

# 1 **Robustness of RADseq for evolutionary network reconstruction** 2 **from gene trees**

3  
4 José Luis Blanco-Pastor<sup>a,b</sup>, Yann J.K. Bertrand<sup>a,c</sup>, Isabel María Liberal<sup>d</sup>, Yanling Wei<sup>e</sup>, E.  
5 Charles Brummer<sup>e</sup> and Bernard E. Pfeil<sup>a</sup>

- 6  
7 a. Department of Biological and Environmental Sciences, University of Gothenburg,  
8 Box 461 40530, Göteborg, Sweden.  
9 b. INRA, Centre Nouvelle-Aquitaine-Poitiers, UR4 (URP3F), 86600 Lusignan, France  
10 c. Institute of Botany, Czech Academy of Sciences, Zámek 1, 25243 Průhonice, Czech  
11 Republic  
12 d. Real Jardín Botánico de Madrid (RJB-CSIC), Madrid, Spain  
13 e. Plant Breeding Center, Department of Plant Sciences, University of California, Davis,  
14 Davis, CA, USA.

15  
16 Correspondance

17 José Luis Blanco-Pastor,

18 INRA, Centre Nouvelle-Aquitaine-Poitiers, UR4 (URP3F), 86600 Lusignan, France.

19 Email: [jose-luis.blanco-pastor@inra.fr](mailto:jose-luis.blanco-pastor@inra.fr)

20  
21 **Keywords:** RADseq; gene trees; evolutionary networks; hybridization; *Medicago*; alfalfa

## 22 **Competing interest statement**

23 Authors have no competing interest to declare

## 24 **Abstract**

25 Although hybridization has played an important role in the evolution of many species,  
26 phylogenetic reconstructions that include hybridizing lineages have been historically  
27 constrained by the available models and data. Recently, the combined development of high-  
28 throughput sequencing and evolutionary network models offer new opportunities for  
29 phylogenetic inference under complex patterns of hybridization in the context of incomplete  
30 lineage sorting. Restriction site associated DNA sequencing (RADseq) has been a popular  
31 sequencing technique for evolutionary reconstructions of close relatives in the Next  
32 Generation Sequencing (NGS) era. However, the utility of RADseq data for the  
33 reconstruction of complex evolutionary networks has not been thoroughly discussed. Here, we  
34 used new molecular data collected from diploid perennial *Medicago* species using single-  
35 digest RADseq to reconstruct evolutionary networks from gene trees, an approach that is  
36 computationally tractable with datasets that include several species and complex patterns of  
37 hybridization. Our analyses revealed that complex network reconstructions from RADseq-  
38 derived gene trees were not robust under variations of the assembly parameters and filters.  
39 Filters to exclusively select loci with high phylogenetic information created datasets that  
40 retrieved the most anomalous topologies. Conversely, alternative clustering thresholds or  
41 filters on the number of samples *per locus* affected the level of missing data but had a lower  
42 impact on networks. When most anomalous networks were discarded, all remaining network  
43 analyses consistently supported a hybrid origin for *M. carstiensis* and *M. cretacea*.

44  
45  
46

## 47 **1. Introduction**

48 The reconstruction of a reticulate history in evolutionary close relatives has been considered  
49 from three different analytical perspectives: i) population genetic models including:  
50 approximate Bayesian Computation (Beaumont et al., 2002), full-likelihood genealogical  
51 samplers that make use of DNA sequences (Gronau et al., 2011; Hey, 2010; Sethuraman and  
52 Hey, 2016) and likelihood or pseudo-likelihood methods based on the joint allele frequency  
53 spectrum (Excoffier et al., 2013; Gutenkunst et al., 2009; Pickrell and Pritchard, 2012); ii) *D*-  
54 statistics (Durand et al., 2011; Eaton and Ree, 2013; Green et al., 2010; Meyer et al., 2012;  
55 Pease and Hahn, 2015); and iii) evolutionary network models (Solís-Lemus and Ané, 2016;  
56 Wen and Nakhleh, 2016; Yu et al., 2014, 2013; Yu and Nakhleh, 2015; Zhang et al., 2018;  
57 Zhu et al., 2017). The first two perspectives assume a previously known backbone phylogeny  
58 to formulate a hypothesis of hybridization. This backbone tree is usually constructed either  
59 using i) a total evidence approach with concatenation of full sequence information (Eaton and  
60 Ree, 2013; Escudero et al., 2014; Fernández-Mazuecos et al., 2017; Hipp et al., 2014; Wagner  
61 et al., 2013) or ii) coalescent based methods (Eaton and Ree, 2013; Fernández-Mazuecos et  
62 al., 2017; Rheindt et al., 2014) that reconcile individual gene trees. Despite being a standard  
63 approach, the construction of a backbone tree could be an incorrect representation of the main  
64 evolutionary history of the species under complex reticulate evolution (Clark and Messer,  
65 2015; Huson et al., 2010; Yu et al., 2011), or molecular data can show a different “main”  
66 phylogeny when hybridization is first accounted for (Sousa et al., 2017).

67  
68 Evolutionary networks provide an explicit model of evolutionary relationships that extends  
69 the tree model to allow for reticulations with internal nodes representing ancestral species.  
70 Recently developed phylogenetic network reconstruction methods are based on maximum  
71 parsimony (Yu et al., 2013), maximum likelihood (ML) (Yu et al., 2014), maximum pseudo-  
72 likelihood (Solís-Lemus and Ané, 2016; Yu and Nakhleh, 2015) and Bayesian inference (BI)  
73 methods (Wen et al., 2016; Wen and Nakhleh, 2016; Zhang et al., 2018; Zhu et al., 2017).  
74 Although ML and BI methods show promise, they are still limited to small datasets (usually  
75 fewer than 10 individuals and less than 3 reticulations, Yu and Nakhleh, 2015). In contrast,  
76 maximum pseudo-likelihood (summary) methods (Solís-Lemus and Ané, 2016; Yu and  
77 Nakhleh, 2015) are nowadays a convenient alternative for complex empirical datasets (Wen et  
78 al., 2017).

79  
80 RADseq approaches (reviewed in Andrews et al., 2016) are widely used sequencing  
81 techniques for evolutionary reconstructions in the Next Generation Sequencing era. RADseq  
82 was first envisioned as a technique to find intraspecific genetic variation (Baird et al., 2008;  
83 Elshire et al., 2011; Hohenlohe et al., 2011). Later RADseq methods have been considered  
84 suitable for phylogenetic studies from shallow to deep timescales (Cariou et al., 2013; Eaton,  
85 2014; Harvey et al., 2016; Rubin et al., 2012). RADseq are particularly appealing for  
86 systematics because they are easily applied to non-model organisms for which no reference  
87 genome or previous genomic information is available (Cariou et al., 2013; Fernández-  
88 Mazuecos et al., 2017; Rubin et al., 2012). For that reason RADseq has become a very  
89 popular technique for hybridization studies across a diversity of organisms and timescales  
90 (Escudero et al., 2014; Fernández-Mazuecos et al., 2017). Nevertheless, because RADseq  
91 datasets are limited in sequence length, contain relatively few variable sites, and do not  
92 generally yield resolved gene trees (Rubin et al., 2012), it is unknown if they are appropriate  
93 for maximum pseudo-likelihood phylogenetic network reconstruction methods. Their intrinsic

94 characteristics suggest that these datasets are limited for network inference from gene trees,  
95 but an in-depth evaluation of their utility is still lacking.

96

97 *Medicago* L. (Fabaceae) is a genus comprising 87 species (Small, 2011) and includes the  
98 economically important forage crop alfalfa (*M. sativa*, section *Medicago*) in addition to the  
99 model legume *M. truncatula* (Barker et al., 1990; Benedito et al., 2008; Branca et al., 2011;  
100 Cook, 1999; Young et al., 2011). The genus *Medicago* L. exhibits severe phylogenetic gene  
101 tree incongruence that has been mainly attributed to hybridization and ILS (Eriksson et al.,  
102 2018, 2017; Maureira-Butler et al., 2008; Sousa et al., 2017, 2014; Steele et al., 2010; Yoder  
103 et al., 2013). We collected new molecular data from diploid perennial *Medicago* species using  
104 single-digest RADseq (Genotyping-By-Sequencing, Elshire et al., 2011). We investigated the  
105 ability of RADseq data to unveil the evolutionary history of diploid species of *Medicago*  
106 section *Medicago* using a network reconstruction method that uses gene trees (Yu and  
107 Nakhleh, 2015). Specifically, we investigated the robustness of this method (i.e. the  
108 propensity to retrieve a set of optimal networks with similar topologies) under a variety of  
109 RADseq data assembly parameters and filters.

