et al. 2010; Yoder et al. 2013; Sousa et al. 2014; Sousa et al. 2016) that

103 has been attributed, at least partially, to hybridisation (Maureira-Butler et al. 2008; Yoder et al.
104 2013; Sousa et al. 2016), although paralogy has not been formally excluded. A survey of published
105 *Medicago* phylogenies indicated that *M. orbicularis* is one of the taxa putatively involved in ancient
106 hybridisation, possibly involving the lineages that include *M. truncatula*, *M. ciliaris* and *M. arabica*
107 (Sousa et al. 2016). The present study investigates gene tree incongruence involving the
108 phylogenetic position of *M. orbicularis* using a method that untangles topological variation caused
109 by ILS, paralogy and hybridisation. We explore the hypothesis that genomic location can be used to
110 identify possible cases of paralogy by comparing groups of genes that are physically associated
111 (tightly linked), and to recognize hybridisation by comparing groups of genes that are physically
112 distant (unlinked). We use target-enrichment methods (Guschanski et al. 2013; Smith et al. 2014;
113 Stull et al. 2013; Sousa et al. 2014) to obtain sequences of multiple loci from each sampled
114 specimen.

115 To evaluate the performance of our method, we simulate data under a species network in the
116 presence of paralogy, hybridisation and ILS. We then (1) identify genes likely to be paralogy-
117 affected and remove them from further analyses, (2) group genes with a shared parental origin (due
118 to hybrid parentage) and (3) reconstruct species phylogenies, while accounting for ILS, and confirm
119 the different parental origins of taxa with hybrid histories through multi-species coalescent
120 inference. Both the paralogy detection and hybridisation detection steps are based on the test
121 developed by B.E.P. and described in Maureira-Butler et al. (2008) that assumes a null hypothesis
122 of ILS and derives null distributions based on gene trees used as surrogates for the unknown species
123 tree. We show that likely cases of paralogy in *Medicago* that could mislead species tree inference
124 are successfully detected. We confirm that a hybridisation signal can be discerned in *Medicago*
125 *orbicularis* only if genomic location is taken into account.

126

127 **METHODS**

128 SAMPLING, ALIGNMENT ASSEMBLY AND GENE TREE INFERENCE

129 We sampled seven species of *Medicago* (*M. arabica* (L.) Huds., *M. ciliaris* (L.) All., *M.*
130 *granadensis* Willd., *M. intertexta* (L.) Mill., *M. littoralis* Rohde ex Loisel., *M. orbicularis* (L.)
131 Bartal., *M. truncatula* Gaertn.) and two outgroups (based on e.g., Steele et al. 2010), *Melilotus*
132 *neapolitanus* Ten. and *Mel. sulcatus* Desf. We used 12 markers (Supp. Table 1), of 2-3 Kb in
133 length, located in four unlinked blocks of three genes (hereafter referred to as blocks A, B, C and
134 D), with < 30 kb span between genes inside a block and each block located on a different
135 chromosome, based on the *M. truncatula* genome (Young et al. 2011) and previously tested for
136 phylogenetic utility (Sousa et al. 2014). We considered that four blocks of three genes each was an

137 adequate dataset to explore the use of genomic location for the identification of reticulation using
138 our newly proposed method. The 12 genes selected showed, in a preliminary gene tree analysis,
139 evidence for alternative positions of *Medicago orbicularis*. DNA extraction, target-enrichment and
140 Illumina sequencing were performed as reported previously (Sousa et al. 2014). Raw sequence
141 reads were submitted to the European Nucleotide Archive repository with the following accession
142 numbers: *Medicago arabica*: ERS719973, *M. ciliaris*: ERS719974, *M. intertexta*: ERS719975, *M.*
143 *granadensis*: ERS719976, *M. littoralis*: ERS719977, *M. orbicularis*: ERS719978, *Melilotus*
144 *sulcatus*: ERS511667, *Melilotus neapolitanus*: ERS511668. *Medicago truncatula* sequences were
145 obtained from the reference genome at www.medicagohapmap.org.

146 CLC Assembly Cell v.4.0.13 software (CLC Bio, Aarhus, Denmark) was used to remove
147 adapter sequences and filter reads for quality, with a phred-score threshold of 20. CLC-mapper was
148 employed to map reads from individual samples to the genomic reference sequence for each gene.
149 Alleles were phased using the program Samtools phase (Li et al. 2009) and phased reads were re-
150 assembled into contigs using CLC-assembler. Sequences were aligned in Geneious Pro (v.5.3.6),
151 using the Geneious alignment tool with default parameters and five iterations. Alignments were
152 tested for recombination with RDP4 (Martin et al. 2010). Recombination events were accepted
153 when detected by at least two methods at P-values <0.01 (with Bonferroni correction) in presence
154 of topological conflict; the corresponding recombinant fragments were removed.

155 In order to estimate appropriate substitution models, we analysed the alignments, previously
156 filtered for recombination, in MrBayes v.3.1.2 (Huelsenbeck and Ronquist 2005) using model
157 averaging over the GTR model family (Huelsenbeck et al. 2004) and gamma rate variation. Each
158 analysis was run for 3M generations, mixing and ESS values were controlled in Tracer v1.6.0
159 (Rambaut and Drummond 2007). A run was deemed successful if examination of the parameter
160 files in Tracer showed no indications of a lack of convergence and displayed high ESS values
161 (above 100) for all parameters of interest. The “sump” function in MrBayes was then used on each
162 run, with a specified burn-in (specific to each gene based on the Tracer results, not reported here),
163 to verify which substitution model had the highest posterior probability. Ultrametric gene trees were
164 then obtained in BEAST v. 1.8.0 (Drummond and Rambaut 2007) using the inferred substitution
165 model and gamma rate variation, an uncorrelated lognormal relaxed clock with a ucl.d.mean rate
166 normal prior with mean= 3.6×10^{-9} (Sousa et al. 2014) and a normal root age prior with
167 mean= 15.9×10^6 (Lavin et al. 2005; Maureira-Butler et al. 2008).

168

169 UNDERLYING ASSUMPTIONS FOR DETECTING ILS, HYBRIDISATION AND PARALOGY

170 *i) Chromosomal blocks become fixed in a hybrid genome*

171 We assume that lineages that have undergone hybridisation eventually fix large
172 chromosomal blocks (= genomic blocks) derived from each parental lineage. Once blocks are fixed,
173 i.e. become homozygous, recombination acts only on segments with the same parental origin.
174 Coupled with the emergence of new mutations, this causes alleles at each locus within fixed blocks
175 to accumulate differences and potentially undergo independent ILS with respect to speciation events
176 subsequent to the hybridisation event. Despite the effect of ILS, for each gene tree within a block,
177 the alleles will still display a closer relationship to the parental lineage of the block, rather than to
178 the other parental lineage (i.e., of other genomic blocks). Each alternative parental relationship is
179 henceforth referred to as a principal tree (PT) (Holland et al. 2008). The result of a single
180 hybridisation episode and subsequent fixation of genomic blocks is a genome that is a mosaic of
181 blocks, and within a block all the genes are evolving according to one of two principal trees (PT1 or
182 PT2).

183 *ii) Fixed genomic blocks are typically large*

184 The two main mechanisms that produce mosaic genomes (homoploid hybrid speciation
185 (HHS) and introgression) are expected to promote the rapid fixation of some genomic blocks: by
186 population bottlenecks and drift (e.g., Buerkle and Rieseberg 2008) for HHS, and by selection in
187 the case of introgression, where the majority of introgressed regions are lost due to drift, but a few
188 that carry adaptive mutations may be rapidly fixed by selection (e.g., Liu et al. 2015). Genomes
189 containing large hybrid blocks have been observed in cases of introgression as well as HHS. The
190 introgression of anticoagulant-resistance alleles into some house mouse populations from the
191 Algerian mouse have resulted in a few large introgressed blocks (some >10 Mb) (Song et al. 2011;
192 Liu et al. 2015). The average genomic block size of recently introgressed (c. 200 years ago) Asian-
193 descendant alleles following hybridisation of Asian domestic pig into a background of European
194 breeds is over 100 000 bp (Bosse et al. 2014), although fixation has not occurred at every locus.
195 Another case where the results of introgression have not completely stabilised is that of Neanderthal
196 alleles into modern humans, as blocks with a median of 129 000 bp have been inferred to occur
197 (Sankararaman et al. 2014). Three homoploid hybrid North American sunflower species display
198 large blocks (Ungerer et al. 1998; Rieseberg et al. 2003), and the homoploid hybrid Italian sparrow
199 has inherited most of its Z chromosome from the house sparrow (Trier et al. 2014).

200 *iii) Genomic rearrangements do not disrupt small genomic blocks*

201 Gene synteny within small genomic blocks is unlikely to be disrupted by genomic
202 rearrangements, which are relatively rare. Large, chromosome-scale translocations or inversions
203 will carry entire genomic blocks intact, and so can also be ignored for our purposes.

204 *iv) Paralogous retention outlasts polymorphism retention*

205 Duplicated gene copies can arise within genomic blocks, e.g., by tandem duplication (e.g.,
206 Innes et al. 2008). We expect that when paralogues are maintained for long periods, and loss is
207 random in different lineages, the corresponding gene topologies will be more discordant than
208 expected by ILS alone (e.g., see Supp. Fig. 1). Indeed, paralogous gene copies can persist for much
209 longer than selectively neutral allelic polymorphisms. For example, many duplicated genes (20-
210 50% of total gene content) originating by whole genome duplication (WGD) can be preserved for
211 hundreds of millions of years (Lynch and Force 2000; Maere et al. 2005). In the legume *Glycine*,
212 77% of low copy genes in a 1 Mb duplicated region have been retained for at least 5 Ma and maybe
213 as long as 10 Ma (Innes et al. 2008). In contrast, the typical duration for ancestral polymorphism
214 retention ($< 5 \times Ne$ generations, Rosenberg 2003) is much less than 5 Ma under typical plant
215 generation times and effective population sizes ($Ne \times$ generation time is usually < 1 million:
216 Maureira-Butler et al. 2008; Strasburg and Rieseberg 2008; Foxe et al. 2009; Lundemo et al. 2009;
217 Gossman et al. 2010; Blanco-Pastor et al. 2012; Ramadugu et al. 2013; Pérez-Collazos et al.
218 2015). However, discordance caused by paralogy can also be minimal and approach the level of
219 incongruence obtained through ILS. Such cases of paralogy are difficult to detect but are also
220 unlikely to create problematic violations of the multispecies coalescent model (see Results). Thus,
221 paralogy identification is most important when phylogenetic discordance is higher than that
222 expected due to ILS, but is also most easily uncovered under these conditions.

223 It should be noted that any other cause of very discordant gene trees that operates on a single
224 gene at a time (and not on entire blocks), whether biological (e.g. recombination, base composition
225 bias) or methodological (e.g. model misspecification, alignment errors) could be mistaken for
226 paralogy by our method. Therefore, an inference of a “paralogy-affected” gene should be treated
227 with caution. However, the use of long informative alignments, recombination testing, quality
228 control in laboratory procedures, dense taxon sampling, among other things, will reduce the chance
229 that other causes will produce extremely discordant gene trees, leaving gene duplication and loss
230 dynamics as (arguably) the most likely contender. For the purposes of describing our method, we
231 assume that best practice procedures have been followed at every stage of analysis, and thus use the
232 inference of “paralogy-affected” genes as a shorthand to indicate that this is the most likely cause.
233 However, we acknowledge that other causes cannot be ruled out given our methodology.

234

235 OUTLINE OF THE METHOD TO DETECT ILS, HYBRIDISATION AND PARALOGY

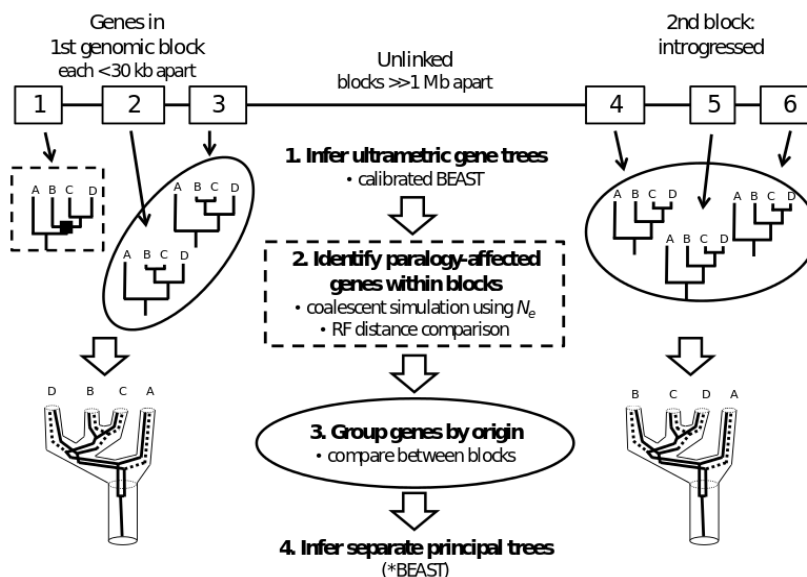
236 The proposed method consists of four main steps: 1) Gene tree inference; 2) Paralogy test -
237 identification and exclusion of highly-discordant trees within genomic blocks; 3) Hybridisation test
238 - identification of hybridisation between genomic blocks and grouping of blocks sharing the same

239 parental history; 4) Inference of principal trees. Steps 1 and 4 rely on established methods for tree
 240 inference. Steps 2 (Paralogy test) and 3 (Hybridisation test) as well as the coalescent simulations
 241 tests used for gene pairwise comparisons, are explained below.

