

1 **Excessive and asymmetrical removal of heterozygous sites**
2 **by *maxSH* biases downstream population genetic inference:**
3 **Implications for hybridization between two primroses**

4 Running title: Exploring biases from bioinformatics processes

5 **Jie Zhang^{1, 2, 3}, Francisco Pina-Martins⁴, Zu-Shi Jin⁵, Yong-**
6 **Peng Cha^{1,3}, Zu-Yao Liu⁶, Jun-Chu Cha^{1,3} Jian-Li Zhao^{1, 3},**
7 **Qing-Jun Li^{1,3}**

8

9 1. Laboratory of Ecology and Evolutionary Biology, School of Ecology and Environmental
10 Sciences, Yunnan University, Kunming, 650500, China

11 2. Institute of International Rivers and Eco-security, Yunnan University, Kunming, 650500, China

12 3. Yunnan Key Laboratory of Plant Reproductive Adaptation and Evolutionary Biology, Yunnan
13 University, Kunming, 650500, China

14 4. Computational Biology and Population Genomics Group, Departamento de Biologia Animal,
15 Faculdade de Ciências, Centre for Ecology, Evolution and Environmental Changes, Universidade
16 de Lisboa, 1649-019 Lisboa, Portugal

17 5. School of Food Science, Tibet Agriculture and Animal Husbandry University, Nyingchi,
18 860000, China

19 6. Division of Evolutionary Ecology, Institute of Ecology and Evolution, University of Bern, 3012
20 Bern, Switzerland.

21 **Abstract**

22 Techniques of reduced-representation sequencing (RRS) have revolutionized ecological and
23 evolutionary genomics studies. Precise establishment of orthologs is a critical challenge for RRS,
24 especially when a reference genome is absent. The proportion of shared heterozygous sites across
25 samples is an alternative criterion for filtering paralogs, as divergent lineages should be less likely
26 to share heterozygosity. In the prevailing pipeline for variant calling of RRS data -
27 PYRAD/IPYRAD, *maxSH* is an often overlooked parameter with implications to detecting and
28 filtering paralogs according to shared heterozygosity. Using empirical GBS data of two primroses
29 (*Primula alpicola* Stapf and *Primula florindae* Ward) and their putative hybrids, and extra datasets
30 of Californian golden cup oaks, we explore the impact of *maxSH* on filtering paralogs and further
31 downstream analyses. Our study sheds light on the simultaneous validity and risk of filtering
32 paralogs using *maxSH*, and its significant effects on downstream analyses of outlier detection,
33 population assignment, and demographic modelling, emphasizing the importance of attention to
34 detail during bioinformatics processes. The mutual confirmation between results of population
35 assignment and demographic modelling in this study suggested *maxSH* = 0.10 has a potentially
36 excessive and asymmetrical effect on the removal of truly shared heterozygous sites as paralogs.
37 These results indicate that hybridization origin hypotheses of putative hybrids represented by
38 results with *maxSH* = 0.25 and 0.50 are more credible. In conclusion, we revealed the critical
39 hazard of paralogs filtration according to sharing heterozygosity at first, so that we propose to use
40 specific protocols, rather than *maxSH*, to filter potential paralogs for closely related lineages.

41 **Key words**

42 Heterozygous site, bioinformatics, demographic modelling, paralog, *Primula*

43

44 **1. Introduction**

45 Techniques of reduced-representation sequencing (RRS) such as restriction site-
46 associated DNA sequencing (RADseq; Baird et al., 2008) and genotyping by
47 sequencing (GBS; Davey et al., 2011) are increasingly prevalent, especially in species
48 lacking reference genomes, having revolutionized ecological and evolutionary
49 genomics studies for their effective generating of genome-wide molecular markers
50 with a desirable cost (Rodríguez-Ezpeleta et al., 2016; McKinney et al., 2017). Such
51 markers allow inferring not only basic statistics of population genetics, but also
52 phylogenetic reconstruction (e.g. Escudero et al., (2014)), population clustering (e.g.
53 Ortego et al., (2017)), local adaptation (e.g. Pina-Martins et al., (2019)), or
54 demographic inference (e.g. Excoffier et al., (2013)). Nevertheless, RRS has obvious
55 shortcomings, such as precise establishment of homologous loci; a key challenge for
56 all sequencing techniques, which is particularly difficult in RRS, since sequence
57 similarity is often a crucial criterion for assembling demultiplexed reads into
58 orthologous loci (Ilut et al., 2014; Harvey et al., 2015; McCartney-Melstad et al.,
59 2019).