## 110 **2. Materials and Methods**

### 111 **2.1 Sampling**

112 Our choice of species (Table 1) was based on results of previous studies grouping diploid  
113 perennial *Medicago* taxa (Bena, 2001; Maureira-Butler et al., 2008; Sousa, 2015; Yoder et al.,  
114 2013). These includes: *M. marina*, *M. cretacea*, *M. rhodopea*, *M. prostrata*, *M. daghestanica*  
115 and *M. sativa* (section *Medicago* subsection *Medicago*), *M. hybrida* and *M. suffruticosa*  
116 (section *Medicago* subsection *Suffruticosae*); *M. carstiensis* (section *Carstiensae*); *M. rugosa*  
117 and *M. scutellata* (section *Spirocarpos* subsection *Rotatae*). As outgroup we used the annual  
118 species *M. truncatula* (section *Spirocarpos* subsection *Pachyspirae*).  
119

### 120 **2.2 Sequence preparation**

121 We extracted genomic DNA with a custom CTAB DNA Extraction Protocol and constructed  
122 a genotyping-by-sequencing (GBS) library following the library preparation protocol of  
123 Elshire et al. (2011) with minor modifications as described by Annicchiarico et al. (2017). In  
124 brief, GBS library was prepared using the frequent cutter ApeKI (R0643L; NEB) restriction  
125 enzyme. Sets of 8-bp barcoded adapters were ligated to restriction fragments for multiplex  
126 sequencing. The QIAquick PCR purification kit (28104; QIAGEN) was used to purify equal  
127 volumes of the pooled ligated products previous to the final PCR amplification step with the  
128 Kapa Library Amplification Readymix (Kapa Biosystems KK2611. Sequences were obtained  
129 at the Genomic Core Facility of the UT Southwestern Medical Center (Dallas, TX) with an  
130 Illumina HiSeq 2500 system that generated 100-bp single-end reads. This protocol was  
131 chosen based on comparisons made among a few protocols and different enzymes, including  
132 the two-enzyme protocol by Poland et al. (2012) and the 2b-RAD protocol by Wang et al.  
133 (2012). The decision was made based on the number of sites genotyped that were shared  
134 among representative individuals (Annicchiarico et al., 2017).

135 Raw single-end sequence reads were trimmed of adapter sequence and filtered with a  
136 minimum quality score of 20 using trimmomatic (Bolger et al., 2014). Assembly was then  
137 performed using ipyrad v. 0.7.19 (<http://ipyrad.readthedocs.io/>), a toolbox for reproducible  
138 assembly and analysis of RADseq type genomic data sets based on the pyRAD pipeline  
139 (Eaton, 2014). Assembly consisted of seven sequential steps, with parameters based on those  
140 recommended for single-end GBS data in the ipyrad documentation. We used the *de novo* +  
141 reference method, with the *M. truncatula* genome sequence (Mt4.0,

142 <http://www.medicagohapmap.org>) as a reference. Briefly, the steps of the ipyrad pipeline are  
143 described as follows: In step 1, sequences were demultiplexed according to barcode  
144 sequences. In step 2, low quality reads and Illumina adapters were filtered out. Step 3  
145 removed amplification duplicates and then clustered reads within each sample according to a  
146 clustering threshold. This step tries to identify all the reads that map to the same locus within  
147 each sample. As we used the *de novo* + reference method, the *M. truncatula* reference was  
148 used to identify homology, and then the remaining unmatched sequences were clustered with  
149 the standard *de novo* ipyrad pipeline. Because phylogenetic results are known to be sensitive  
150 to the similarity threshold employed in step 3 for within-sample and step 6 (see below) for  
151 across-sample sequence clustering (Fernández-Mazuecos et al., 2017; Leaché et al., 2015;  
152 Shafer et al., 2017; Takahashi et al., 2014), five assemblies of GBS loci were generated using  
153 a range of clustering thresholds (clust parameter) from  $c=0.75$  to  $c=0.95$  (Table 2). Step 4  
154 jointly estimated the error rate and heterozygosity to differentiate “good” reads from  
155 sequencing errors. Step 5 called the consensus of sequences within each cluster. Step 6  
156 clustered consensus sequences across samples. Step 7 filtered the data and wrote output files.  
157 In step 7 we applied filters for the maximum number of indels *per locus* (8), max  
158 heterozygosity *per locus* (50% of samples) and max number of SNPs *per locus* (20).  
159 To evaluate the effect of missing data on network inference, for each assembly we generated  
160 datasets with two alternative values for the minimum number of samples *per locus*  
161 (“minimum taxon coverage” -min- parameter, 4 and 10). The effect of locus variation on  
162 networks was tested by generating datasets with two alternative values for the minimum  
163 number of parsimony-informative sites (PIS parameter, 4 and 10). We saved the data in the  
164 ipyrad format (\*.loci) that was later on transformed in individual alignment files *per locus* in  
165 the phylip format using a custom R script. We obtained 20 RADseq datasets under different  
166 combinations of assembly parameters and filters described above (Table 2).  
167

### 168 **2.3 Network inference**

169 We analyzed RADseq datasets alignments with PhyloNet (Than et al., 2008; Wen et al.,  
170 2017). Within PhyloNet we applied the method that infers species networks from gene trees  
171 using maximum pseudo-likelihood (InferNetwork\_MPL command; Yu and Nakhleh, 2015)  
172 which is computationally fast. First, we analyzed separate sets of genes from each of the 20  
173 RADseq datasets with RaxML v.7.2.8 (Stamatakis, 2006) using the GTRCAT substitution  
174 model and using *M. truncatula* as outgroup. We sampled several individuals/alleles for some  
175 species (see Table 1) that were mapped to single taxa with the -a parameter. Ten optimal  
176 networks were returned with the -n parameter. We chose 5 maximum allowed number of  
177 reticulation events. Remaining parameter values were set as default.  
178

### 179 **2.4 Network distances**

180 To investigate dissimilarities between evolutionary networks computed with alternative  
181 RADseq datasets we used multidimensional scaling. We first calculated a matrix of distances  
182 among networks computed with the topological dissimilarity measure of Nakhleh (2010)  
183 (normalized to get values within [0, 1]), which is implemented in PhyloNet. Then we applied  
184 a Principal Coordinate Analysis (PCoA) to transform the distance matrix into a set of  
185 coordinates that were plotted to display network distances. We performed the PCoA using all  
186 pairwise distances between every network returned by the PhyloNet analyses.

## 187 **3. Results**

### 188 **3.1 Sequence capture and RADseq data**

189 Among the 20 RADseq datasets the number of loci ranked from 4 (clust95.min10.PIS10) to  
190 3,405 (clust85.min4.PIS4), concatenated length (bp) ranged from 367 (clust95.min10.PIS10)  
191 to 303,272 (clust85.min4.PIS4) and missing data (%) ranged from 16.2 (clust95.min10.PIS10)  
192 to 56.6 (clust75.min4.PIS10).

193

### 194 **3.2 Phylogenetic networks**

195 Best networks (networks with highest likelihood scores) for each of the 20 RADseq datasets  
196 showed marked differences (Fig. 1). A hybrid origin was recovered for all species (excluding  
197 the outgroup species, *M. truncatula*) at least in one of the 20 best species networks (Table 3):  
198 *M. carstiensis* (observed as hybrid in 18 networks), *M. cretacea* (in 16 networks), *M.*  
199 *rhodopea* (in 10 networks), *M. marina* (in 6 networks), *M. rugosa* (in 4 networks), *M.*  
200 *scutellata* (in 4 networks), *M. daghestenica* (in 4 networks), *M. suffruticosa* (in 4 networks),  
201 *M. prostrata* (in 2 networks), *M. hybrida* (in 2 networks) and *M. sativa* (in 2 networks).

202

### 203 **3.3 Network distances**

204 The RADseq datasets that retrieved the highest distances from the “core” set of networks  
205 were those that were computed with datasets filtered to contain only the most variable loci  
206 (PIS10 filter, see Fig. 2). In general, these datasets contained a low number of loci and short  
207 concatenated sequence lengths. The PCoA did not show a marked effect of the filter on the  
208 minimum number of samples *per* locus (min filter) or the use alternative clustering thresholds  
209 (clust parameter).

210 After excluding datasets with the PIS10 filter, a hybrid origin was recovered for eight species  
211 at least in one of the remaining 10 best species networks: *M. carstiensis* (observed as hybrid  
212 in all 10 networks), *M. cretacea* (in all 10 networks), *M. rugosa* (in 4 networks), *M. rhodopea*  
213 (in 3 networks), *M. marina* (in 2 networks), *M. suffruticosa* (in 2 network), *M. scutellata* (in 1  
214 network) and *M. sativa* (in 1 network).

## 215 **4. Discussion**

216 Our empirical comparison among networks computed from the RADseq datasets reveal some  
217 general patterns in how assembly parameters and filters influence complex evolutionary  
218 network reconstructions from gene trees. Our study shows that RADseq datasets with a low  
219 number of loci retrieve the most atypical network topologies, regardless the high phylogenetic  
220 information contained in the loci. The RADseq networks that were the closest to the core set  
221 of networks were those that assembled the highest number of loci with very little impact on  
222 the clustering threshold or the minimum number of samples *per* locus and therefore with very  
223 little impact on the level of missing data. In general RADseq datasets showed low robustness  
224 (different best network topologies) under variation of the assembly parameters and filters.  
225 But, after excluding the most divergent networks, all remaining analyses supported a hybrid  
226 origin for two species: *M. carstiensis* and *M. cretacea*.