242

243 *Paralogy Test*

244 We considered a sampling scheme where each genomic block has three genes (triplet) (Fig.
 245 1). Pairwise comparisons (Supp. Fig. 1) were performed between the observed gene trees within
 246 each genomic block, using the coalescent simulation test, for a total of three comparisons (Supp.
 247 Fig. 2). A gene tree was inferred as likely to be affected by paralogy if ILS was rejected in both
 248 pairwise comparisons with the other two gene trees of the triplet (Fig. 1, also Supp. Fig. 2). All
 249 genes identified as paralogy-affected were removed from further analysis. When ILS was rejected
 250 in all three pairwise comparisons, the whole block was excluded, because this case does not allow
 251 the identification of the paralogy-affected gene nor does it produce any remaining compatible pairs
 252 of genes. In contrast, if the block was not excluded, genes that differ by ILS alone were further
 253 subjected to the hybridization test.



254

255 Figure 1.

256

257 *Coalescent Simulation Tests*

258 All pairwise comparisons used in our method are based on a parametric bootstrapping

259 procedure that compares the topological distance between two observed gene trees with the
260 topological distances expected, under a coalescent process, between the simulated gene trees from
261 two modelled species trees (Maureira-Butler et al. 2008; Blanco-Pastor et al. 2012; Ramadugu et al.
262 2013). In the original implementation (Maureira-Butler et al. 2008), the observed distance was
263 calculated directly from a point estimate of each gene tree, and thus did not account for tree
264 inference uncertainty. In order to take this source of uncertainty into account, we instead derived a
265 distribution of distances from a posterior distribution of Bayesian trees. All tree distances in this test
266 are calculated using the unweighed Robinson-Foulds (RF) metric (Robinson and Foulds 1981),
267 a.k.a symmetric distance.

268 The test can be summarised as follows: Let G_1 and G_2 be two observed gene trees, and let D_1
269 and D_2 be the respective post burn-in gene tree posterior distributions. The topological distance
270 between G_1 and G_2 (observed gene tree distances: Supp. Fig 2) is represented by a distribution, with
271 a 95% credibility interval, of the distances between 100 trees drawn randomly from D_1 (D_1^1 to
272 D_1^{100}) and 100 trees drawn randomly from D_2 (D_2^1 to D_2^{100}) – 100 x 100 RF distances. As the true
273 species trees that contain G_1 and G_2 are unknown, they are modelled by treating G_1 and G_2 as
274 surrogate species trees. From each tree drawn from D_1 (D_1^1 to D_1^{100}) and from D_2 (D_2^1 to D_2^{100}), 100
275 gene trees are simulated, under the multispecies coalescent, using Kingman's neutral coalescent
276 process, resulting in simulated sets, S_1^1 to S_1^{100} and S_2^1 to S_2^{100} , each with 100 trees. The ILS null is
277 obtained by calculating the RF distance between each tree drawn from D_1 and D_2 and the respective
278 simulated set ($D_1^1 \times S_1^1, D_1^2 \times S_1^2, \dots, D_1^{100} \times S_1^{100}$), generating two distributions (for G_1 and G_2) of
279 100 x 100 RF distances each (referred to as coalescent distance distributions, see Supp. Fig 2). Gene
280 tree distances are then used as a test statistic for an upper-tailed test with a given critical value. We
281 define the p% critical value to be the value of the simulated null distributions that (1-p)% of the
282 simulated values exceed. The null hypothesis of ILS is rejected if the test statistic (gene tree
283 distances) exceeds the p% critical value of both coalescent distance distributions.

284 The effect of population sizes in simulations for a similar coalescent-based test was assessed
285 in Bertrand et al. (2015), where it was found that the performance of the test improved when the
286 population sizes used for simulating the gene trees were significantly smaller than the real
287 population size. One explanation that can account for this observation is that gene lineages have
288 coalescent times extending farther back than species lineages (Maddison and Knowles 2006).
289 Consequently, using a gene tree as a surrogate species tree for gene tree simulation will tend to
290 further push back in time the divergence ages on the simulated trees. It ensues that the use of a
291 surrogate tree introduces more sequence variation in the simulations than what would be obtained in
292 an actual species tree. Instead of evaluating the population size needed for satisfactory ILS

293 acceptance/rejection rates in our coalescent simulation test, we opted for changing the critical value
294 of the one-tailed test, and thus performed it with p% critical values of 65%, 75%, 85% and 95%.

295

296

297 *Hybridisation Test*

298 After the exclusion of putative paralogy-affected genes, we tested for gene tree discordance
299 between genomic blocks (Fig. 1). The coalescent simulation test was used for pairwise comparisons
300 between genes from different blocks (maximum of 3 x 3 gene comparisons). We accept
301 hybridisation as the source of the observed discordance only when ILS was rejected in at least two
302 pairwise comparisons between genes from different blocks. Other decision rules were tested
303 (accepting hybridisation if ILS was rejected in at least one pairwise comparison, or only if all
304 pairwise comparisons rejected ILS) but these proved to have worse rates of type I or type II error
305 (not shown). Genes from genomic blocks that share the same PT were pooled for species tree
306 inference (Fig. 1).

307

308 VALIDATION OF THE METHOD

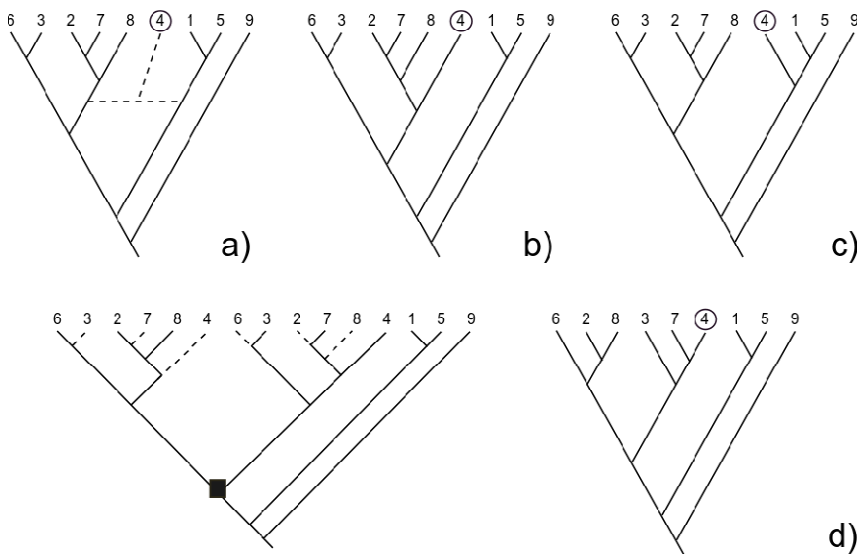
309 *Data Simulation*

310 The simulation protocol is outlined in Figure 2 and Supplementary Figure 3. We used a
311 simulated nine-taxon network, including one species of hybrid origin and one gene duplication (Fig.
312 2a). Principal tree 1 (PT1) was generated in Mesquite v.2.75 (Maddison and Maddison 2006),
313 assuming a total tree depth of 2 000 000 generations, though a pure birth (Yule) process (Fig. 2b).
314 Principal tree 2 (PT2) was generated from PT1 by grafting the hybrid taxon to the alternative
315 parental lineage, at the same tree depth (Fig. 2c).

316 Our approach consisted in analysing replicates (Supp. Fig. 3) that contain data simulated
317 under a range of conditions. Each replicate modelled four unlinked genomic blocks, each with three
318 non-homologous genes. For the first two blocks, the genes evolved according to PT2, according to
319 PT1 for the third block and in the last block two genes evolved according to PT1 and one gene was
320 paralogy-affected. Paralogous trees were obtained from PT1, by duplication of the clade containing
321 species 2, 3, 4, 6, 7 and 8, followed by the random pruning of one copy from each taxon descended
322 from the gene duplication (i.e., the locus was returned to single copy in all descendants, Fig. 2d).

323 We simulated gene trees in each block using Kingman's neutral coalescent process
324 implemented in the Dendropy python library (Sukumaran and Holder 2010). On paralogous trees,
325 the random pruning of one copy of each terminal taxon produced trees with RF distance from PT1
326 ranging from a minimum of 0 (identical) to a maximum of 8. We aimed at evaluating the influence

327 of the topological effect introduced by paralogy on our method. To this purpose, we simulated 100
 328 trees for each of the possible distances and grouped them into five paralogy bins corresponding to
 329 the RF values of [0, 2, 4, 6, 8], in order to observe the general trends across the results.



330

331 Figure 2

332

333 Tree depth was also included as a source of variation among replicates. The simulation of
 334 gene trees was performed under a fixed PT depth in number of generations, but varied in terms of
 335 coalescent units (CU; i.e., species tree branch length/ N_e) by modifying the effective population size
 336 (N_e) in the simulations. To this effect, we selected a branch in the trees and varied the global N_e in
 337 order to obtain the desired length for this reference branch. This branch extended from the root to
 338 the node where the hybrid lineage (taxon 4) attached to the parental lineage in each PT (equivalent
 339 in length in PT1 and PT2). The length of this reference branch in CUs is the key parameter that
 340 determines the difficulty of detecting hybridisation and distinguishing it from ILS. ILS will, on
 341 average, be more prevalent when the branch length in CUs is short. Deep coalescence of alleles
 342 from taxon 4 in either PT (i.e., where polymorphisms were held throughout the reference branch)
 343 can result in gene trees where this taxon is placed sister to the remaining species, and thus ILS alone
 344 can produce identical topologies from gene trees derived from either PT. Given that the entire tree
 345 was varied in a systematic way (by changing N_e during the simulations), the reference branch length
 346 also determines the relative difficulty of detecting paralogy (along with the paralogy bin, i.e.,
 347 detection should be easier with greater RF distances). Five different levels of CU, ranging from 2 to
 348 6, with a unit increment, were tested. For a species with one generation/year and N_e of 50 000, this
 349 represents a realistic interval between speciation events, i.e., 100 000 to 300 000 years.

350 In order to include realistic parameters values in our simulation, we acquired nucleotide
351 substitution models and rates from 50 empirical nucleotide alignments of low copy nuclear markers
352 (Sousa et al. 2014), each with the same six species (four *Medicago* species + two species of its
353 sister genus, *Melilotus*). Substitution models were inferred in jModelTest v. 2.1.4 (Darriba et al.
354 2012). The empirical distributions of mean rates were modelled with a lognormal distribution, from
355 which the mu and sigma parameters were retrieved, using the distribution fitting software EasyFit
356 5.3 (MathWave Technology). For each simulated tree, we randomly sampled one of the 50
357 empirical genes to obtain the nucleotide substitution model and the ‘meanrate’ lognormal
358 distribution parameters. Trees were converted from generations to substitutions by sampling a rate
359 for each branch from the associated lognormal distribution. Random selection of substitution
360 models and rates simulated the complexity expected from empirical data. Their effect was not
361 evaluated systematically and we treated them as nuisance parameters downstream, focusing instead
362 on the effect of ILS and duplication. We simulated sequence alignments of 1 500 bp using SeqGen
363 (Rambaut and Grass 1997) for a total of 30 000 alignments: 100 replicates \times 12 genes \times 5 CUs \times 5
364 paralogy bins. For each replicate of 12 genes, we therefore simulated alignments comparable to our
365 empirical data with respect to sequence length per locus and number of loci, as well as model of
366 sequence evolution and tree lengths.

367

368 *Modelling of the Tree Root Prior*

369 In our simulation we used a tree depth of 2 M generations and assumed a generation time of
370 one year. The synthetic alignments were analysed in BEAST, which requires either time or rate
371 prior information. In an application of our test to empirical data, a fossil (for example) could be
372 used to calibrate the root of the species tree. However in our simulation, because we calibrated gene
373 trees instead of species trees, there is a mismatch between species tree and gene tree node ages that
374 needs to be taken into account. We therefore used an exponential prior distribution on the crown
375 node to calibrate the gene trees estimated from the simulated data sets, offset at 2 Ma (as though we
376 had a fossil date for the node) and with a mean equal to $2N_e$ (to include the variation expected due to
377 the coalescent).