60 Several bioinformatics pipelines and programs have been developed for RRS loci
61 assembly and data analyses, designed with various algorithms, logic and computer
62 languages, such as STACKS (Catchen et al., 2013), IPYRAD (Eaton & Overcast, 2020),

63 dDOCENT (Puritz et al., 2014), RAINBOW (Chong et al., 2012) and so on. While
64 considerable attention has been paid to reveal and reduce biases associated with wet
65 laboratory, sequencing process and species-specific genome properties, biases from
66 downstream bioinformatics analyses are less concentrated (but see Illut, et al., 2014;
67 Mastretta-Yanes et al., 2015; Shafer et al., 2017; O’Leary et al., 2018).

68 PYRAD/IPYRAD is one of the most widely-used programs for quality filtering and
69 variant calling of RRS data with or without a reference genome, and is superior in
70 handling paired-end sequencing data and INDEL variation (Eaton, 2014; Eaton &
71 Overcast, 2020). In PYRAD/IPYRAD, “clustering threshold” is a key parameter with
72 respect to establishing homology, with a default value 85%. Several articles have
73 soundly evaluated its influence and proposed advice on choosing its optimal value
74 (Illut et al., 2014; McCartney-Melstad et al., 2019). But another important parameter,
75 perhaps severely overlooked in most studies, is the maximum number (or proportion)
76 of shared heterozygous sites in a locus. It can be set as both integer and decimal,
77 abbreviated as *maxSH* or *maxsharedH* in PYRAD and *max_shared_Hs_locus* in IPYRAD
78 (Eaton, 2014, Eaton & Overcast, 2020, in order to avoid confusion, we unified its
79 name as *maxSH* hereafter). This parameter allows identification of paralogs according
80 to the extent of shared heterozygosity across samples. It is worth mentioning that, in
81 PYRAD the default value of *maxSH* is 4, which means heterozygous site shared by
82 more than four individuals will be removed. While in IPYRAD the default value of
83 *maxSH* is 0.50 (50%), allowing half of tested samples to share any given
84 heterozygous site. Given bi-allelic SNPs were only generated by mutation, sharing

85 heterozygosity among species could only be originated by incomplete lineage sorting
86 (ILS), gene flow and sharing gene/genome duplication (Spofford, 1969; Innan and
87 Kondrashov, 2010; Fijarczyk and Babik, 2015). On this basis, highly divergent
88 species (lineages) should be less likely to share heterozygous sites regarding their
89 respective evolutionary history, and heterozygous sites simultaneously presented in
90 many samples are more likely to represent clustering of paralogs with a fixed
91 difference rather than a true heterozygous site (Eaton et al. 2015; Eaton & Overcast,
92 2020). PYRAD was primitively designed for phylogenetic reconstruction, thus extreme
93 low default value are benefit and reasonable. However, numerous studies have used
94 this default value or some extreme low values to generate datasets for population
95 genomics study on close related lineages (e.g. Cavender-Bares et al., 2015; Eaton et
96 al. 2015; Ortego et al., 2017; Tonzo et al., 2020). Closely related lineages possess
97 relatively high proportion of heterozygous sites (Sota and Vogler 2003, Lischer et al.,
98 2014). Besides, the high probability of introgression, hybridization, and sharing
99 ancestral polymorphism can potentially contribute to extensively sharing
100 heterozygosity among close lineages. Therefore, truly shared heterozygous sites will
101 be simultaneously filtered if we set a lower threshold of *maxSH* in population
102 genomics studies. In this regard, tuning *maxSH* is essential for unbiased inference,
103 and its influence on downstream analyses is to be expected. However, it remains
104 undiscovered to what extent *maxSH* affects such downstream analyses and further
105 biological inference in population genomics frameworks.