227

228 Recently Fernández-Mazuecos et al. (2017) showed a high robustness of coalescent  
229 approaches for RADseq-based species trees reconstructions. Contrastingly, here we observed  
230 variation among network topologies under variations of the assembly parameters and filters  
231 underlining the importance of RADseq data preparation on the final results. RADseq is a  
232 particularly appealing technique for systematics because of their potential for detecting both  
233 current and historical hybridization (Escudero et al., 2014; Twyford and Ennos, 2012) and

234 because they are easily applied with no previous genomic information and reduced lab costs.  
235 The pseudo-likelihood method of Yu and Nakhleh (2015) is also attractive because it does not  
236 require heavy computational resources. Nevertheless its use on RADseq data may produce  
237 misleading results without a proper evaluation of the optimal assembly parameters and filters.  
238 In phylogenetic analysis with RADseq, it is particularly challenging to establish general  
239 criteria for determining the assembly parameters that maximize the number of orthologous  
240 RAD sequences between samples and filtering parameters that retain loci with the optimal  
241 level of missing data or phylogenetic information. It has been suggested that low phylogenetic  
242 resolution of loci may constrain the identification of hybrids because poorly resolved gene  
243 trees, constructed from markers with limited sequence divergence between species, are likely  
244 to be uninformative in tracing the reticulate history of species (Linder and Rieseberg, 2004;  
245 Twyford and Ennos, 2012). In contrast our study suggests that high loci number increases the  
246 power for network inference from RADseq-gene trees despite the low phylogenetic  
247 information contained within each individual locus. Additionally, selectively choosing the  
248 most variable RADseq dataset may be detrimental as these loci may introduce potential biases  
249 typical of hypervariable regions of the genome. Indeed, the most variable regions could be  
250 those retaining ancestral polymorphisms or those representing regions of introgressed DNA  
251 (Eaton and Ree, 2013).

252  
253 In recent years RADseq has been applied for the evolutionary reconstruction of complex  
254 taxonomic groups (Eaton and Ree, 2013; Escudero et al., 2014; Fernández-Mazuecos et al.,  
255 2017; Hipp et al., 2014; Wagner et al., 2013). Most previous studies using RADseq data relied  
256 on a “backbone tree” and placed a limited number of hybridization events upon it.  
257 Nevertheless it is known that this approach could provide an incorrect representation of the  
258 evolutionary history of the species under complex reticulate evolution with multiple  
259 hybridization events (Huson et al., 2010; Yu et al., 2011). New tools for evolutionary network  
260 reconstructions (Solís-Lemus and Ané, 2016; Wen et al., 2016; Wen and Nakhleh, 2016; Yu  
261 et al., 2014, 2013; Yu and Nakhleh, 2015; Zhang et al., 2018; Zhu et al., 2017) now offer the  
262 opportunity to study reticulate evolution including cases with multiple hybridization events  
263 and with no previous information on the “backbone tree” or where such main tree is  
264 potentially non-existent. Development of such evolutionary network models are now in full  
265 swing and should become standard methods for phylogenetic inference under incomplete  
266 lineage sorting (ILS) and hybridization. Despite these remarkable methodological advances,  
267 in the most complex cases computational limitations reduce the set of methods to those using  
268 maximum pseudo-likelihood inference of networks from gene trees (Solís-Lemus and Ané,  
269 2016; Yu and Nakhleh, 2015). These methods have a great potential but there is no  
270 information in the literature about the adequacy of the commonly used RADseq datasets for  
271 the estimation of evolutionary networks using these type of analyses were a previous  
272 computation of gene trees is required. In general using RADseq for gene tree reconstruction  
273 poses a number of potential problems: the orthology relationships among sequences are  
274 unknown, mutations on restriction sites is expected to yield missing data that increases with  
275 evolutionary time and the genetic linkage relationships among loci are unknown (see Rubin et  
276 al., 2012). Additionally, given the short length of sequences, the phylogenetic information of  
277 each locus is very scarce and recombination detection is not straightforward.

278  
279 Despite the varied result obtained with different RADseq datasets, general patterns emerged  
280 regarding the identification of hybrid species which were more evident when the PIS10  
281 datasets were excluded. A hybrid origin was retrieved by all remaining PIS4 datasets for *M.*  
282 *carstiensis* and *M. cretacea*. This signal was clearly stronger than the hybridization signal  
283 detected for the remaining species (hybrid origin detected in  $\leq 4$  datasets for *M. rugosa*, *M.*

284 *rhodopea*, *M. suffruticosa*, *M. marina*, *M. scutellata* and *M. sativa*). *M. carstiensis* is the only  
285 *Medicago* species exclusively with rhizomes (which are found only sporadically in a few  
286 other species, especially *M. sativa*). Phylogenetic relationships around *M. carstiensis* has been  
287 enigmatic as it forms a monospecific section (Carstiensae, Small, 2011) and previous  
288 phylogenetic studies did not provide well-supported information on the relationships of this  
289 species (Maureira-Butler et al., 2008; Small, 2011). It has been speculated that *M. carstiensis*  
290 is a relic species that is ancestral to the much more widespread *M. orbicularis* (Bennett et al.,  
291 2006). But its particular characteristics could be also explained by speciation after disruptive  
292 selection on hybrids (Seehausen, 2004). Regarding *M. cretacea*, Urban (1873) (the first to  
293 prepare a comprehensive analysis of the genus *Medicago*) already considered that this species  
294 had controversial affinities. Lesins and Lesins (Lesins and Lesins, 1979), in the second  
295 comprehensive systematic analysis of *Medicago*, already included *M. cretacea* in the  
296 monotypic section *Cretaceae*. Later on, analyses by Bena (2001) and Maureira-Butler et al.  
297 (2008) showed alternative inconsistent phylogenetic relationships for *M. cretacea*. These  
298 contentious taxonomic and phylogenetic placement of *M. carstiensis* and *M. cretacea*  
299 observed in previous studies are consistent with hybridisation.

## 300 **5. Conclusions**

301 Here we inferred a hybrid origin for *M. carstiensis* and *M. cretacea* using RADseq data and a  
302 maximum pseudo-likelihood approach for network inference from gene trees. We observed  
303 that loci number had an important impact on network reconstruction from RADseq-gene trees,  
304 whereas the clustering threshold used in the data assembly or a filter on taxon coverage had a  
305 lower impact on network inference. Future research on methods that explore the parameter  
306 space for optimal assembly parameters and filters may be required to obtain a clear  
307 phylogenetic picture of all diploid perennial *Medicago* species and to consider these  
308 approaches sufficiently robust for their standard use in the phylogenetics community.  
309

310 **Tables**

311 Table 1. Information on the *Medicago* samples used in the present study.

Species	Accession No	Germplasm Bank	Country of origin	Accession European Nucleotide Database	Chromosome number <sup>d</sup>
<i>M. carstiensis</i>	PI641414	USDA	Russian Federation		16 <sup>a</sup>
<i>M. carstiensis</i>	MED152/91	Ösnabruck Botanic Garden	NA		16 <sup>a</sup>
<i>M. cretacea</i>	PI631721	USDA	Russian Federation		16 <sup>b</sup>
<i>M. daghestanica</i>	Eoo337403; 1653*1	GB herbarium	Russian Federation		16 <sup>a</sup>
<i>M. daghestanica</i>	Eoo337402; 1653*2	GB herbarium	Russian Federation		16 <sup>a</sup>
<i>M. hybrida</i>	PI538998	USDA	Russian Federation		<sup>d</sup>
<i>M. marina</i>	PI516711	USDA	Morocco		16 <sup>a</sup>
<i>M. marina</i>	PI419391	USDA	Greece		<sup>d</sup>
<i>M. prostrata</i>	PI577445	USDA	Italy		16 <sup>b</sup>
<i>M. rhodopea</i>	SA43026	SARDI	NA		16 <sup>a</sup>
<i>M. rhodopea</i>	W619154	USDA	Bulgaria		16 <sup>a</sup>
<i>M. rugosa</i>	PI487386	USDA	Tunisia		<sup>d</sup>
<i>M. sativa</i> <sup>c</sup>	PI577567	USDA	Italy		16 <sup>a</sup>
<i>M. scutellata</i>	PI505433	USDA	Spain		16 <sup>a</sup>
<i>M. suffruticosa</i>	PI516913	USDA	Morocco		<sup>d</sup>
<i>M. suffruticosa</i>	W64952	USDA	Spain		16 <sup>a</sup>
<i>M. truncatula</i>	W64996	USDA	Greece		<sup>d</sup>

312 <sup>a</sup>, flow cytometry; <sup>b</sup>, chromosome counts; <sup>c</sup>, five individuals used; <sup>d</sup>, samples confirmed as  
313 diploid by inspection of the phased samtools files (sam files) for all eight genes using Tablet  
314 (Milne et al., 2013, 2010).