378 However, we could not infer N_e directly from our data, since we had simulated a single
379 allele per taxon. To determine the error associated with N_e estimates from allelic data, we simulated
380 100 alleles per taxon from PT1 with SeqGen at population sizes of 10^3 and 10^4 . Using 100
381 simulated alleles for one species and inferring N_e in DNAsp (Librado and Rozas 2009), we
382 observed up to 50% of deviation in the recovered N_e from the values used for allele simulation.
383 Thus, the mean of the exponential distribution for the tree root prior was obtained by first

384 calculating the mean N_e used for the gene tree simulations (across all CUs), and adding 50% to
385 incorporate the expected degree of uncertainty. This value was then multiplied by two ($2N_e$ being
386 the expected mean of exponentially distributed coalescent times) and added to the constraint, which
387 equated to a mean root age of 2.87 Ma.

388

389 *Gene Tree Inference from Simulated Alignments*

390 Each simulated alignment was analysed in BEAST v1.8.0. under a substitution model
391 obtained from jModelTest, a Yule prior on the branching process (a species branching process,
392 because we simulated one sequence per species), and an uncorrelated lognormal relaxed clock with
393 a ucl.d.mean rate prior of 3.6×10^{-9} (Sousa et al. 2014). We used a starting tree with a topology
394 equivalent to PT1 (to speed up the search, given the large number of replicates to analyse) and a
395 root age of 2 Ma. Each dataset was analysed with an MCMC of 5 M generations and the maximum
396 clade credibility tree was summarized in TreeAnnotator (Drummond and Rambaut 2007) with 10%
397 burn-in samples discarded. The resulting gene trees were subjected to paralogy and hybridisation
398 tests, and to the subsequent grouping of genes according to their underlying PT. All tests were
399 implemented in a Python 2.7 pipeline relying on methods from the Dendropy (Sukumaran and
400 Holder 2010), Networkx (Hagberg et al. 2008) and Biopython (Cock et al. 2009) libraries and
401 parallelized using MPI for Python (Dalcín et al. 2008).

402

403 *Grouping of Genomic Blocks and Species Tree Inference*

404 Genomic blocks differing by larger values than what can be obtained from ILS alone are
405 considered to originate from different PTs. In each replicate, the expected result of the hybridisation
406 test is the identification of two PTs (PT1 and PT2, as per our simulations). We evaluated whether
407 the genes categorized into two principal tree groups would recover PT1 and PT2. To this effect we
408 estimated species trees from each group of genes, in ten replicates drawn from the bin RF=8 at
409 CU=2 (the hardest tree length for the test), using *BEAST. A lognormal distribution with
410 mean= 3.61×10^{-9} was set for substitution rates prior (Sousa et al. 2014). For population size we
411 chose a normal distribution for prior with a mean= 4.34×10^5 (the mean population size value used
412 for simulation with a 50% error increase, as used for data analysis) and sd= 5×10^4 , truncated to
413 2×10^6 . A normal prior was applied to the species tree root height (mean= 2.87×10^6 , sd= 0.5×10^6).
414 Analyses were run for 50M MCMC generations with three separate chains for each analysis.

415

416 APPLYING THE METHOD TO THE *MEDICAGO* DATASET

417 Paralogy and hybridisation tests were run on the 12 inferred gene trees, in the same way as

418 for the simulated data, with acceptance at p% critical values of 95%, 85%, 75%, 65% of the null
419 distributions. To generate the ILS null distribution we used the effective population size of $N_e = 240$
420 000 (Maureira-Butler et al. 2008). After running the tests and removing putative paralogy-affected
421 genes, we grouped the genes that shared the same principal tree. Genes within each group were
422 subjected to species tree analysis in *BEAST. Chains were run for 50 M generations using the
423 inferred substitution model for each gene, an uncorrelated lognormal relaxed clock with a
424 ucl.d.mean rate normal prior with mean= 3.6×10^{-9} , a Yule model species tree prior and a normal
425 prior on the species-tree root height with mean= 15.9×10^6 and sd= 2.7×10^6 .

426 In order to determine how important the new part of the testing procedure was – the
427 paralogy test – we also ran the hybridisation test by itself on the *Medicago* data. Thus, inside a
428 block, genes were presumed to differ by ILS alone, and we did not test for within-block
429 incongruence. Each gene in a block was tested against all the genes in the other blocks. Blocks were
430 handled separately if two pairwise rejections between blocks occurred and the resulting groups of
431 blocks were analysed separately in *BEAST, as before. We also compared these analyses with two
432 other approaches that do not use any paralogy or hybridisation test. Firstly, we ran a 12 gene
433 *BEAST analysis (i.e., with gene trees assumed to be unlinked and thus independent due to
434 recombination among loci). Secondly, we ran a concatenation of the 12 genes in BEAST. In both
435 approaches each gene was endowed with its own substitution and clock model.

436

437 RESULTS

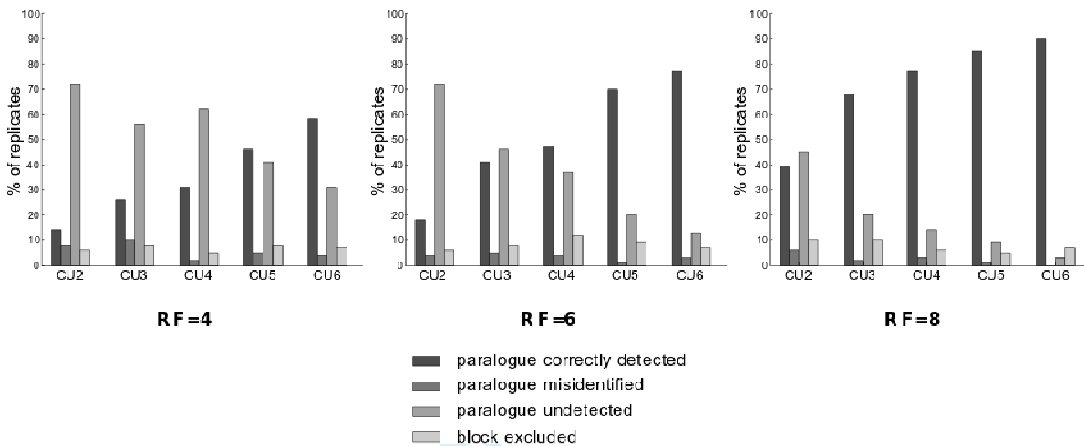
438 SIMULATED DATA

439 *Paralogy Test*

440 We explored the effect of two variables on the ability to detect paralogy with our method,
441 namely: 1) the amount of topological variation introduced by a gene duplication with subsequent
442 random losses of one paralogous locus, such that all taxa returned to single copy status, and 2) the
443 length of a reference branch measured in CU units.

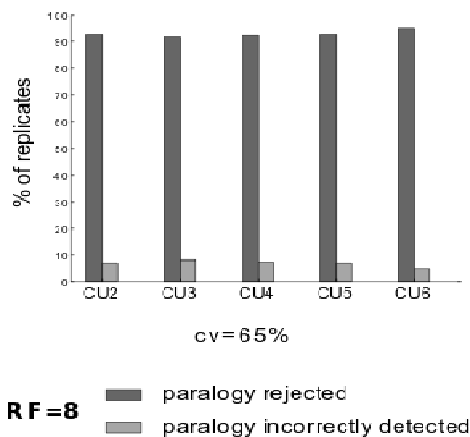
444 The power of the method increased with the depth of the tree (Fig. 3). For settings with large
445 paralogy induced topological variation (RF=8) we found that lowering the detection critical value
446 was effective to identify the paralogy-affected gene (Supp. Fig. 4). At lower cv, type II error (failure
447 to detect real paralogy) decreased and the rate of successful detection of genuine paralogy increased
448 (Supp. Fig. 4). This occurred with only a minor increase in type I error (either due to paralogous
449 misidentification or to the incorrect exclusion of the whole genomic block). For the shortest tree
450 (CU=2) at 95% cv, type I error ranged from zero to ~10% for block exclusion and ~5% for
451 misidentification (Supp. Fig. 4). At 65% cv and for the shortest tree (CU=2), successful detection of

452 paralogy-affected genes in the RF=8 bin dropped to ~ 40% (Fig. 3), highlighting that shallow tree
 453 depth is the most challenging situation for coalescent-based methods (see also Maddison and
 454 Knowles 2006). Without paralogy-affected genes present in a block, the test performed well at all
 455 cv and CUs, showing less than 10% type I error in the most difficult case (CU=2) at 65% cv (Fig.
 456 4).



457 Figure 3.

458



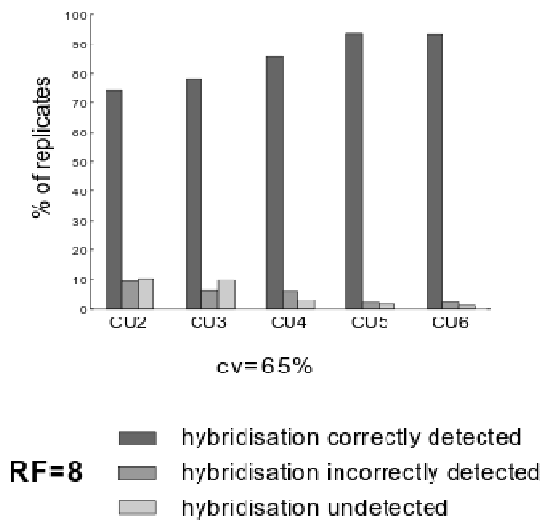
459

460 Figure 4.

461

462 *Hybridisation Test*

463 We focus on the results for the paralogy setting that introduces the most severe discordance
 464 (RF=8). The hybridisation test, like the paralogy test, performed better as the cv was lowered (Supp.
 465 Fig. 5). The rate of correct hybridisation detection at 65% cv approached 95% in the most
 466 favourable cases, CU=6, (Fig. 5) and was still of ~ 75% for the most difficult case, CU =2 (Fig. 5),
 467 with the failure rate to detect hybridisation not exceeding 10%. The false detection of hybridisation
 468 between blocks that share the same PT, occurred at less than 10% at any cv.



469 Figure 5.

470

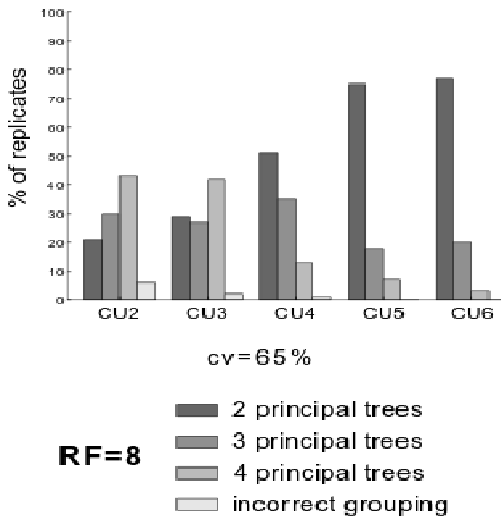
471 *Grouping of Genomic Blocks*

472 Blocks containing genes that differed due to ILS alone (in the hybridisation test) were
 473 assumed to be derived from the same origin. Genes in these blocks were then pooled for species tree
 474 inference. In our simulation settings, with four genomic blocks tracking two PTs, the hybridisation
 475 test can potentially output one to four PTs. Failure to recover the two initial grouping occurs when
 476 blocks are incorrectly clustered or wrongly segregated. Lower critical values produced more
 477 accurate results (Supp. Fig. 6). At 65% cv and CU=6, genomic blocks were correctly pooled into
 478 two groups corresponding to the two initial PTs, in over 75% of cases (Fig. 6). At this cv, the rate of
 479 wrong segregation increased at lower CUs (Fig. 6). However, the incorrect grouping of blocks that
 480 were simulated from different PTs – which represents the most severe kind of error – was only ~
 481 5% in the most challenging case at the lowest CU (Fig. 6).

482 In the more difficult cases, the method sometimes failed to group blocks originating from
 483 the same principal tree (i.e., it predicted either three or four PTs when only two PTs were present):
 484 this over-segregation rose from ~25% at CU=6 to ~75% at CU=2 across the 65% critical value
 485 cases in bin RF=8 (Fig. 6). Failure to group blocks from the same PT is not of major concern,
 486 however, as gene blocks that have been separated by our method can still cluster a posteriori when
 487 they are shown share the same tree in separate species tree analyses.

488 Species trees inferred from each group of blocks in replicates at CU=2 (RF=8 and cv=65%)
 489 recovered the hybrid taxon in the correct parental lineage in 100% of cases, with full branch support
 490 in 79% of cases and at the correct position in the parental lineage in 64% of the trees (not shown).
 491 Three out of 32 analyses did not converge for any of the three chains and were dismissed. The exact
 492 topology of the input PT for the gene simulation was recovered in 45% of cases. Both PTs were

493 recovered in eight out of 10 replicates, with full support in seven replicates. Failure to recover both
 494 PTs in two replicates resulted from lack of convergence of the MCMC.