106 We mainly performed our empirical test on the influence of *maxSH* values using
107 two distylous primroses, *Primula alpicola* Stapf and *P. florindae* Ward co-occurring in
108 the Shergyla mountains, southeast of Qinghai-Tibet region. *Primula L.* (Primulaceae)
109 is arguably a high-profile genus for heterostyly since Darwin's seminal book (Darwin,
110 1877), exhibiting extreme species richness and diversity in the eastern Sino-
111 Himalayan region (Richards, 2003). Although frequent hybridization and
112 introgression have been considered critical factors for the complex phylogenetic
113 relationships within *Primula* (Richards, 2003; Ren et al., 2018), only a few natural
114 hybridization events were well documented in the complex's distribution center (see
115 Zhu et al., 2009; Ma et al., 2014; Xie et al., 2017). Natural hybridization between *P.*
116 *alpicola* and *P. florindae* was reported according to field observations by Ward, and
117 artificial crossing is compatible in the garden (Richards, 2003). These two species are
118 somewhat difficult to distinguish in sympatry due to morphological similarity, but still
119 can be identified according to some critical differences on leaves and flowers.
120 Moreover, their putative hybrids have been identified in middle elevation area of the
121 Shergyla mountains, where *P. alpicola* and *P. florindae* share most pollinators with
122 relative long flowering phenology overlapping (personal observation).

123 In this work, we sequence two primrose species and their putative hybrids using
124 GBS to mainly disentangle 1) whether *maxSH* is effective on handling potential
125 paralogs; 2) How, and to what extent could downstream bioinformatics analyses be
126 influenced regarding different *maxSH* thresholds; 3) whether biological inference can
127 tolerate potential excessive removal of heterozygous sites resulting from low *maxSH*

128 thresholds. In order to further verify the influence of *maxSH* on population genomics
129 studies, we also conducted extra demographic modelling for another datasets
130 discussing the introgression between two Californian golden cup oaks.

131 **2. Materials and Methods**

132 **2.1 Sample collection**

133 In 2015 and 2018, we sampled *P. alpicola* and *P. florindae* in the putative hybrid zone
134 (elevation = 3672.71 m, latitude = 29.6704°N, longitude = 94.7157°E) from the
135 Shergyla mountains. We collected ten individuals of *P. alpicola*, nine individuals of *P.*
136 *florindae*, and fifteen putative hybrids identified by their flower color and leaf shape
137 traits from this zone. Furthermore, three *P. sikkimensis* individuals collected from
138 higher elevation of Shergyla mountains were sampled as outgroup species. Fresh leaf
139 samples were quickly dried and stored with silica gel until DNA extraction.

140 **2.2 DNA extraction and sequencing**

141 Total DNA extraction followed a modified CTAB protocol (Doyle & Doyle 1987).
142 The purity and amount of all extracted DNA was assessed using Nanodrop 1000 and
143 Agarose gel. “Genotyping by sequencing” (GBS) technique was used for genotyping
144 DNA samples and obtaining high density SNPs. In brief, DNA was double digested
145 using *MseI*+*HaeII* restriction enzymes, and ligated Illumina adapters and barcodes.
146 After libraries were constructed, Qubit 2.0 was used for preliminary quantification,
147 then DNA samples were uniformly diluted to about 1 ng/μl. At last, libraries were

148 pooled, and then paired-end sequenced following the standard protocol using Illumina
149 HiSeq PE150 platform by Novogene Bioinformatics Technology Co., Ltd., Beijing,
150 China (www.novogene.cn).