315  
316



317 Table 2. Characteristics of RADseq datasets generated in ipyrad and used for gene tree and  
318 network inference.

Dataset	Clustering threshold	Minimum taxon coverage	Minimum PIS <i>per</i> loci	Number of loci	Concatenated length (bp)	Missing data (%)
clust75.min4.PIS4	0.75	4	4	3194	283,753	49.7
clust75.min4.PIS10	0.75	4	10	346	29,974	56.6
clust75.min10.PIS4	0.75	10	4	1756	153,733	34.4
clust75.min10.PIS10	0.75	10	10	140	12,274	41.5
clust80.min4.PIS4	0.80	4	4	3364	299,468	49.5
clust80.min4.PIS10	0.80	4	10	317	28,744	55.3
clust80.min10.PIS4	0.80	10	4	1856	162,504	34.0
clust80.min10.PIS10	0.80	10	10	134	11,846	40.4
clust85.min4.PIS4	0.85	4	4	3405	303,272	48.8
clust85.min4.PIS10	0.85	4	10	205	18,789	51.8
clust85.min10.PIS4	0.85	10	4	1910	167,477	33.8
clust85.min10.PIS10	0.85	10	10	98	8,711	38.4
clust90.min4.PIS4	0.90	4	4	3036	271,409	48.0
clust90.min4.PIS10	0.90	4	10	81	7,451	45.5
clust90.min10.PIS4	0.90	10	4	1729	151,806	33.3
clust90.min10.PIS10	0.90	10	10	40	3,616	33.2
clust95.min4.PIS4	0.95	4	4	1506	135,056	46.1
clust95.min4.PIS10	0.95	4	10	5	455	24.2
clust95.min10.PIS4	0.95	10	4	912	80,449	32.7
clust95.min10.PIS10	0.95	10	10	4	367	16.2

319

320  
321

Table 3. Positive hybridization signal detected for each taxon in each network. Only strict hybridization signal is considered, i.e. a taxa nested within a hybrid clade but represented with a single branch is not considered of hybrid origin.

Dataset	<i>M. carstiensis</i>	<i>M. cretacea</i>	<i>M. daghestanica</i>	<i>M. hybrida</i>	<i>M. marina</i>	<i>M. prostrata</i>	<i>M. rhodopea</i>	<i>M. rugosa</i>	<i>M. sativa</i>	<i>M. scutellata</i>	<i>M. suffruticosa</i>	<i>M. truncatula</i>
clust75.min4.PIS4	✓	✓						✓				
clust75.min4.PIS10	✓						✓			✓		
clust75.min10.PIS4	✓	✓					✓	✓				
clust75.min10.PIS10	✓	✓	✓		✓		✓					
clust80.min4.PIS4	✓	✓									✓	
clust80.min4.PIS10	✓						✓			✓		
clust80.min10.PIS4	✓	✓					✓					
clust80.min10.PIS10	✓	✓				✓	✓				✓	
clust85.min4.PIS4	✓	✓									✓	
clust85.min4.PIS10	✓		✓		✓	✓						
clust85.min10.PIS4	✓	✓										
clust85.min10.PIS10		✓		✓	✓				✓			
clust90.min4.PIS4	✓	✓									✓	
clust90.min4.PIS10	✓	✓	✓		✓							
clust90.min10.PIS4	✓	✓			✓		✓					
clust90.min10.PIS10	✓						✓				✓	
clust95.min4.PIS4	✓	✓			✓			✓				
clust95.min4.PIS10	✓	✓		✓			✓					
clust95.min10.PIS4	✓	✓						✓	✓			
clust95.min10.PIS10		✓	✓				✓			✓		

322 **Figure Legends**

323 Fig. 1 - Best networks (networks with highest likelihood scores) for each of the 20 RADseq  
324 datasets.

325

326 Fig. 2 - PCoA showing pairwise network distances calculated with the topological  
327 dissimilarity measure of Nakhleh (2010). The figure show pairwise distances between the 10  
328 best networks of each of the 20 RADseq datasets. Transparency represents the filter on  
329 parsimony informative sites, shapes represent the filter on min. samples locus, and colors  
330 represent the clustering threshold used to generate the RADseq dataset.

331

332 Fig. 3 - Bar chart representing number of PIS4 RADseq datasets supporting a hybrid origin  
333 for the *Medicago* species analyzed in this study.

334

## 335 **Acknowledgements**

336 The authors thank Luay Nakhleh and Jiafan Zhu for their assistance with PhyloNet analyses  
337 and Filipe de Sousa for providing plant material.

## 338 **Funding**

339 This work was supported by a grant from the Swedish Research Council (grant reference  
340 2009-5206) and by the Marie Curie Intra-European Fellowship “AlfalfaEvolution” (FP7-  
341 PEOPLE-2013-IEF, project reference 625308).

## 342 **References**

- 343 Andrews, K.R., Good, J.M., Miller, M.R., Luikart, G., Hohenlohe, P.A., 2016. Harnessing the  
344 power of RADseq for ecological and evolutionary genomics. *Nat. Rev. Genet.* 17, 81–  
345 92. <https://doi.org/10.1038/nrg.2015.28>
- 346 Annicchiarico, P., Nazzicari, N., Wei, Y., Pecetti, L., Brummer, E.C., 2017. Genotyping-by-  
347 Sequencing and Its Exploitation for Forage and Cool-Season Grain Legume Breeding.  
348 *Front. Plant Sci.* 8, 679. <https://doi.org/10.3389/fpls.2017.00679>
- 349 Baird, N.A., Etter, P.D., Atwood, T.S., Currey, M.C., Shiver, A.L., Lewis, Z.A., Selker, E.U.,  
350 Cresko, W.A., Johnson, E.A., 2008. Rapid SNP discovery and genetic mapping using  
351 sequenced RAD markers. *PLoS One* 3, e3376.  
352 <https://doi.org/10.1371/journal.pone.0003376>
- 353 Barker, D.G., Bianchi, S., Blondon, F., Dattée, Y., Duc, G., Essad, S., Flament, P., Gallusci,  
354 P., Génier, G., Guy, P., Muel, X., Tourneur, J., Dénarié, J., Huguet, T., 1990. *Medicago*  
355 *truncatula*, a model plant for studying the molecular genetics of the Rhizobium-legume  
356 symbiosis. *Plant Mol. Biol. Report.* 8, 40–49. <https://doi.org/10.1007/BF02668879>
- 357 Beaumont, M.A., Zhang, W., Balding, D.J., 2002. Approximate Bayesian Computation in  
358 Population Genetics. *Genetics* 162.
- 359 Bena, G., 2001. Molecular phylogeny supports the morphologically based taxonomic transfer  
360 of the “medicagoid” *Trigonella* species to the genus *Medicago* L. *Plant Syst. Evol.* 229,  
361 217–236. <https://doi.org/10.1007/s006060170012>
- 362 Benedito, V.A., Torres-Jerez, I., Murray, J.D., Andriankaja, A., Allen, S., Kakar, K.,  
363 Wandrey, M., Verdier, J., Zuber, H., Ott, T., Moreau, S., Niebel, A., Frickey, T., Weiller,  
364 G., He, J., Dai, X., Zhao, P.X., Tang, Y., Udvardi, M.K., 2008. A gene expression atlas  
365 of the model legume *Medicago truncatula*. *Plant J.* 55, 504–513.  
366 <https://doi.org/10.1111/j.1365-313X.2008.03519.x>
- 367 Bennett, S.J., Broughton, D.A., Maxted, N., 2006. Ecogeographical analysis of the perennial  
368 *Medicago*. *CRC for Plant-Based Management of Dryland Salinity*.
- 369 Bolger, A.M., Lohse, M., Usadel, B., 2014. Trimmomatic: A flexible trimmer for Illumina  
370 sequence data. *Bioinformatics* 30, 2114–2120.  
371 <https://doi.org/10.1093/bioinformatics/btu170>
- 372 Branca, A., Paape, T.D., Zhou, P., Briskine, R., Farmer, A.D., Mudge, J., Bharti, A.K.,  
373 Woodward, J.E., May, G.D., Gentzittel, L., Ben, C., Denny, R., Sadowsky, M.J.,  
374 Ronfort, J., Bataillon, T., Young, N.D., Tiffin, P., 2011. Whole-genome nucleotide  
375 diversity, recombination, and linkage disequilibrium in the model legume *Medicago*  
376 *truncatula*. *Proc. Natl. Acad. Sci.* 108, E864–E870.  
377 <https://doi.org/10.1073/pnas.1104032108>
- 378 Cariou, M., Duret, L., Charlat, S., 2013. Is RAD-seq suitable for phylogenetic inference? An  
379 in silico assessment and optimization. *Ecol. Evol.* 3, 846–852.