495 Figure 6.

496

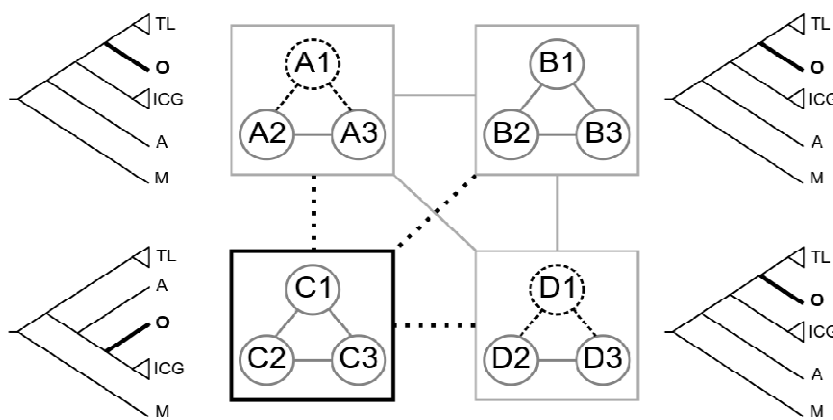
497 *MEDICAGO ORBICULARIS* DATASET

498 BEAST trees for the individual markers (not shown) revealed alternative supported
 499 placements (posterior probability –PP– >95%) of *M. orbicularis*: either 1) within the (*M. truncatula*
 500 - *M. littoralis* clade) (genes A1, A3, B1, B2 and D3) or 2) within the (*M. intertexta* - *M. ciliaris* - *M.*
 501 *granadensis*) clade (genes C1, C2). A2 showed no supported placement for *M. orbicularis*. For B3,
 502 the topology was the same as for the other genes of block B, but with marginal support for the
 503 branching of *M. orbicularis*. In C3, *M. orbicularis* was placed with low support as sister to the *M.*
 504 *truncatula* clade; this relationship was not recovered in C1 or C2. In D1, *M. orbicularis* appeared as
 505 sister to the remaining *Medicago* species and in D2 it was placed as sister to the *M. intertexta* clade,
 506 albeit with no support in either case. The (*M. truncatula* - *M. littoralis*) clade was always retrieved,
 507 whereas *M. granadensis* was separated from *M. intertexta* and *M. ciliaris* in gene A2. *Medicago*
 508 *arabica* was recovered, with good support, as sister to *M. orbicularis* (A1), to (*M. orbicularis* + *M.*
 509 *truncatula*) (B2), to *M. granadensis* (B1, C3, D1) and to the remaining *Medicago* species (B3, C2,
 510 D2). Thus, we observed a clear lack of concordance between several well supported topologies.

511 We applied our paralogy and hybridisation detection method to the *Medicago* data at critical
 512 values of 95%, 85% and 75%. We obtained the exclusion of two genes (A1, D1) as probably
 513 paralogy-affected, and the recovery of two groups of blocks (blocks A, B and D versus block C)
 514 tracking two different principal trees. With a cv of 65%, the paralogy test excluded block B. When
 515 species trees were inferred from each block, the topology alternated between one shared by block
 516 A, B and D (*M. orbicularis* + *truncatula* clade) and a different one recovered with C (*M. orbicularis*

517 + *ciliaris* clade) (Fig. 7; Supp. Fig. 7). This result was concordant with the outcome of the
 518 hybridisation test.

519 The grouping of blocks A, B and D was recovered in the hybridisation test for all cvs above
 520 75% and was in agreement with the species tree topologies obtained from each block after the
 521 exclusion of paralogy-affected genes (Fig 7). Since our simulation demonstrated a higher type I
 522 error at 65% cv, we opted for keeping B in the (A, D) group. We also estimated the depth of this
 523 hybridisation problem, based on individual gene trees, to be up to eight CUs (not shown), a tree
 524 depth for which the paralogy-test performed well in the simulated data at cvs above 75%.



525 Figure 7.

526

527 The species tree inferred from group (A, B, D) (Supp. Fig. 7) displayed the same topology
 528 as the one obtained from its individual blocks. Whereas the position of *M. orbicularis* was poorly
 529 supported in the individual block analyses, its support increased in the group analysis, reaching a PP
 530 of 0.94. For block C the support for the alternative position of *M. orbicularis* was high (PP=0.93).
 531 *M. orbicularis* branched as sister to the *M. truncatula* clade in five genes across three blocks, and as
 532 sister to the *M. intertexta* clade in two genes from one block. These alternative relationships were
 533 also obtained in the species tree (after exclusion of paralogy-affected genes) when we analysed the
 534 individual blocks and the grouped blocks. The position of *M. arabica* also differed between the two
 535 recovered species trees: on the tree inferred from group (A, B, D), *M. arabica* appeared as sister to
 536 all other *Medicago* species (PP=0.71), and on the tree inferred from block C, it branched as sister to
 537 the (*M. orbicularis* + *M. intertexta*) clade (PP=0.93). On the gene trees, this species was involved in
 538 four different well supported topologies and no block contained at least two genes that tracked the

539 same history for this taxon. When *M. arabica* was removed from the analyses, the position of *M.*
540 *orbicularis* was unchanged in the two trees obtained from blocks (A, B, D) and C (not shown).

541 When the hybridisation test was performed without filtering for paralogy-affected genes, the
542 block groupings and species trees inference steps were adversely affected. Block A was separated
543 from B and D, which grouped together as before (Supp. Fig. 7). Block A alone returned a different
544 topology than that obtained from (A, B, D), namely that *M. orbicularis* was weakly supported as
545 sister to *M. arabica*, and together these were placed as sister to the *truncatula* clade with marginal
546 support (Supp. Fig. 7); B and D returned the same position for *M. orbicularis* as before, but with
547 lower support (PP=0.84) than in the (A, B, D) species tree with paralogy-affected genes removed.
548 (Supp. Fig. 7). The result from C was unchanged.

549 Gene A1 appeared to drive the placement of *M. orbicularis* as sister to *M. arabica* in the
550 block A, as it is the only gene in the block displaying this relationship, while the co-analysis of A2
551 and A3 inferred *M. orbicularis* as sister to the *truncatula* clade, but with lower support (0.71; Supp.
552 Fig. 7). In fact, each of blocks A, B and D alone and together (with suspected paralogy-affected
553 genes removed in each case) recovered this relationship, although with weak support from
554 individual blocks (Supp. Fig. 7). If *M. orbicularis* sister grouping to the *truncatula* clade is correct
555 (as it appears to be), then including suspected paralogues degraded support, negatively impacted the
556 block grouping and returned a block with a different topology. What was a well-supported
557 relationship became poorly supported or weakly contradicted to the extent that no clear
558 interpretation could be obtained from this group of genes. With respect to hybrid origins, one
559 parentage was recovered as before (block C), but the other parental lineage was not discovered.
560 Thus, hybridisation could not be inferred from the *BEAST analyses *unless* the paralogy test had
561 been performed and suspected paralogy-affected genes removed.

562 The analysis of 12 separated genes in *BEAST recovered the same topology as from blocks
563 (B, D) (Supp. Fig. 7). However, the supports for including *M. orbicularis* in the *M. truncatula* clade
564 and *M. arabica* in the *M. ciliaris* clade, were lower (PPs varied from 0.84 to 0.75, and 0.55 to 0.52,
565 respectively). The concatenation of all 12 genes in BEAST returned a topology identical to blocks
566 (A, B, D), but with 1.0 PP on each node, and again *Medicago orbicularis* was sister to the *M.*
567 *truncatula* clade (Supp. Fig. 7). Both methods suggested a single possible relationship compatible
568 with one of the parental origins, but failed to provide evidence for a hybridization event.

569

570 **DISCUSSION**

571 RECOVERY OF DIFFERENT PRINCIPAL TREES IN *MEDICAGO* AFTER PARALOGY AND HYBRIDISATION

572 TESTING

573 Applying our method to the *Medicago* dataset resulted in the recovery of two principal trees.
574 Our results corroborate the initial hypothesis that hybridisation occurred in the *Medicago*
575 *orbicularis* lineage, as it is recovered in two clear alternative positions in the two principal trees.
576 The paralogy detection step indicated the presence of two putatively paralogy-affected genes: In
577 gene D1, *Medicago orbicularis* branched deeper than in any other gene, which is indeed consistent
578 with a hypothesis of paralogy. The other putative paralogy-affected gene, A1, recovered *M. arabica*
579 sister to *M. orbicularis* with good support at a shallow position. This would be more consistent with
580 recent hybridisation between these two species than with paralogy, but the gene was nevertheless
581 effectively identified within the respective genomic block as a highly-discordant gene tree.
582 Removal of suspected paralogy-affected genes was critical to correctly infer the parentage of the
583 hybrid lineage and recover the two principal trees. The hybridisation detection step indicated that
584 block C evolved under a different principal tree from the remaining blocks and species tree
585 reconstruction confirms the results of the hybridisation test: the tree recovered from block C
586 recovered *M. orbicularis* as sister to the *M. ciliaris* clade, whereas the remaining blocks
587 (individually or together) always placed *M. orbicularis* as sister to the *M. truncatula* clade. This
588 result shows that the use of genomic location can reveal hybridisation signal that would otherwise
589 be overlooked. None of the *Medicago orbicularis* alleles were identical to any of the alleles in the
590 two putative parental lineages, nor appeared at shallow positions within these lineages, which
591 suggests that the two principal trees obtained through our tests are compatible with a hypothesis of
592 ancient (rather than recent) hybridisation involving the *M. orbicularis* lineage. Six previously
593 published dated gene trees suggest that the hybridisation episode(s) probably took place more than
594 1.7 Ma years ago (Sousa et al. 2016).

595

596 *Taxon and Gene Sampling Effects on the Method*

597 One limitation of the metric used in the current method is that the Robinson-Foulds distance
598 could return only a minimum value (RF=2) for an unstable taxon placed at alternative positions on
599 two sister lineages, a situation that could be due to hybridisation. Such a small value is attained
600 when no other taxon branches between the unstable taxon and the divergence of the two sister
601 lineages. Our method would likely mistake this small distance for the effect of ILS in almost any
602 realistic example. To overcome this limitation, after we identified an unstable taxon as a putative
603 hybrid, it was necessary to sample a species sister to the putative hybrid, or one branching earlier
604 from the parental branch. The two putative parental lineages of *M. orbicularis* are represented by
605 the *M. intertexta* clade and the *M. truncatula* clade in our sample.

606 We included *M. arabica* in our sampling as this species appeared to be related to *M.*

607 *orbicularis* in previous phylogenies (e.g., Sousa et al. 2016) and was expected to increase the RF
608 distance between alternative trees with respect to the *M. orbicularis* placement. However, the
609 placement of *M. arabica* in the two recovered principal trees is also unstable, which suggests that
610 *M. arabica* may also be involved in a hybridisation process. On the tree recovered from blocks (A,
611 B, D), the low support for the placement of *M. arabica* (PP=0.71) may be caused by conflicting
612 hybridisation signal in *M. arabica* among the three grouped blocks. This alternative placement of
613 *M. arabica* did not affect the grouping of blocks as the corresponding topological difference, given
614 the current sampling, is minimal (RF=2). Testing for hybridisation in *M. arabica* would, therefore,
615 require a different sampling of taxa. This illustrates another limitation of our method, namely that
616 only one independent hybridisation history can be effectively detected at a time. This limitation is
617 due to the fact that genomic blocks can trace different hybridisation histories for different species,
618 i.e., the genomic blocks tracing the same history in one lineage may trace more than one history in
619 another lineage. If hybridisation occurs in more than one sampled lineage, our method may actually
620 fail to group any genomic blocks. In such cases, a much higher number of genomic blocks need to
621 be sampled in order to increase the chances of obtaining all principal trees that describe all the
622 hybridisation events. Alternatively, lineages of suspected hybrid origin can be included one by one
623 in analyses involving other lineages that do not appear to have hybrid histories (Sousa et al. 2016).
624 Individual signals of hybridisation should be much more easily recovered and interpreted than
625 overlapping signals.

626

627 *Modelling Nuisance Parameters and Coalescent Stochasticity*

628 We added complexity to our simulation with several sources of variation, such as
629 substitution model variation and model selection uncertainty, substitution rate variation, gene
630 branch rate variation, gene tree inference uncertainty, coalescent stochasticity and random gene
631 copy loss after gene duplication. Thus, the simulations covered a realistic and challenging range of
632 parameters for our method. The effect of coalescent stochasticity on tree-distance distributions,
633 which may contain a tail of extreme values driven by very improbable topologies, needs to be
634 overcome by choosing an adequate cv for the one-tailed tests. At high cv, the tests lose power, even
635 at high CUs. This implies that previous applications of the test (for hybrid detection only) were
636 probably conservative, i.e., favoured the ILS null (Maureira-Butler et al. 2008; Blanco-Pastor et al.
637 2012; Ramadugu et al. 2013). We found that decreasing the cv to 65% resulted in a decrease of type
638 II error without a proportional increase of type I error.