151 **2.3 Bioinformatics and data filtering**

152 We assembled *de novo* loci and called SNPs using IPYRAD v.0.7.30. We kept non-
153 target parameters of IPYRAD as default value for prescriptive quality control steps, then
154 designed parameter assemblies to test and compare their influence on downstream
155 analyses: 1) We clustered our quality-filtered reads considering three thresholds: 85%,
156 90% and 95%; 2) As *maxSH* was the main-tested parameter, we set it to two
157 frequently adopted values: 0.10 and 0.50 plus an intermediate value of 0.25. Complete
158 information for each parameter assembly is available in supplemental files as datafile
159 S1. In order quantify filtered paralogous loci, we used the output file “stats.txt” from
160 IPYRAD step 7 following McCartney-Melsted et al. (2019) to plot the percentage of
161 flagged paralogous
162 $([\text{filtered_by_max_indels} + \text{filtered_by_max_snps} + \text{filtered_by_max_shared_het} + \text{filtered_by_max_alleles}] / \text{total_prefiltered_loci})$. We further filtered processed data for
163 keeping only biallelic SNPs, requiring a minimum allele frequency (MAF) > 0.03 ,
164 and setting missing data rate (proportion of samples does not contain data at a given
165 SNP) as 60% using VCFTOOLS v.0.1.14 (Danecek et al., 2011) since the number of
166 retained SNPs sharply decreased with further missing data constrains. In the final
167 step, we kept only the center SNP per locus so that we can minimize the effect of
168 linkage disequilibrium using a python script “*vcf_parser.py*”

170 (https://github.com/CoBiG2/RAD_Tools/blob/master/vcf_parser.py) as of commit
171 “0893296”. In total, 9 datasets were created representing assemblies resulting from all
172 parameter combinations. Each dataset was entitled as the combination of the initial of
173 tested parameter (clustering threshold and *maxSH*) and their representative value, like
174 *c85m10*.

175 **2.4 Outlier detection**

176 Outlier SNP detection was performed to obtain unbiased population structure and
177 further demographic modelling. We detected outlier loci using two programs:
178 BAYESCAN v.2.1 (Foll & Gaggiotti, 2008) based on Bayesian approach and R package
179 *pcadapt* v.4.3.2 (Luu et al., 2017) based on principal component analysis (PCA). For
180 both programs, individuals were preliminarily grouped according to sampling
181 categories. BAYESCAN was run using 20 pilot runs of length 5,000, a burn-in length of
182 50,000, a main output iterations of 10,000, a thinning interval of 10, and a detecting
183 threshold of 0.05. For running *pcadapt*, we firstly evaluated the number of principal
184 components using “score plot” wrapped in *pcadapt* due to the poor resolution of
185 Cattell’s graphic rule in our case, a list of candidate SNPs under an expected FDR
186 $\alpha = 0.05$ were identified as outliers following a standard *pcadapt* workflow. Outliers
187 identified by either program were excluded from population structure inference
188 analyses.

189 **2.5 Population structure and hybrids identification**

190 We then inferred population genetic structure and preliminarily identified hybrids for
191 each data set using STRUCTURE v.2.3.4 (Pritchard et al., 2000) wrapped in the program
192 *Structure_threader* v.1.3.4 (Pina-Martins et al., 2017). This program is characterized
193 by parallelizing multiple runs of genetic clustering software and automatically
194 assessing the best K as well as drawing the “meanQ” plots. We used filtered SNPs
195 with 20 independent runs for K values from 1 to 6 to estimate the optimal number of
196 clusters with a burn-in of 100 000, followed by 200 000 Markov chain Monte Carlo
197 (MCMC) repetitions. The best K was estimated according to the widely used ΔK
198 method (Evanno et al., 2005) implemented in *Structure_threader*. Principal
199 component analysis was also performed using an R script “*snp_pca_static.R*”
200 (https://github.com/CoBiG2/RAD_Tools/blob/master/snp_pca_static.R) as of commit
201 “bb2fc45”, in order to improve presentation, we slightly tweaked this script for our
202 case regarding colours.