- 380 <https://doi.org/10.1002/ece3.512>
- 381 Clark, A.G., Messer, P.W., 2015. Conundrum of jumbled mosquito genomes. *Science*. 347,  
382 27–28. <https://doi.org/10.1126/science.aaa3600>
- 383 Cook, D.R., 1999. *Medicago truncatula* - A model in the making! *Curr. Opin. Plant Biol.*  
384 [https://doi.org/10.1016/S1369-5266\(99\)80053-3](https://doi.org/10.1016/S1369-5266(99)80053-3)
- 385 Durand, E.Y., Patterson, N., Reich, D., Slatkin, M., 2011. Testing for ancient admixture  
386 between closely related populations. *Mol. Biol. Evol.* 28, 2239–2252.  
387 <https://doi.org/10.1093/molbev/msr048>
- 388 Eaton, D.A.R., 2014. PyRAD: assembly of de novo RADseq loci for phylogenetic analyses.  
389 *Bioinformatics* 30, 1844–1849. <https://doi.org/10.1093/bioinformatics/btu121>
- 390 Eaton, D.A.R., Ree, R.H., 2013. Inferring Phylogeny and Introgression using RADseq Data:  
391 An Example from Flowering Plants (Pedicularis: Orobanchaceae). *Syst. Biol.* 62, 689–  
392 706. <https://doi.org/10.5061/dryad.bn281>
- 393 Elshire, R.J., Glaubitz, J.C., Sun, Q., Poland, J.A., Kawamoto, K., Buckler, E.S., Mitchell,  
394 S.E., 2011. A robust, simple genotyping-by-sequencing (GBS) approach for high  
395 diversity species. *PLoS One* 6, e19379. <https://doi.org/10.1371/journal.pone.0019379>
- 396 Eriksson, J.S., De Sousa, F., Bertrand, Y.J.K., Antonelli, A., Oxelman, B., Pfeil, B.E., 2018.  
397 Allele phasing is critical to revealing a shared allopolyploid origin of *Medicago arborea*  
398 and *M. strasseri* (Fabaceae). *BMC Evol. Biol.* 18, 9. <https://doi.org/10.1186/s12862-018-1127-z>
- 400 Eriksson, J.S.S., Blanco-Pastor, J.L.L., Sousa, F., Bertrand, Y.J.K.J.K., Pfeil, B.E.E., 2017. A  
401 cryptic species produced by autopolyploidy and subsequent introgression involving  
402 *Medicago prostrata* (Fabaceae). *Mol. Phylogenet. Evol.* 107, 367–381.  
403 <https://doi.org/10.1016/j.ympev.2016.11.020>
- 404 Escudero, M., Eaton, D.A.R., Hahn, M., Hipp, A.L., 2014. Genotyping-by-sequencing as a  
405 tool to infer phylogeny and ancestral hybridization: A case study in *Carex* (Cyperaceae).  
406 *Mol. Phylogenet. Evol.* 79, 359–367. <https://doi.org/10.1016/j.ympev.2014.06.026>
- 407 Excoffier, L., Dupanloup, I., Huerta-Sánchez, E., Sousa, V.C., Foll, M., 2013. Robust  
408 Demographic Inference from Genomic and SNP Data. *PLoS Genet.* 9, e1003905.  
409 <https://doi.org/10.1371/journal.pgen.1003905>
- 410 Fernández-Mazuecos, M., Mellers, G., Vigalondo, B., Sáez, L., Vargas, P., Glover, B.J.,  
411 2017. Resolving Recent Plant Radiations: Power and Robustness of Genotyping-by-  
412 Sequencing. *Syst. Biol.* <https://doi.org/10.1093/sysbio/syx062>
- 413 Green, R.E., Krause, J., Briggs, A.W., Maricic, T., Stenzel, U., Kircher, M., Patterson, N., Li,  
414 H., Zhai, W., Fritz, M.H.-Y.Y., Hansen, N.F., Durand, E.Y., Malaspinas, A.S., Jensen,  
415 J.D., Marques-Bonet, T., Alkan, C., Prüfer, K., Meyer, M., Burbano, H.A., Good, J.M.,  
416 Schultz, R., Aximu-Petri, A., Butthof, A., Höber, B., Höffner, B., Siegemund, M.,  
417 Weihmann, A., Nusbaum, C., Lander, E.S., Russ, C., Novod, N., Affourtit, J., Egholm,  
418 M., Verna, C., Rudan, P., Brajkovic, D., Kucan, Ž., Gušić, I., Doronichev, V.B.,  
419 Golovanova, L. V., Lalueza-Fox, C., De La Rasilla, M., Fortea, J., Rosas, A., Schmitz,  
420 R.W., Johnson, P.L.F., Eichler, E.E., Falush, D., Birney, E., Mullikin, J.C., Slatkin, M.,  
421 Nielsen, R., Kelso, J., Lachmann, M., Reich, D., Pääbo, S., 2010. A draft sequence of the  
422 neandertal genome. *Science*. 328, 710–722. <https://doi.org/10.1126/science.1188021>
- 423 Gronau, I., Hubisz, M.J., Gulko, B., Danko, C.G., Siepel, A., 2011. Bayesian inference of  
424 ancient human demography from individual genome sequences. *Nat. Genet.* 43, 1031–  
425 1035. <https://doi.org/10.1038/ng.937>
- 426 Gutenkunst, R.N., Hernandez, R.D., Williamson, S.H., Bustamante, C.D., 2009. Inferring the  
427 joint demographic history of multiple populations from multidimensional SNP frequency  
428 data. *PLoS Genet.* 5, e1000695. <https://doi.org/10.1371/journal.pgen.1000695>
- 429 Harvey, M.G., Smith, B.T., Glenn, T.C., Faircloth, B.C., Brumfield, R.T., 2016. Sequence

- 430 Capture versus Restriction Site Associated DNA Sequencing for Shallow Systematics.  
431 *Syst. Biol.* 65, 910–924. <https://doi.org/10.1093/sysbio/syw036>
- 432 Hey, J., 2010. Isolation with Migration Models for More Than Two Populations. *Mol. Biol.*  
433 *Evol.* 27, 905–920. <https://doi.org/10.1093/molbev/msp296>
- 434 Hipp, A.L., Eaton, D.A.R., Cavender-Bares, J., Fitzek, E., Nipper, R., Manos, P.S., 2014. A  
435 framework phylogeny of the American oak clade based on sequenced RAD data. *PLoS*  
436 *One* 9, e93975. <https://doi.org/10.1371/journal.pone.0093975>
- 437 Hohenlohe, P.A., Amish, S.J., Catchen, J.M., Allendorf, F.W., Luikart, G., 2011. Next-  
438 generation RAD sequencing identifies thousands of SNPs for assessing hybridization  
439 between rainbow and westslope cutthroat trout. *Mol. Ecol. Resour.* 11, 117–122.  
440 <https://doi.org/10.1111/j.1755-0998.2010.02967.x>
- 441 Huson, D.H., Rupp, R., Scornavacca, C., 2010. Phylogenetic networks: concepts, algorithms  
442 and applications. Cambridge University Press.
- 443 Leaché, A.D., Chavez, A.S., Jones, L.N., Grummer, J.A., Gottscho, A.D., Linkem, C.W.,  
444 2015. Phylogenomics of phrynosomatid lizards: Conflicting signals from sequence  
445 capture versus restriction site associated DNA sequencing. *Genome Biol. Evol.* 7, 706–  
446 719. <https://doi.org/10.1093/gbe/evv026>
- 447 Lesins, K.A., Lesins, I., 1979. *Genus Medicago (Leguminosae)*, Dr. W. Junk Publishers, The  
448 Hague. Springer Netherlands, Dordrecht. <https://doi.org/10.1007/978-94-009-9634-2>
- 449 Linder, C.R., Rieseberg, L.H., 2004. Reconstructing patterns of reticulate evolution in plants.  
450 *Am. J. Bot.* 91, 1700–1708. <https://doi.org/10.3732/ajb.91.10.1700>
- 451 Maureira-Butler, I.J., Pfeil, B.E., Muangprom, A., Osborn, T.C., Doyle, J.J., 2008. The  
452 reticulate history of *Medicago* (Fabaceae). *Syst. Biol.* 57, 466–482.  
453 <https://doi.org/10.1080/10635150802172168>
- 454 Meyer, M., Kircher, M., Gansauge, M.-T., Li, H., Racimo, F., Mallick, S., Schraiber, J.G.,  
455 Jay, F., Prüfer, K., de Filippo, C., Sudmant, P.H., Alkan, C., Fu, Q., Do, R., Rohland, N.,  
456 Tandon, A., Siebauer, M., Green, R.E., Bryc, K., Briggs, A.W., Stenzel, U., Dabney, J.,  
457 Shendure, J., Kitzman, J., Hammer, M.F., Shunkov, M. V., Derevianko, A.P., Patterson,  
458 N., Andrés, A.M., Eichler, E.E., Slatkin, M., Reich, D., Kelso, J., Pääbo, S., 2012. A  
459 high-coverage genome sequence from an archaic Denisovan individual. *Science* 338,  
460 222–6. <https://doi.org/10.1126/science.1224344>
- 461 Milne, I., Bayer, M., Cardle, L., Shaw, P., Stephen, G., Wright, F., Marshall, D., 2010.  
462 Tablet—next generation sequence assembly visualization. *Bioinforma. Appl. NOTE* 26,  
463 401–402. <https://doi.org/10.1093/bioinformatics/btp666>
- 464 Milne, I., Stephen, G., Bayer, M., Cock, P.J.A., Pritchard, L., Cardle, L., Shaw, P.D.,  
465 Marshall, D., 2013. Using Tablet for visual exploration of second-generation sequencing  
466 data. *Brief. Bioinform.* 14, 193–202. <https://doi.org/10.1093/bib/bbs012>
- 467 Nakhleh, L., 2010. A metric on the space of reduced phylogenetic networks. *IEEE/ACM*  
468 *Trans. Comput. Biol. Bioinforma.* 7, 218–222. <https://doi.org/10.1109/TCBB.2009.2>
- 469 Pease, J.B., Hahn, M.W., 2015. Detection and Polarization of Introgression in a Five-Taxon  
470 Phylogeny. *Syst. Biol.* 64, 651–662. <https://doi.org/10.1093/sysbio/syv023>
- 471 Pickrell, J.K., Pritchard, J.K., 2012. Inference of population splits and mixtures from genome-  
472 wide allele frequency data. *PLoS Genet.* 8, e1002967.  
473 <https://doi.org/10.1371/journal.pgen.1002967>
- 474 Poland, J.A., Brown, P.J., Sorrells, M.E., Jannink, J.-L., 2012. Development of high-density  
475 genetic maps for barley and wheat using a novel two-enzyme genotyping-by-sequencing  
476 approach. *PLoS One* 7, e32253.
- 477 Rheindt, F.E., Fujita, M.K., Wilton, P.R., Edwards, S. V., 2014. Introgression and phenotypic  
478 assimilation in zimmerius flycatchers (Tyrannidae): Population genetic and phylogenetic  
479 inferences from genome-wide SNPs. *Syst. Biol.* 63, 134–152.