639

640 *Relevance of Paralogy Detection*

641 If the proportion of paralogy-affected genes in a sample is low relative to non-affected
642 genes, undetected paralogy may only marginally impact species-tree inference, as enough
643 orthologous genes should swamp the paralogous signal. However, if this proportion is high, which
644 may be the case when many taxa but few genes are sampled, or when WGD has recently occurred
645 in an ancestor, identifying paralogues becomes more important. In such cases our method is able to
646 detect paralogy-affected genes that would be the most detrimental for species tree inference.

647 Importantly, paralogy caused by WGD could likely also be detected by our method. If the
648 reference organism used to determine genomic location has genes X, Y, Z in a genomic block,
649 WGD in a related lineage will produce two blocks with X, Y, Z and X', Y', Z'. With random gene
650 loss in each descendent, the gene trees may differ significantly. For example, a gene tree containing
651 both Z/Z' copies among taxa (but only one copy sampled per taxon) can carry a markedly different
652 signal than the tree produced from pure X and Y sequences (e.g., if the X' and Y' copies were lost
653 before taxon divergence). If the sequences are recovered using primers or gene capture, the fact that
654 some of these gene loci do not reside in the same physical block will be undetected, and would be
655 handled as linked loci (in the absence of further information, i.e., for the majority of samples that
656 would lack a physical map). The gene tree containing a mixture of Z/Z' sequences could be
657 identified by our paralogy-test as being inconsistent with the other genes of the block, although the
658 paralogous copies were not derived from tandem duplication but from WGD instead.

659 If a physical map for these genes were available, then the paralogy of Z' would be already
660 known. So our test is applicable when the physical position of genes is well estimated in the
661 ancestor of the whole lineage, but has been subsequently modified by WGD in some sub-lineages.
662 This may be a fairly common situation, where a physical map is known for one model organism that
663 is used to represent an entire clade that is not excessively old, which would imply small departure
664 from the ancestral genomic map. Given how common WGD is in some lineages (e.g., flowering
665 plants, Soltis et al. 2015 and references therein), our test may be especially useful for detecting
666 WGD-derived paralogy.

667

668 *Future Prospects*

669 The use of genomic blocks to investigate hybridisation can be improved upon, and larger
670 blocks containing a higher number of genes should be considered, albeit with caution. The genome-
671 wide SNP study of several *Medicago* species by Yoder et al. (2013) is instructive with regard to this
672 point. That study found that different genomic regions contained some strongly supported
673 alternative trees compared to the entire genome-wide data set of over 82 000 SNPs. This is
674 congruent with the idea that different genomic blocks contain different histories due to ancient

675 hybridisation. The genomic blocks used by Yoder et al. (2013) were fairly large (containing 500
676 SNPs spread over an average of 1.7 Mb of contiguous genomic sequence) and can clearly include
677 many coalescent histories, given that linkage disequilibrium decays after only a few thousand
678 nucleotides in *M. truncatula* (Branca et al. 2011). However, our results here (and with further taxon
679 sampling, not shown) suggest that genomic blocks as small as 60 kb can contain more than one
680 historical signal due to hybridisation. It is therefore possible that most of the large genomic regions
681 used by Yoder et al. (2013) contain several hybridisation histories. This might explain why few
682 large genomic blocks (defined as four consecutive windows, covering around 4.2 Mb) were found
683 that displayed a single clear conflicting signal with the genome-wide tree. We contend that any
684 genomic block as large or larger than 1.7 Mb (the smallest genomic block size use by Yoder et al.
685 2013) is highly likely to contain more than one parental origin. If so, then the sparse sampling of
686 SNPs across these genomic blocks may not be the best strategy to unravel different histories due to
687 ancient hybridisation. A much denser sampling of characters over the scale of 10s to 100s of kbp
688 might be a better approach to infer resolved gene trees of much smaller genomic partitions that
689 could be compared with other genomic partitions. This suggests that our test method can be
690 improved by including, e.g., two, three or four times as many genes per genomic block, but
691 extending the genomic blocks to maintain a similar density of sampled genes per length of DNA.
692 The paralogy detection step in our methods could then serve to identify physical boundaries within
693 a genomic block that separates two sets of genes with alternative hybrid histories.

694 Other future developments for our method should be centred on the use of alternative
695 metrics that take into account other parameters besides tree topology, such as branch length or
696 support values for each branch. A more sensitive metric would help overcome possible limitations
697 caused by taxon sampling, which results in lack of power in gene tree pairwise comparisons. Other
698 gene sampling approaches, which do not require such comprehensive a priori genomic information
699 as the one available to us (a completely annotated genome), should also be pursued. Given the
700 current developments in enrichment and sequencing techniques, it should not be difficult to
701 implement sampling strategies that extend on ours, to include the sequencing of long contiguous
702 regions that would enable the use of genomic location for a better understanding of complex
703 phylogenetic problems.

704

705 CONCLUSIONS

706 Using linked and unlinked loci we found evidence of ancient hybridisation in *Medicago*
707 *orbicularis*, confirming that the evolutionary history of this genus is affected by reticulation
708 (Maureira-Butler et al. 2008). It is likely that introgressive hybridisation has occurred in many

709 lineages within the genus (Sousa et al. 2016). Thus, any attempt to represent the phylogeny of
710 *Medicago* as a single dichotomous tree requires caution. The method used here to investigate
711 incongruence in *Medicago* introduces a new angle to species tree inference in the presence of
712 paralogy or other sources of highly discordant trees in addition to ILS and hybridisation. So far,
713 phylogenetic inference methods have dealt with these two processes separately. Here, we tackle
714 both processes sequentially and introduce the use of genomic location as an indicator of reticulate
715 gene inheritance. Our approach does not rely on information about the (as yet) unknown species
716 tree in order to test for paralogy and hybridisation. This is an advantage compared to other methods,
717 because we do not need to infer, a priori, a species tree that might itself be misled by paralogy or
718 hybridisation in order to assess the presence of misleading factors (e.g., Joly et al. 2009). We also
719 do not need to specify in advance which lineages might be affected by hybridisation (unlike e.g.,
720 Kubatko et al. 2009; Jones et al. 2013). Finally, rather than attempting to reconcile paralogy- or
721 hybridisation-affected tree topologies with the coalescent process (Rasmussen and Kellis 2012; Yu
722 et al. 2013), we sort and deal with each process separately, as paralogy-affected genes are simply
723 excluded and the remaining loci are analysed according to a shared history, allowing for a more
724 complete perspective of speciation and diversification processes.

725

726 ACKNOWLEDGEMENTS

727 This work was supported by grants from the Swedish Research Council (grant 2009-5206),
728 the Royal Swedish Academy of Sciences, Lars Hiertas Minne fund, The Royal Physiographic
729 Society in Lund, Helge Ax:son Johnsons fund, and the Lundgrenska fund to B.E.P; from the P. A.
730 Larssons fund and Lars Hiertas Minne fund to F.S.

731

732 REFERENCES

- 733 Ané C., Larget B., Baum D.A., Smith S.D., Rokas A. 2007. Bayesian estimation of concordance
734 among gene trees. *Mol. Biol. Evol.* 24:412-426.
- 735 Blanc G., Wolfe K.H. 2004. Functional divergence of duplicated genes formed by polyploidy
736 during *Arabidopsis* evolution. *Pl. Cell* 16:1679-1691.
- 737 Blanco-Pastor J.L., Vargas P., Pfeil B.E. 2012. Coalescent simulations reveal hybridization,
738 incomplete lineage sorting in Mediterranean *Linaria*. *PloS One* 7:e39089.
- 739 Bena G. 2001. Molecular phylogeny supports the morphologically based taxonomic transfer of the
740 “medicagoid” *Trigonella* species to the genus *Medicago* L. *Pl. Syst. Evol.* 229:217-236.

- 741 Bertrand Y.J.K., Scheen A.-C., Marcussen T., Pfeil B.E., Sousa F., Oxelman B. 2015. Assignment
742 of homoeologues to parental genomes in allopolyploids for species tree inference, with an example
743 from *Fumaria* (Papaveraceae). *Syst. Biol.* 64: 448-471.
- 744 Bosse M., Megens H.-J., Madsen O., Frantz L.A.F., Paudel Y., Crooijmans R.P.M.A., Groenen
745 M.A.M. 2014. Untangling the hybrid nature of modern pig genomes: a mosaic derived from
746 biogeographically distinct, highly divergent *Sus scrofa* populations. *Mol. Ecol.* 23:4089-4102.
- 747 Branca A., Paape T.D., Zhou P., Briskine R., Farmer A.D., Mudge J., Bharti A.K., Woodward J.E.,
748 May G.D., Gentzittel L., Ben C., Denny R., Sadowsky M.J., Ronfort J., Bataillon T., Young N.D.,
749 Tiffin P. 2011. Whole-genome nucleotide diversity, recombination, linkage disequilibrium in the
750 model legume *Medicago truncatula*. *Proc. Natl. Acad. Sci. U.S.A.* 108:E864-E870.
- 751 Buckley T.R., Cordeiro M., Marshall D.C., Simon C. 2006. Differentiating between hypotheses of
752 lineage sorting, introgression in New Zealand alpine cicadas (*Maoricicada* Dugdale). *Syst. Biol.*
753 55:411-425.
- 754 Buerkle C.A., Rieseberg L.H. 2008. The rate of genome stabilization in homoploid hybrid species.
755 *Evolution* 62:266-275.
- 756 Carstens B.C., Knowles L.L. 2007. Estimating species phylogeny from gene-tree probabilities
757 despite incomplete lineage sorting: an example from *Melanoplus* grasshoppers. *Syst. Biol.* 56:400-
758 411.
- 759 Chen J., Huang Q., Gao D., Wang J., Lang Y., Liu T., Li B., Bai Z., Luis Goicoechea J., Liang C.,
760 Chen C., Zhang W., Sun S., Liao Y., Zhang X., Yang L., Song C., Wang M., Shi J., Liu G., Liu J.,
761 Zhou H., Zhou W., Yu Q., An N., Chen Y., Cai Q., Wang B., Liu B., Min J., Huang Y., Wu H., Li
762 Z., Zhang Y., Yin Y., Song W., Jiang J., Jackson S.A., Wing R.A., Wang J., Chen M. 2013. Whole-
763 genome sequencing of *Oryza brachyantha* reveals mechanisms underlying *Oryza* genome
764 evolution. *Nat. Commun.* 4:1595 doi: 10.1038/ncomms2596.
- 765 Cock P.J, Antao T., Chang J.T., Chapman B.A., Cox C.J., Dalke A., Friedberg I., Hamelryck T.,
766 Kauff F., Wilczynski B., de Hoon M.J. 2009. Biopython: freely available Python tools for
767 computational molecular biology, bioinformatics. *Bioinformatics* 25:1422-1423.
- 768 Cotton J.A., Page R.D. 2005. Rates, patterns of gene duplication, loss in the human genome. *Proc.*
769 *Roy. Soc. London, Ser. B, Biol. Sci.* 272:277-283.

- 770 Cui L., Wall P.K., Leebens-Mack J.H., Lindsay B.G., Soltis D.E., Doyle J.J., Soltis P.S., Carlson
771 J.E., Arumuganathan K., Barakat A., Albert V.A., Ma H., dePamphilis C.W. 2006. Widespread
772 genome duplications throughout the history of flowering plants. *Genome Res.* 16:738-749.
- 773 Dalcín L., Paz R., D'Elia M.S.J. 2008. MPI for Python: Performance improvements, MPI-2
774 extensions. *J. Parallel Distrib. Comput.* 68:655-662.
- 775 Darriba D., Taboada G.L., Doallo R., Posada D. 2012. jModelTest 2: more models, new heuristics,
776 parallel computing. *Nature methods* 9:772-772.
- 777 Degnan J.H., Rosenberg N.A. 2009. Gene tree discordance, phylogenetic inference, the
778 multispecies coalescent. *Trends Ecol. Evol.* 24:332-340.
- 779 Degnan J.H., Salter L.A. 2005. Gene tree distributions under the coalescent process. *Evolution*
780 59:24-37.
- 781 Doyle J.J. 1992. Gene trees and species trees: molecular systematics as one-character taxonomy.
782 *Syst. Bot.* 17:144-163.
- 783 Drummond A.J., Rambaut A. 2007. BEAST: Bayesian evolutionary analysis by sampling trees.
784 *BMC Evol. Biol.* 7:214.
- 785 Edwards S.V., Liu L., Pearl D.K. 2007. High-resolution species trees without concatenation. *Proc.*
786 *Natl. Acad. Sci. U.S.A.* 104:5936-5941.
- 787 Foxe J.P., Slotte T., Stahl E.A., Neuffer B., Hurka H., Wright S.I. 2009. Recent speciation
788 associated with the evolution of selfing in *Capsella*. *Proc. Natl. Acad. Sci. U.S.A.* 106:5241-5245.
- 789 Good J.M., Vanderpool D., Keeble S., Bi K. 2015. Negligible nuclear introgression despite
790 complete mitochondrial capture between two species of chipmunks. *Evolution* 69: 1961-1972.
- 791 Gossman T.I., Song B.-H., Windsor A.J., Mitchell-Olds T., Dixon C.J., Kapralov M.V., Filatov
792 D.A., Eyre-Walker A. 2010. Genome wide analyses reveal little evidence for adaptive evolution in
793 many plant species. *Mol. Biol. Evol.* 27:1822–1832.
- 794 Green R.E., Krause J., Briggs A.W., Maricic T., Stenzel U., Kircher M., Patterson N., Li H., Zhai
795 W., Fritz M.H., Hansen N.F., Durand E.Y., Malaspinas A.S., Jensen J.D., Marques-Bonet T., Alkan
796 C., Prüfer K., Meyer M., Burbano H.A., Good J.M., Schultz R., Aximu-Petri A., Butthof A., Höber
797 B., Höffner B., Siegemund M., Weihmann A., Nusbaum C., Lander E.S., Russ C., Novod N.,

- 798 Affourtit J., Egholm M., Verna C., Rudan P., Brajkovic D., Kucan Z., Gusic I., Doronichev V.B.,
799 Golovanova L.V., Lalueza-Fox C., de la Rasilla M., Fortea J., Rosas A., Schmitz R.W., Johnson
800 P.L., Eichler E.E., Falush D., Birney E., Mullikin J.C., Slatkin M., Nielsen R., Kelso J., Lachmann
801 M., Reich D., Pääbo S. 2010. A draft sequence of the Neandertal genome. *Science* 328:710-722.
- 802 Guschanski K., Krause J., Sawyer S., Valente L.M., Bailey S., Finstermeier K., Sabin R., Gilissen
803 E., Sonet G., Nagy Z.T., Lenglet G., Mayer F., Savolainen V. 2013. Next-generation museomics
804 disentangles one of the largest primate radiations. *Syst. Biol.* 62: 539–554.
- 805 Hagberg A., Swart P., Chult D. 2008. Exploring network structure, dynamics, function using
806 networkx. Technical report. Los Alamos, NM: Los Alamos National Laboratory (LANL).
- 807 Heled J., Drummond A.J. 2010. Bayesian inference of species trees from multilocus data. *Mol.*
808 *Biol. Evol.* 27:570-580.
- 809 Hermansen J.S., Saether S.A., Elgvin T.O., Borge T., Hjelle E., Sætre, G.P. 2011. Hybrid speciation
810 in sparrows I: phenotypic intermediacy, genetic admixture, barriers to gene flow. *Molec. Ecol.*
811 20:3812-3822.
- 812 Holland B.R., Benthin S., Lockhart P.J., Moulton V., Huber K.T. 2008. Using supernetworks to
813 distinguish hybridization from lineage-sorting. *BMC Evol. Biol.* 8:202.
- 814 Huelsenbeck J.P., Larget B., Alfaro M.E. 2004. Bayesian phylogenetic model selection using
815 reversible jump Markov chain Monte Carlo. *Molec. Biol. Evol.* 21:1123-1133.
- 816 Huelsenbeck J.P., Ronquist F. 2005. Bayesian analysis of molecular evolution using MrBayes. In
817 *Statistical methods in molecular evolution*. New York: Springer. p. 183-226
- 818 Innes R.W., Ameline-Torregrosa C., Ashfield T., Cannon E., Cannon S.B., Chacko B., Chen
819 N.W.G., Couloux A., Dalwani A., Denny R., Deshpande S., Doyle J.J., Egan A., Geffroy V.,
820 Glover N., Hans C.S., Howell S., Ilut D., Jackson S., Lai H., Mammadov J., Martin del Campo S.,
821 Metcalf M., Nguyen A., O’Bleness M., Pfeil B.E., Podicheti R., Ratnaparkhe M.B., Roe B.A.,
822 Saghai Maroof M.A., Samain S., Sanders I., Séguens B., Sévignac M., Sherman-Broyles S.,
823 Thureau V., Tucker D.M., Walling J., Wawrzynski A., Yi J., Young N.D. 2008. Differential
824 accumulation of retroelements, diversification of NB-LRR disease resistance genes in duplicated
825 regions following polyploidy in the ancestor of soybean. *Plant Physiol.* 148:1740-1759.
- 826 Joly S., McLenachan P.A., Lockhart P.J. 2009. A statistical approach for distinguishing

- 827 hybridization, incomplete lineage sorting. *Am. Nat.* 174:E54-E70.
- 828 Joly S., Pfeil B.E., Oxelman B., McLenachan P.A., Lockhart P.J. 2010. A statistical approach for
829 distinguishing hybridization, incomplete lineage sorting: correction. *Am. Nat.* 175:621-622.
- 830 Jones G., Sagitov S., Oxelman B. 2013. Statistical inference of allopolyploid species networks in
831 the presence of incomplete lineage sorting. *Syst. Biol.* 62:467-478.
- 832 Kingman J.F.C. 1982. The coalescent. *Stochastic processes, their applications* 13:235-248.
- 833 Kubatko L.S., Degnan J.H. 2007. Inconsistency of phylogenetic estimates from concatenated data
834 under coalescence. *Syst. Biol.* 56:17-24.
- 835 Kubatko, L. S. 2009. Identifying hybridization events in the presence of coalescence via model
836 selection. *Syst. Biol.* 58:478-488.
- 837 Kubatko L.S., Carstens B.C., Knowles L.L. 2009. STEM: species tree estimation using maximum
838 likelihood for gene trees under coalescence. *Bioinformatics* 25:971-973.
- 839 Lavin M., Herendeen P.S., Wojciechowski M.F. 2005. Evolutionary rates analysis of Leguminosae
840 implicates a rapid diversification of lineages during the Tertiary. *Syst. Biol.* 54:575-594.
- 841 Librado P., Rozas J. 2009. DnaSP v5: a software for comprehensive analysis of DNA
842 polymorphism data. *Bioinformatics* 25:1451-1452.
- 843 Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin
844 R., 1000 Genome Project Data Processing Subgroup. 2009. The sequence alignment/map format,
845 SAMtools. *Bioinformatics* 25:2078-2079.
- 846 Liu K., Dai J., Truong K., Song Y., Kohn M.H., Nakhleh L. 2014. An HMM-based comparative
847 genomic framework for detecting introgression in eukaryotes. *PLoS Comp. Biol.* 10:e1003649
- 848 Liu K.J., Steinberg E., Yozzo A., Song Y., Kohn M.H., Nakhleh L. 2015. Interspecific introgressive
849 origin of genomic diversity on the house mouse. *Proc. Natl. Acad. Sci. U.S.A.* 112:196-201.
- 850 Liu L., Pearl D.K. 2007. Species trees from gene trees: reconstructing Bayesian posterior
851 distributions of a species phylogeny using estimated gene tree distributions. *Syst. Biol.* 56:504-514.
- 852 Liu L., Yu L., Edwards S.V. 2010. A maximum pseudo-likelihood approach for estimating species
853 trees under the coalescent model. *BMC Evol. Biol.* 10:302.

- 854 Lundemo S., Falahati-Anbaran M., Stenøien H.K. 2009. Seed banks cause elevated generation
855 times, effective population sizes of *Arabidopsis thaliana* in northern Europe. *Mol. Ecol.* 18:1795-
856 2811.
- 857 Lynch M., Connery J.S. 2000. The evolutionary fate, consequences of duplicate genes. *Science*
858 290:1151-1155.
- 859 Lynch M., Force A. 2000. The probability of duplicate gene preservation by subfunctionalisation.
860 *Genetics* 154:459-473.
- 861 Maddison W.P. 1997. Gene trees in species trees. *Syst. Biol.* 46:523-536.
- 862 Maddison W.P., Knowles L.L. 2006. Inferring phylogeny despite incomplete lineage sorting. *Syst.*
863 *Biol.* 55:21-30.
- 864 Maddison W.P., Maddison D.R. 2006. Mesquite: a modular system for evolutionary analysis,
865 version 1.12.
- 866 Maere S., Bodt S.D., Raes J., Casneuf T., Montagu M.V., Kuiper M., Van de Peer, Y. 2005.
867 Modelling gene, genome duplications in eukaryotes. *Proc. Natl. Acad. Sci. U.S.A.* 102:5454–5459.
- 868 Mallet J. 2007. Hybrid speciation. *Nature* 446:279-283.
- 869 Martin D.P., Lemey P., Lott M., Moulton V., Posada D., Lefeuve P. 2010. RDP3: a flexible, fast
870 computer program for analyzing recombination. *Bioinformatics* 26:2462-2463.
- 871 Maureira-Butler I.J., Pfeil B.E., Muangprom A., Osborn T.C., Doyle J.J. 2008. The reticulate
872 history of *Medicago* (Fabaceae). *Syst. Biol.* 57:466-482.
- 873 Meng C., Kubatko L.S. 2009. Detecting hybrid speciation in the presence of incomplete lineage
874 sorting using gene tree incongruence: A model. *Theor. Pop. Biol.* 75:35-45.
- 875 Mirarab S., Reaz R., Bayzid M.S., Zimmermann T., Swenson M.S., Warnow T. 2014. ASTRAL:
876 genome-scale coalescent-based species tree estimation. *Bioinformatics* 30:i541-i548.
- 877 Nei M., Kumar S. 2000. *Molecular evolution, phylogenetics.* Oxford University Press, New York.
- 878 Nei M., Rooney A.P. 2005. Concerted, birth-and-death evolution of multigene families. *Annual*
879 *Rev. Genet.* 39:121.
- 880 Oxelman B., Yoshikawa N., McConaughy B.L., Luo J., Denton A.L., Hall B.D. 2004. RPB2 gene

- 881 phylogeny in flowering plants, with particular emphasis on asterids. *Mol. Phylogenet. Evol.* 32:462-
882 479.
- 883 Pamilo P., Nei M. 1988. Relationships between gene trees, species trees. *Mol. Biol. Evol.* 5:568-
884 583.
- 885 Pérez-Collazos E., Segarra-Moragues J.G., Villar L., Catalán P. 2015. Ant pollination promotes
886 spatial genetic structure in the long-lived plant *Borderea pyrenaica* (Dioscoreaceae). *Biol. J. Linn.*
887 *Soc.* 116:144-155.
- 888 Peters J.L., Zhuravlev Y., Fefelov I., Logie A., Omland K.E. 2007. Nuclear loci, coalescent
889 methods support ancient hybridization as cause of mitochondrial paraphyly between gadwall,
890 falcated duck (*Anas* spp.). *Evolution* 61:1992-2006.
- 891 Phillips M.J., Haouchar D., Pratt R.C., Gibb G.C., Bunce M. 2013. Inferring kangaroo phylogeny
892 from incongruent nuclear, mitochondrial genes. *PloS One* 8:e57745.
- 893 Pollard D.A., Iyer V.N., Moses A.M., Eisen M.B. 2006. Widespread discordance of gene trees with
894 species tree in *Drosophila*: evidence for incomplete lineage sorting. *PLoS Genet.* 2:1634-1647.
- 895 Ramadugu C., Pfeil B.E., Manjunath K.L., Lee R.F., Maureira-Butler I.J., Roose M.L. 2013.
896 Coalescence simulation testing of hybridization versus lineage sorting in *Citrus* (Rutaceae) using
897 six nuclear genes. *PloS One* 8:e68410.
- 898 Rambaut A., Drummond A.J. 2007. Tracer v1.4.
- 899 Rambaut A., Grass N.C. 1997. Seq-Gen: an application for the Monte Carlo simulation of DNA
900 sequence evolution along phylogenetic trees. *Computer applications in the biosciences: CABIOS*
901 13:235-238.
- 902 Rannala B., Yang Z. 2003. Bayes estimation of species divergence times, ancestral population sizes
903 using DNA sequences from multiple loci. *Genetics* 164:1645-1656.
- 904 Rasmussen M.D., Kellis M. 2012. Unified modeling of gene duplication, loss,, coalescence using a
905 locus tree. *Genome research*, 22:755-765.
- 906 Reid N., Demboski J.R., Sullivan J. 2012. Phylogeny estimation of the radiation of western north
907 American chipmunks (*Tamias*) in the face of introgression using reproductive protein genes. *Syst.*
908 *Biol.* 61:44-62.