- 480 <https://doi.org/10.1093/sysbio/syt070>
- 481 Rubin, B.E.R., Ree, R.H., Moreau, C.S., 2012. Inferring phylogenies from RAD sequence  
482 data. *PLoS One* 7, 1–12. <https://doi.org/10.1371/journal.pone.0033394>
- 483 Seehausen, O., 2004. Hybridization and adaptive radiation. *Trends Ecol. Evol.* 19, 198–207.  
484 <https://doi.org/10.1016/j.tree.2004.01.003>
- 485 Sethuraman, A., Hey, J., 2016. IMA2p - parallel MCMC and inference of ancient demography  
486 under the Isolation with migration (IM) model. *Mol. Ecol. Resour.* 16, 206–215.  
487 <https://doi.org/10.1111/1755-0998.12437>
- 488 Shafer, A.B.A., Peart, C.R., Tusso, S., Maayan, I., Brelsford, A., Wheat, C.W., Wolf, J.B.W.,  
489 2017. Bioinformatic processing of RAD-seq data dramatically impacts downstream  
490 population genetic inference. *Methods Ecol. Evol.* 8, 907–917.  
491 <https://doi.org/10.1111/2041-210X.12700>
- 492 Small, E., 2011. *Alfalfa and Relatives: Evolution and Classification of Medicago*. NRC  
493 Research Press, Ottawa, Ontario, Canada. <https://doi.org/doi:10.1139/9780660199795>
- 494 Solís-Lemus, C., Ané, C., 2016. Inferring Phylogenetic Networks with Maximum  
495 Pseudolikelihood under Incomplete Lineage Sorting. *PLoS Genet.* 12, e1005896.  
496 <https://doi.org/10.1371/journal.pgen.1005896>
- 497 Sousa, F., 2015. *Next-generation Molecular Systematics and Evolution: Insights into*  
498 *Medicago*. University of Gothenburg, Gothenburg (Sweden).
- 499 Sousa, F., Bertrand, Y.J.K., Doyle, J.J., Oxelman, B., Pfeil, B.E., 2017. Using Genomic  
500 Location and Coalescent Simulation to Investigate Gene Tree Discordance in *Medicago*  
501 *L.* *Syst. Biol.* 66, 934–949. <https://doi.org/10.1093/sysbio/syx035>
- 502 Sousa, F., Bertrand, Y.J.K., Nylinder, S., Oxelman, B., Eriksson, J.S., Pfeil, B.E., 2014.  
503 Phylogenetic Properties of 50 Nuclear Loci in *Medicago* (Leguminosae) Generated  
504 Using Multiplexed Sequence Capture and Next-Generation Sequencing. *PLoS One* 9,  
505 e109704. <https://doi.org/10.1371/journal.pone.0109704>
- 506 Stamatakis, A., 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses  
507 with thousands of taxa and mixed models. *Bioinformatics* 22, 2688–2690.  
508 <https://doi.org/10.1093/bioinformatics/btl446>
- 509 Steele, K.P., Ickert-Bond, S.M., Zarre, S., Wojciechowski, M.F., 2010. Phylogeny and  
510 character evolution in *Medicago* (Leguminosae): Evidence from analyses of plastid  
511 *trnK/matK* and nuclear *GA3ox1* sequences. *Am. J. Bot.* 97, 1142–1155.  
512 <https://doi.org/10.3732/ajb.1000009>
- 513 Takahashi, T., Nagata, N., Sota, T., 2014. Application of RAD-based phylogenetics to  
514 complex relationships among variously related taxa in a species flock. *Mol. Phylogenet.*  
515 *Evol.* 80, 77–81. <https://doi.org/10.1016/j.ympev.2014.07.016>
- 516 Than, C., Ruths, D., Nakhleh, L., *Bioinformatics*, B., Than, C., Ruths, D., Nakhleh, L., 2008.  
517 *PhyloNet*: a software package for analyzing and reconstructing reticulate evolutionary  
518 relationships. *BMC Bioinformatics* 9, 322. <https://doi.org/10.1186/1471-2105-9-322>
- 519 Twyford, A.D., Ennos, R.A., 2012. Next-generation hybridization and introgression. *Heredity*  
520 (Edinb). <https://doi.org/10.1038/hdy.2011.68>
- 521 Urban, I., 1873. *Prodomus einer Monographie der Gattung Medicago L.* Rudolph Gaertner.
- 522 Wagner, C.E., Keller, I., Wittwer, S., Selz, O.M., Mwaiko, S., Greuter, L., Sivasundar, A.,  
523 Seehausen, O., 2013. Genome-wide RAD sequence data provide unprecedented  
524 resolution of species boundaries and relationships in the Lake Victoria cichlid adaptive  
525 radiation, in: *Molecular Ecology*. Wiley/Blackwell (10.1111), pp. 787–798.  
526 <https://doi.org/10.1111/mec.12023>
- 527 Wang, S., Meyer, E., McKay, J.K., Matz, M. V., 2012. 2b-RAD: a simple and flexible method  
528 for genome-wide genotyping. *Nat. Methods* 9, 808–810.
- 529 Wen, D., Nakhleh, L., 2016. Co-estimating Reticulate Phylogenies and Gene Trees from

- 530 Multi-locus Sequence Data. *bioRxiv* 26, 1–13. <https://doi.org/10.1101/095539>
- 531 Wen, D., Yu, Y., Nakhleh, L., 2016. Bayesian Inference of Reticulate Phylogenies under the  
532 Multispecies Network Coalescent. *PLoS Genet.* 12, 1–17.  
533 <https://doi.org/10.1371/journal.pgen.1006006>
- 534 Wen, D., Yu, Y., Zhu, J., Nakhleh, L., 2017. Inferring Phylogenetic Networks Using  
535 PhyloNet. *Syst. Biol.* 00, 197–204. <https://doi.org/10.1093/sysbio/syy015>
- 536 Yoder, J.B., Briskine, R., Mudge, J., Farmer, A., Paape, T., Steele, K., Weiblen, G.D., Bharti,  
537 A.K., Zhou, P., May, G.D., Young, N.D., Tiffin, P., 2013. Phylogenetic signal variation  
538 in the genomes of medicago (Fabaceae). *Syst. Biol.* 62, 424–438.  
539 <https://doi.org/10.1093/sysbio/syt009>
- 540 Young, N.D., Debelle, F., Oldroyd, G.E.D., Geurts, R., Cannon, S.B., Udvardi, M.K.,  
541 Benedito, V.A., Mayer, K.F.X., Gouzy, J., Schoof, H., Van de Peer, Y., Proost, S., Cook,  
542 D.R., Meyers, B.C., Spannagl, M., Cheung, F., De Mita, S., Krishnakumar, V.,  
543 Gundlach, H., Zhou, S., Mudge, J., Bharti, A.K., Murray, J.D., Naoumkina, M.A.,  
544 Rosen, B., Silverstein, K.A.T., Tang, H., Rombauts, S., Zhao, P.X., Zhou, P., Barbe, V.,  
545 Bardou, P., Bechner, M., Bellec, A., Berger, A., Bergès, H., Bidwell, S., Bisseling, T.,  
546 Choisine, N., Couloux, A., Denny, R., Deshpande, S., Dai, X., Doyle, J.J., Dudez, A.-M.,  
547 Farmer, A.D., Fouteau, S., Franken, C., Gibelin, C., Gish, J., Goldstein, S., González,  
548 A.J., Green, P.J., Hallab, A., Hartog, M., Hua, A., Humphray, S.J., Jeong, D.-H., Jing,  
549 Y., Jöcker, A., Kenton, S.M., Kim, D.-J., Klee, K., Lai, H., Lang, C., Lin, S., Macmil,  
550 S.L., Magdelenat, G., Matthews, L., Mccorrison, J., Monaghan, E.L., Mun, J.-H., Najjar,  
551 F.Z., Nicholson, C., Noirot, C., O’Bleness, M., Paule, C.R., Poulain, J., Prion, F., Qin,  
552 B., Qu, C., Retzel, E.F., Riddle, C., Sallet, E., Samain, S., Samson, N., Sanders, I.,  
553 Saurat, O., Scarpelli, C., Schiex, T., Segurens, B., Severin, A.J., Sherrier, D.J., Shi, R.,  
554 Sims, S., Singer, S.R., Sinharoy, S., Sterck, L., Viollet, A., Wang, B.-B., Wang, K.,  
555 Wang, M., Wang, X., Warfsmann, J., Weissenbach, J., White, D.D., White, J.D., Wiley,  
556 G.B., Wincker, P., Xing, Y., Yang, L., Yao, Z., Ying, F., Zhai, J., Zhou, L., Zuber, A.,  
557 Dénarié, J., Dixon, R.A., May, G.D., Schwartz, D.C., Rogers, J., Quétier, F., Town,  
558 C.D., Roe, B.A., O’bleness, M., Paule, C.R., Poulain, J., Prion, F., Qin, B., Qu, C.,  
559 Retzel, E.F., Riddle, C., Sallet, E., Samain, S., Samson, N., Sanders, I., Saurat, O.,  
560 Scarpelli, C., Schiex, T., Segurens, B., Severin, A.J., Sherrier, D.J., Shi, R., Sims, S.,  
561 Singer, S.R., Sinharoy, S., Sterck, L., Viollet, A., Wang, B.-B., Wang, K., Wang, M.,  
562 Wang, X., Warfsmann, J., Weissenbach, J., White, D.D., White, J.D., Wiley, G.B.,  
563 Wincker, P., Xing, Y., Yang, L., Yao, Z., Ying, F., Zhai, J., Zhou, L., Zuber, A., Dénarié,  
564 J., Dixon, R.A., May, G.D., Schwartz, D.C., Rogers, J., Quétier, F., Town, C.D., Bruce,  
565 &, 2011. The Medicago genome provides insight into the evolution of rhizobial  
566 symbioses. *Nature* 480, 520–524. <https://doi.org/10.1038/nature10625>
- 567 Yu, Y., Barnett, R.M., Nakhleh, L., 2013. Parsimonious Inference of Hybridization in the  
568 Presence of Incomplete Lineage Sorting. *Syst. Biol.* 62, 738–751.  
569 <https://doi.org/10.1093/sysbio/syt037>
- 570 Yu, Y., Cuong, T., Degnan, J.H., Nakhleh, L., 2011. Coalescent Histories on Phylogenetic  
571 Networks and Detection of Hybridization Despite Incomplete Lineage Sorting. *Syst.*  
572 *Biol.* 60, 138–149. <https://doi.org/10.1093/sysbio/syq084>
- 573 Yu, Y., Dong, J., Liu, K.J., Nakhleh, L., 2014. Maximum likelihood inference of reticulate  
574 evolutionary histories. *Proc. Natl. Acad. Sci. U. S. A.* 111, 16448–16453.  
575 <https://doi.org/10.1073/pnas.1407950111>
- 576 Yu, Y., Nakhleh, L., 2015. A maximum pseudo-likelihood approach for phylogenetic  
577 networks. *BMC Genomics* 16, S10. <https://doi.org/10.1186/1471-2164-16-S10-S10>
- 578 Zhang, C., Ogilvie, H.A., Drummond, A.J., Stadler, T., 2018. Bayesian inference of species  
579 networks from multilocus sequence data. *Mol. Biol. Evol.* 35, 504–517.