- 909 Rieseberg L.H. 1997. Hybrid origins of plant species. *Annu. Rev. Ecol. Syst.* 28:359-389.
- 910 Rieseberg L.H., Raymond O., Rosenthal D.M., Lai Z., Livingstone K., Nakazato T., Durphy J.L.,
911 Schwarzbach A.E., Donovan L.A., Lexer C. 2003. Major ecological transitions in wild sunflowers
912 facilitated by hybridization. *Science* 301:1211-1216.
- 913 Robinson D., Foulds L. 1981. Comparison of phylogenetic trees. *Math. Biosci.* 53:131-147.
- 914 Rosenberg N.A. 2003. The shapes of neutral gene genealogies in two species: probabilities of
915 monophyly, paraphyly, polyphyly in a coalescent model. *Evolution* 57:1465-1477.
- 916 Sanderson M.J., Doyle J.J. 1992. Reconstruction of organismal, gene phylogenies from data on
917 multigene families: concerted evolution, homoplasy, confidence. *Syst. Biol.* 41:4-17.
- 918 Sankararaman S., Mallick S., Dannemann M., Prufer K., Kelso J., Paabo S., Patterson N., Reich D.
919 2014. The genomic landscape of Neanderthal ancestry in present-day humans. *Nature* 507:354–357.
- 920 Sousa F., Bertrand Y.J.K., Nylinder S., Oxelman B., Eriksson J.S., Pfeil B.E. 2014. Phylogenetic
921 properties of 50 nuclear loci in *Medicago* (Leguminosae) generated using multiplexed sequence
922 capture, next-generation sequencing. *PloS One* 9:e109704.
- 923 Sousa F., Bertrand Y.J., Pfeil B.E. 2016. Patterns of phylogenetic incongruence in *Medicago* found
924 among six loci. *Pl. Syst. Evol.* 302:493-513.
- 925 Smith B.T., Harvey M.G., Faircloth B.C., Glenn T.C., Brumfield R.T. 2014. Target capture,
926 massively parallel sequencing of ultraconserved elements for comparative studies at shallow
927 evolutionary time scales. *Syst. Biol.* 63: 83–95.
- 928 Soltis P.S., Marchant D.B., Van de Peer Y., Soltis D.E. 2015. Polyploidy, genome evolution in
929 plants, 35:119-125.
- 930 Song Y., Endepols S., Klemann N., Richter D., Matuschka F.-R., Shih C.-H., Nachman M.W.,
931 Kohn M.H. 2011. Adaptive introgression of anticoagulant rodent poison resistance by hybridization
932 between Old World mice. *Cur. Biol.* 21:1296-1301.
- 933 Steele K.P., Wojciechowski M.F. 2003. Phylogenetic analyses of tribes Trifolieae, Viciae, based
934 on sequences of the plastid gene matK (Papilionoideae: Leguminosae). *Advances Legume Syst.*
935 Part 10:355-370.
- 936 Steele K.P., Ickert-Bond S.M., Zarre S., Wojciechowski M.F. 2010. Phylogeny, character evolution

- 937 in *Medicago* (Leguminosae): Evidence from analyses of plastid trnK/matK, nuclear GA3ox1
938 sequences. *Amer. J. Bot.* 97:1142-1155.
- 939 Strasburg J.L., Rieseberg L.H. 2008. Molecular demographic history of the annual sunflowers
940 *Helianthus annuus*, *H. petiolaris* - Large effective population sizes, rates of long-term gene flow.
941 *Evolution* 62:1936-1950.
- 942 Stull G.W., Moore M.J., Mandala V.S., Douglas N.A., Kates H.R., Qi X., Brockington S.F., Soltis
943 P.S., Soltis D.E., Gitzendanner M.A. 2013. A targeted enrichment strategy for massively parallel
944 sequencing of angiosperm plastid genomes. *Applications in Plant Sciences* 1:1-7.
- 945 Sukumaran J., Holder M.T. 2010. DendroPy: a Python library for phylogenetic computing.
946 *Bioinformatics* 26:1569-1571.
- 947 Syring J., Farrell K., Businsky R., Cronn R., Liston A. 2007. Widespread genealogical
948 nonmonophyly in species of *Pinus* subgenus *Strobus*. *Syst. Biol.* 56:163-181.
- 949 Takahata N., Nei M. 1985. Gene genealogy, variance of interpopulational nucleotide differences.
950 *Genetics* 110:325-344.
- 951 Trier C.N., Hermansen J.S., Sætre G.-P., Bailey R.I. 2014. Evidence for mito-nuclear and sex-
952 linked reproductive barriers between the hybrid Italian sparrow and its parent species. *PLoS*
953 *Genetics* 10: e1004075. doi:10.1371/journal.pgen.1004075
- 954 Ungerer M.C., Baird S.J.E., Pan J., Rieseberg L.H. 1998. Rapid hybrid speciation in wild
955 sunflowers. *Proc. Natl. Acad. Sci. U.S.A.* 95:11757-11762.
- 956 Van Zee J.P., Schluter J.A., Schluter S., Dixon P., Brito Sierra C.A., Hill C.A. 2016. Paralog
957 analyses reveal gene duplication events and genes under positive selection in *Ixodes scapularis* and
958 other ixodid ticks. *BMC Gen.* 17:241 doi: 10.1186/s12864-015-2350-
- 959 Yoder J.B., Briskine R., Mudge J., Farmer A., Paape T., Steele K., Weiblen G.D., Bharti A.K.,
960 Zhou P., May G.D., Young N.D., Tiffin P. 2013. Phylogenetic signal variation in the genomes of
961 *Medicago* (Fabaceae). *Syst. Biol.* 62:424-438.
- 962 Young N.D., Debelle F., Oldroyd G.E.D., Geurts R., Cannon S.B., Udvardi M.K., Benedito V.A.,
963 Mayer K.F.X., Gouzy J., Schoof H., Van de Peer Y., Proost S., Cook D.R., Meyers B.C., Spannagl
964 M., Cheung F., De Mita S., Krishnakumar V., Gundlach H., Zhou S., Mudge J., Bharti A.K.,
965 Murray J.D., Naoumkina M.A., Rosen B., Silverstein K.A.T., Tang H., Rombauts S., Zhao P.X.,

- 966 Zhou P., Barbe V., Bardou P., Bechner M., Bellec A., Berger A., Berges H., Bidwell S., Bisseling
967 T., Choisine N., Couloux A., Denny R., Deshpande S., Dai X., Doyle J.J., Dudez A., Farmer A.D.,
968 Fouteau S., Franken C., Gibelin C., Gish J., Goldstein S., Gonzalez A.J., Green P.J., Hallab A.,
969 Hartog M., Hua A., Humphray S.J., Jeong D., Jing Y., Jocker A., Kenton S.M., Kim D., Klee K.,
970 Lai H., Lang C., Lin S., Macmil S.L., Magdelenat G., Matthews L., McCorrison J., Monaghan E.L.,
971 Mun J., Najjar F.Z., Nicholson C., Noirot C., O'Bleness M., Paule C.R., Poulain J., Prion F., Qin B.,
972 Qu C., Retzel E.F., Riddle C., Sallet E., Samain S., Samson N., Sanders I., Saurat O., Scarpelli C.,
973 Schiex T., Segurens B., Severin A.J., Sherrier D.J., Shi R., Sims S., Singer S.R., Sinharoy S., Sterck
974 L., Viollet A., Wang B., Wang K., Wang M., Wang X., Warfsmann J., Weissenbach J., White D.D.,
975 White J.D., Wiley G.B., Wincker P., Xing Y., Yang L., Yao Z., Ying F., Zhai J., Zhou L., Zuber A.,
976 Denarie J., Dixon R.A., May G.D., Schwartz D.C., Rogers J., Quetier F., Town C.D., Roe B.A.
977 2011. The *Medicago* genome provides insight into the evolution of rhizobial symbioses. *Nature*
978 480:520-524.
- 979 Yu Y., Barnett R.M., Nakhleh L. 2013. Parsimonious inference of hybridization in the presence of
980 incomplete lineage sorting. *Syst. Biol.* 62:738-751.

981

982 **Figure Legends**

983

984 **Figure 1:** Cartoon illustration of the overall test scheme. Top: the physical position of genes within
985 genomic blocks and between genomic blocks is indicated on a chromosome. Middle: individual
986 ultrametric gene trees (with sequences from species A–D at the gene tree tips) are inferred (step 1)
987 and compared within blocks (step 2). Significant differences among them are explained by gene
988 duplication (black rectangle in gene tree 1 on left) and losses (not shown). Paralogy-affected gene
989 trees (dashed box on left around gene tree 1) are discarded. Bottom: compatible genes are grouped
990 by a common origin from their parental source (if hybridisation has occurred) (step 3), both within
991 blocks (shown) and among blocks (not shown). Incompatible blocks of genes (left-hand pair versus
992 right-hand group of three genes in this example) are kept separate for principal tree estimation using
993 the multi-species coalescent model (step 4).

994

995 **Figure 2:** Hypothetical species phylogeny with a hybrid taxon, taxon 4 (panel a), that can be
996 decomposed into principal tree 1 (panel b) and principal tree 2 (panel c). Paralogous gene trees are
997 obtained by removing taxa from the clade descending from a gene duplication event, such that each
998 taxon retains only one copy of the duplicated gene (panel d). Thus, the mismatch between the gene
999 tree (panel d) and the principal tree (panel b) is due to differential gene loss after the gene
1000 duplication event.

1001

1002 **Figure 3:** Paralogy test for bins RF=4, RF=6 and RF=8 at critical value 65%, for blocks that
1003 contain a single paralogy-affected gene. When the paralogy-affected gene is *misidentified*, another
1004 gene (not harbouring paralogues) is removed instead. When the paralogy-affected gene is
1005 *undetected*, all three genes in the block are retained for down-stream analyses. *Block excluded* refers
1006 to when all three gene trees show significant differences to one another and are all excluded,
1007 because no pair of compatible genes could be identified.

1008

1009 **Figure 4:** Paralogy test at critical value 65% for blocks without paralogy-affected genes. Paralogy
1010 *rejected* means that all three genes are correctly retained for down-stream analyses. Paralogy
1011 *incorrectly detected* means that one gene was falsely identified as paralogy-affected and incorrectly
1012 excluded from further analyses.

1013

1014 **Figure 5:** Hybridisation test using only those genes and blocks that were not identified as being
1015 paralogy-affected in the previous step (cv 65% paralogy bin RF=8). Both hybridisation and earlier

1016 paralogy tests used the same critical value. Hybridisation *correctly detected* separates blocks of
1017 genes into the correct groups based on the principal tree of origin. When hybridisation is *incorrectly*
1018 *detected*, some separation of blocks of genes occurs (compatible with hybridisation), but the
1019 groupings are incorrect. *Undetected* hybridisation means that no separation of blocks of genes has
1020 occurred, even though two principal trees have generated the blocks of genes.