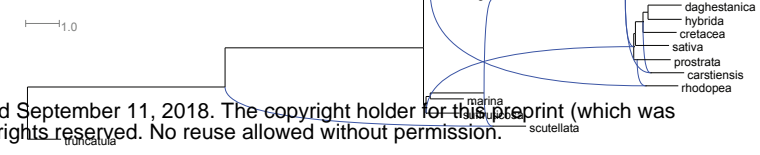


580 <https://doi.org/10.1093/molbev/msx307>  
581 Zhu, J., Wen, D., Yu, Y., Meudt, H.M., Nakhleh, L., 2017. Bayesian Inference Of  
582 Phylogenetic Networks From Bi-allelic Genetic Markers. PLoS Comput Biol 14,  
583 e1005932. <https://doi.org/10.1371/journal.pcbi.1005932>  
584

clust75.min4.PIS4



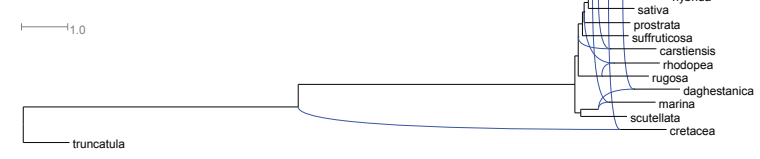
clust75.min4.PIS10



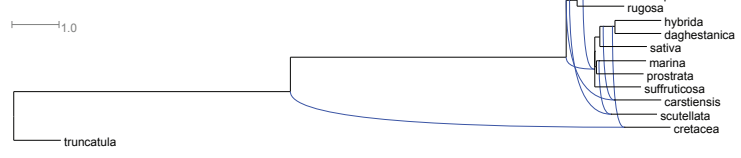
clust75.min10.PIS4



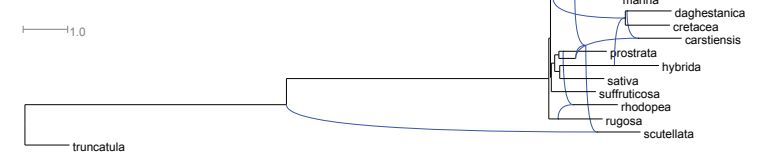
clust75.min10.PIS10



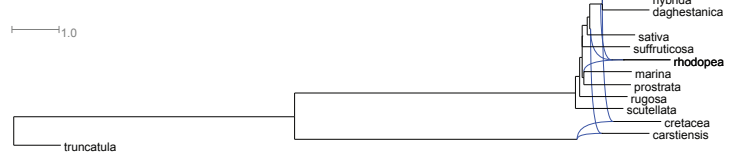
clust80.min4.PIS4



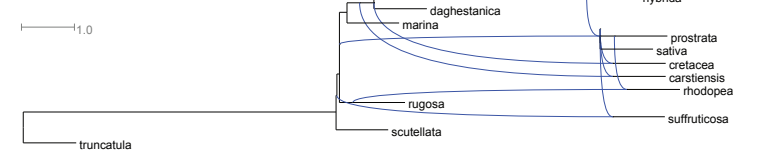
clust80.min4.PIS10



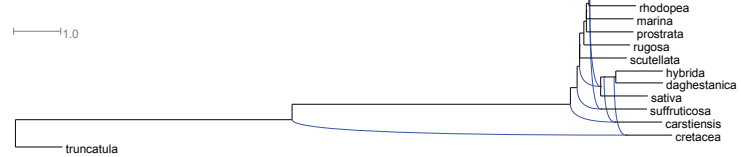
clust80.min10.PIS4



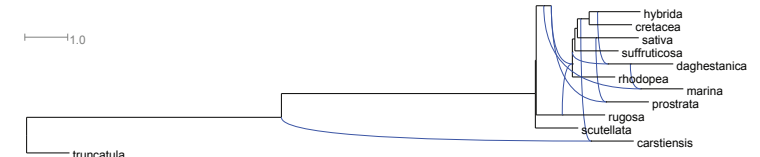
clust80.min10.PIS10



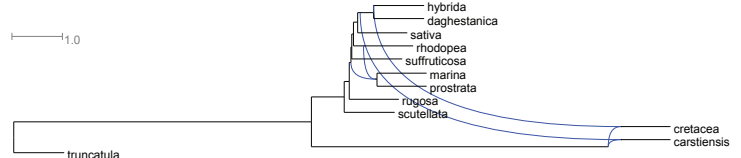
clust85.min4.PIS4



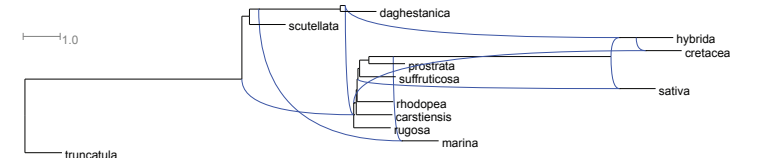
clust85.min4.PIS10



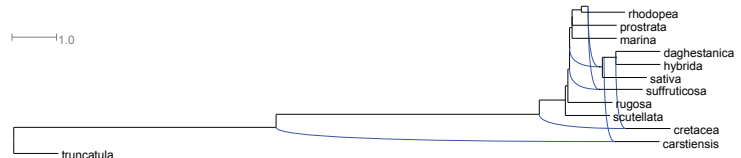
clust85.min10.PIS4



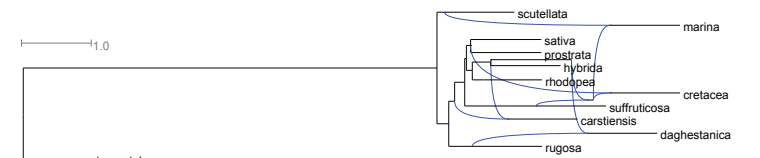
clust85.min10.PIS10



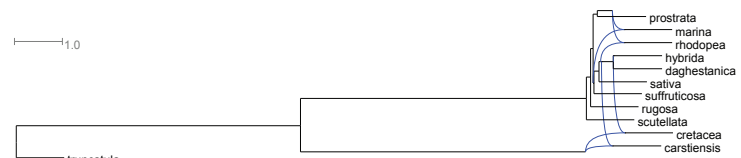
clust90.min4.PIS4



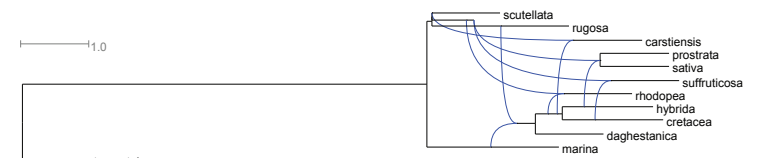
clust90.min4.PIS10



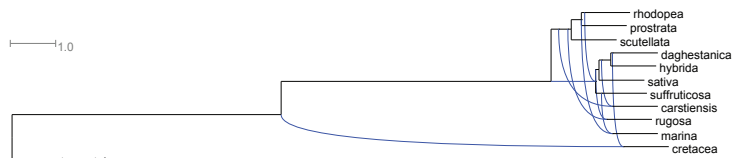
clust90.min10.PIS4



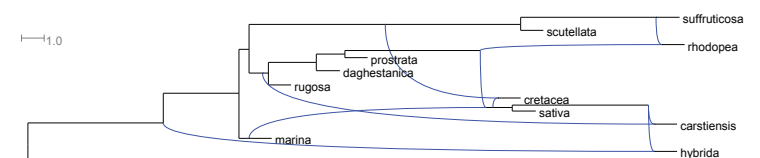
clust90.min10.PIS10



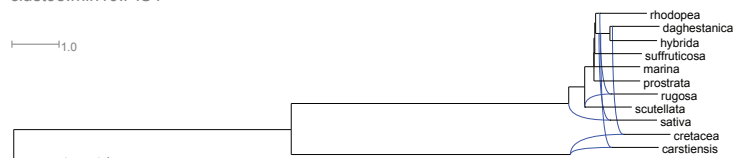
clust95.min4.PIS4



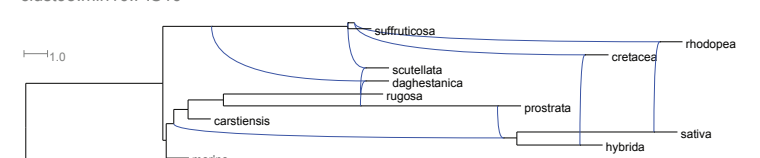
clust95.min4.PIS10

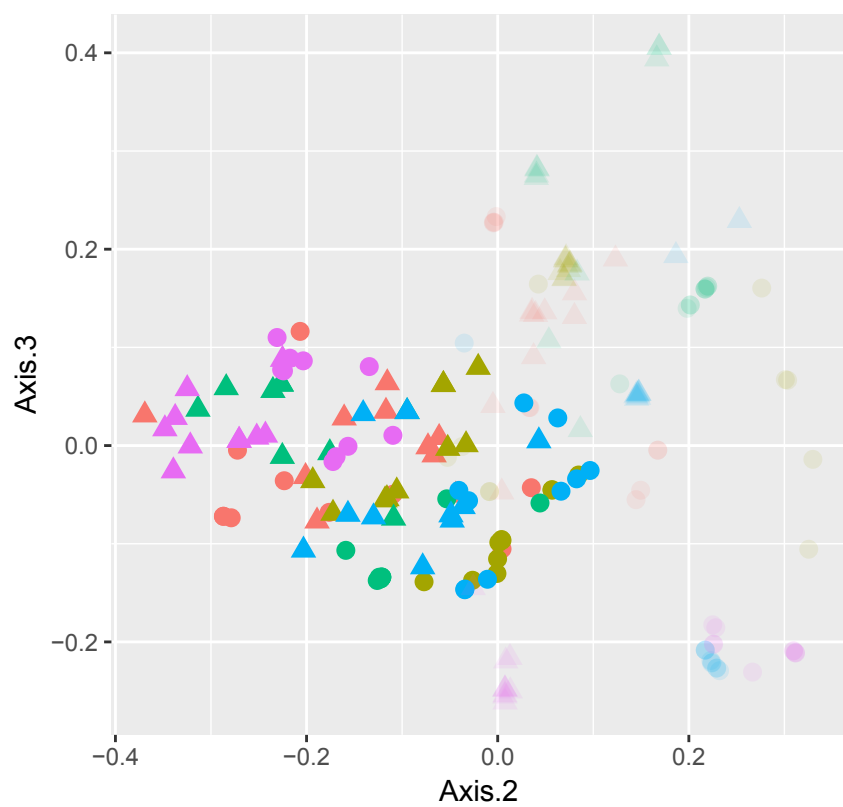
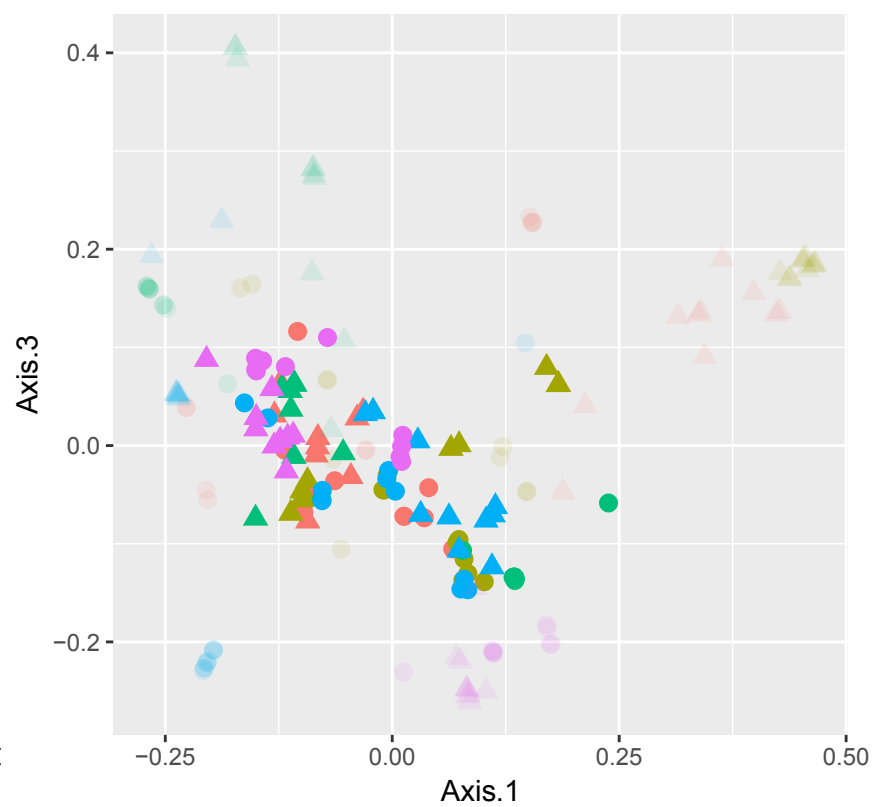
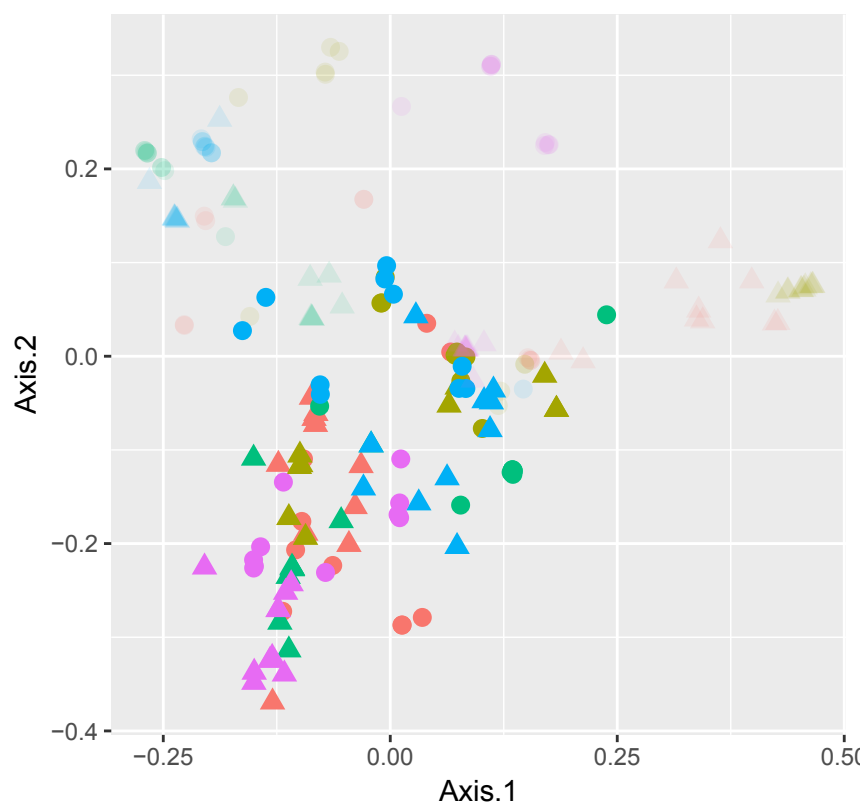


clust95.min10.PIS4



clust95.min10.PIS10





#### Parsimony informative sites



#### Min. samples locus



#### Clustering threshold