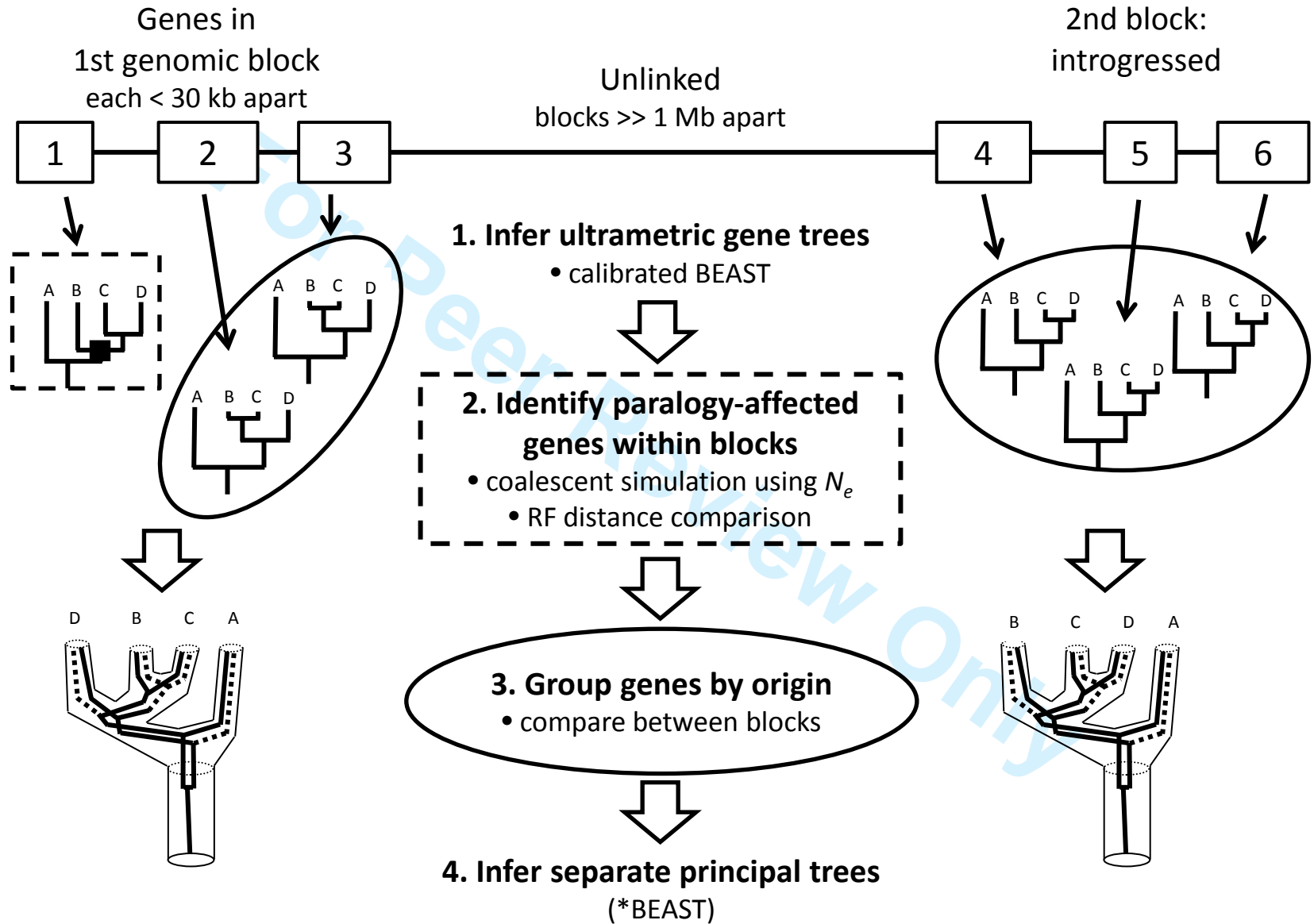
1021

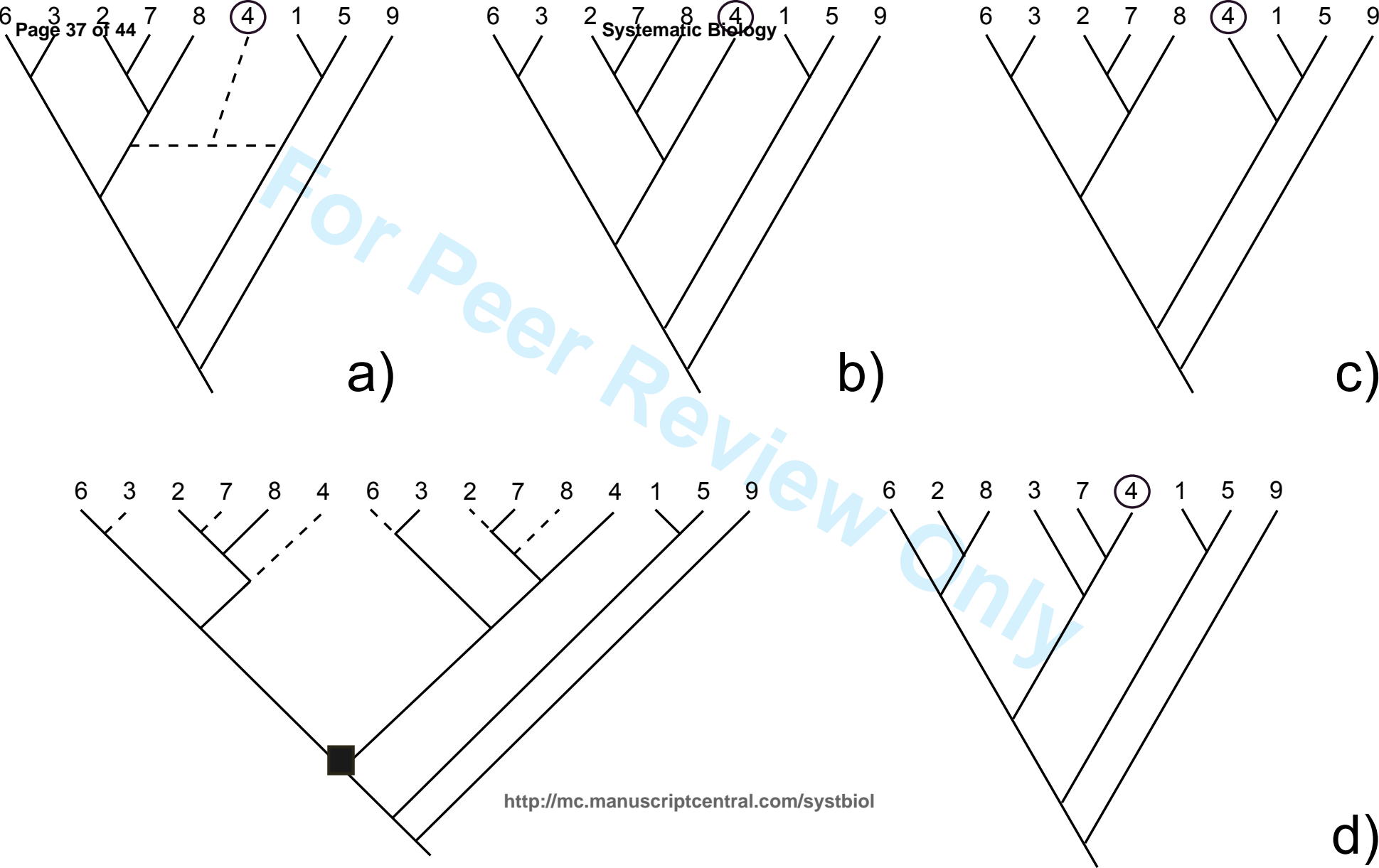
1022 **Figure 6:** The number of principal trees found after the hybridisation test (incorporating the
1023 outcome of the paralogy test) for cv 65% and paralogy bin RF=8. Both tests use the same critical
1024 value. Correct detection (Fig. 5) sometimes results in the over-separation of genomic blocks into
1025 more than two implied principal trees, although each group contains blocks of genes derived from
1026 only one principal tree each. An *incorrect grouping* places blocks together than do not originate
1027 from the same single principal tree.

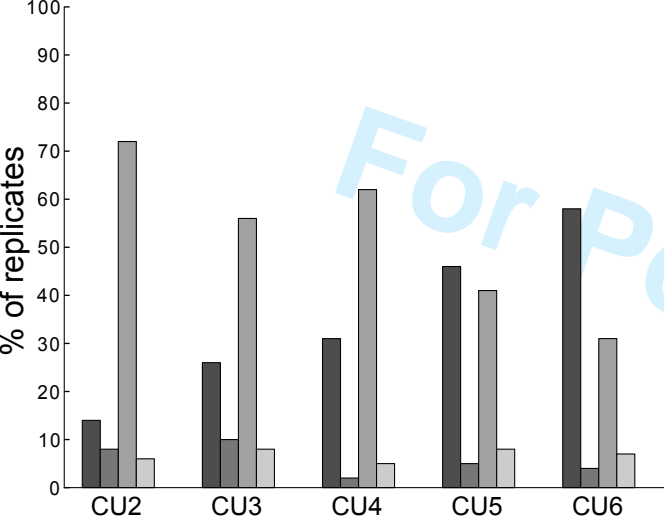
1028

1029 **Figure 7:** A schematic representation of the paralogy and hybridisation test results on genes from
1030 four genomic blocks sampled from *Medicago orbicularis* and relatives. Circles represent individual
1031 genes and dashed circles represent excluded genes (following the paralogy test). Squares represent
1032 genomic blocks. Dashed lines represent ILS rejection between genes within blocks (implying
1033 paralogy) and between blocks (implying hybridisation). Cladograms next to each block illustrate
1034 which of the two principal tree topologies were recovered with *BEAST for that block. Labels on
1035 cladograms correspond to the following taxa: TL – *M. truncatula*, *M. littoralis*; O – *M. orbicularis*;
1036 ICG – (*M. intertexta*, *M. ciliaris*, *M. granadensis*) clade; A – *M. arabica*; M – *Melilotus* (outgroup).

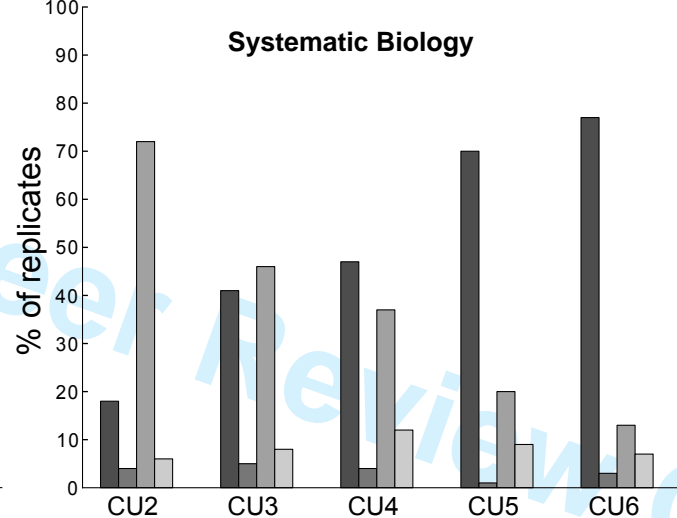
1037



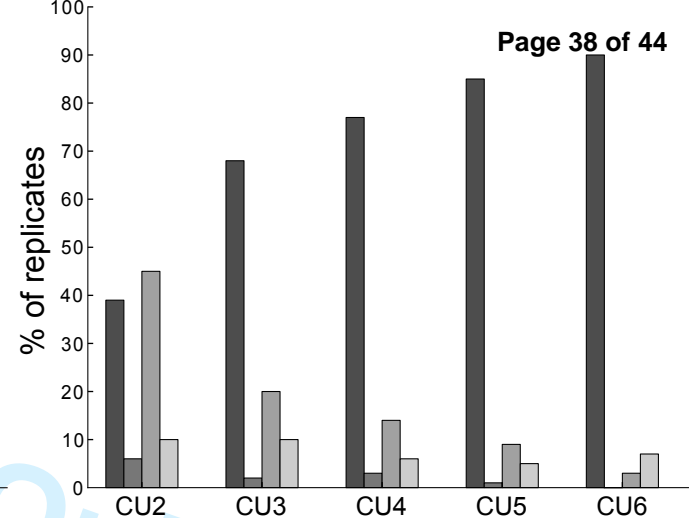




RF = 4



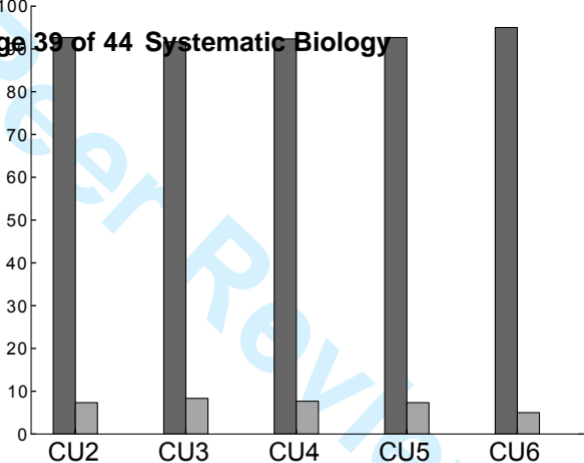
RF = 6



RF = 8

- paralogue correctly detected
- paralogue misidentified
- paralogue undetected
- block excluded

% of replicates



cv=65%

<http://mc.manuscriptcentral.com/systbiol>

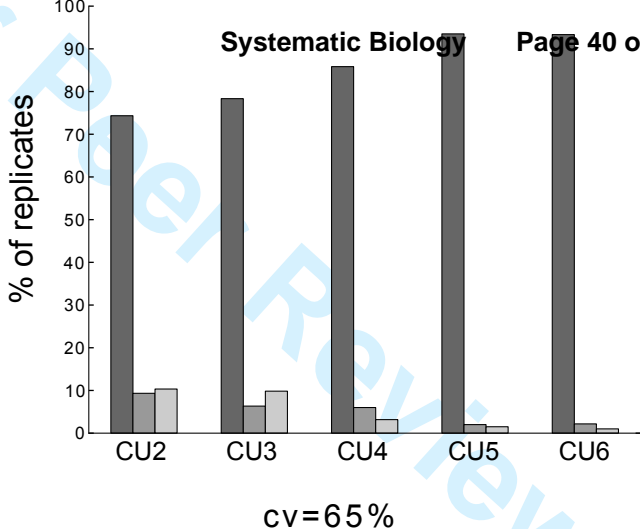
RF = 8



paralogy rejected



paralogy incorrectly detected



cv=65%

RF=8

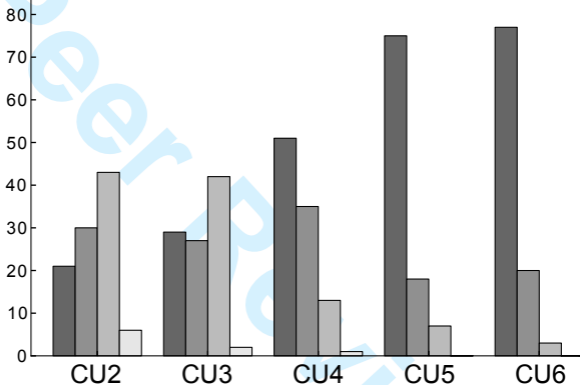
■ hybridisation correctly detected

<http://mc.manuscriptcentral.com/systbiol>

■ hybridisation incorrectly detected

■ hybridisation undetected

% of replicates



cv=65%

■ 2 principal trees

■ 3 principal trees

■ 4 principal trees

■ incorrect grouping