203 **2.6 Demographic modelling**

204 We used FASTSIMCOAL2 v.2.6.0.2 (Excoffier et al., 2013) for comparing different
205 demographic models via coalescent simulations. Because demographic modelling
206 depends on well-defined population structure, three populations (demes) defined
207 based on STRUCTURE results were prepared for modelling. The folded joint SFS and
208 unbiased estimation of allele frequency were performed using the Python script
209 *easySFS.py* (<https://github.com/isaacovercast/easySFS>) as of commit “aaf80ea”,
210 which can effectively downsample populations for generating input “.obs” files for
211 FASTSIMCOAL2. Because no invariable loci were involved in our SFS, we enabled

212 demographic modelling by introducing a calculated effective population size of *P.*
213 *alpicola* (Papadopoulou & Knowles, 2015; Ortego et al., 2017). As $Ne = \pi/4\mu$, we
214 inferred average mutation rate per site per generation μ from *Arabidopsis thaliana*
215 (Nordborg et al., 2005) following Gossmann et al. (2012) as it is the genetically
216 closest species with known μ value. π value was computed in DNASP v.6.12.03 (Rozas,
217 et al., 2017) using “.allele.loci” file containing both polymorphic and non-
218 polymorphic loci generated by IPYRAD. Average generation time was set as 1 yr
219 (personal observation). Finally, we roughly bounded upper limit of divergence time
220 between *P. alpicola* and *P. florindae* as 25 Ma (million years ago) according to
221 estimated time when *Primula* diverged from *Soldanella* (de Vos et al., 2014).

222 To disentangle how putative hybrids speciated, three models were designed,
223 respectively describing putative hybrids diverged from *P. florindae*, putative hybrids
224 diverged from *P. alpicola*, and putative hybrids originated from hybridization between
225 *P. alpicola* and *P. florindae*, considering post-divergence asymmetric gene flow (Fig.
226 S1). Meanwhile, three alternative models describing similar divergence scenarios but
227 without gene flow were also prepared as comparison (Fig. S1). Each model was run
228 100 independent replicates following 250,000 simulations with 60 expectation-
229 conditional maximization (ECM) cycles, a stop criterion of 0.001, and zero SFS
230 removed using FASTSIMCOAL2. Akaike’s information criterion (AIC) was used to select
231 the best model. The replicate with the maximum estimated likelihood of each model
232 was selected for AIC and Δ AIC calculation.

233 Because Ortego et al. (2017) has conducted the comparison of genetic clustering
234 results on different parameter assemblies of both stacks and IPYRAD, we mainly
235 performed the comparison of coalescent analyses results among the STACKS datasets
236 author used for his downstream analyses and two available PYRAD datasets with
237 $maxSH = 0.1$ available as supplementary material. We named two PYRAD datasets as
238 oak_c85 and oak_c90 respectively according to their clustering threshold, and the
239 STACKS datasets as oak_stacks. As demographic modelling has been done for STACKS
240 datasets (Ortego et al., 2017, Table1, Table2, Fig 2). For two PYRAD datasets, we
241 followed Ortego's filtering steps to extract the unlinked neutral bi-allelic SNPs, then
242 we used *easySFS.py* for downsampling and generating input files for FASTSIMCOAL2
243 containing folded joint SFS information, respectively. The alternative models and
244 execution of FASTSIMCOAL2 were also in line with Ortego et al. (2017). At last, we
245 compared the difference of best model and corresponding parameter estimates among
246 these three datasets.

247 **3. Results**

248 **3.1 Sequencing output and variation in data processing**

249 The number of usable paired-end sequence reads ranged from 3,108,740 to 7,799,062
250 with an average of 5,291,571 per sample (SD=1,063,312). Total loci assembled by
251 IPYRAD increased with both clustering threshold and $maxSH$ values, ranging from
252 33,328 to 88,630 (Table 1). Total SNPs called primarily by IPYRAD varied similarly to
253 total loci, from 226,965 to 605,829, with an average of 422,815 (Table 1). The

254 number of filtered SNPs also varied similarly, ranging from 5276 to 24166 (Table 1),
255 but its differences among different *maxSH* values are further enlarged despite
256 changing clustering thresholds.

257 Changing either clustering threshold or *maxSH* values effectively alters the
258 proportion of flagged paralogs, suggesting their remarked association to detect and
259 filter potential paralogs (Fig. 1). The percentage of flagged paralogs steeply decreased
260 (~3% to ~5%) when changing *maxSH* from 0.10 to 0.25, and then tend to be
261 approximate (< 0.8%) between *maxSH* = 0.25 and 0.50. Changing clustering threshold
262 values from 0.85 to 0.95 resulted in a stepwise reduction of flagged paralogs. In
263 dataset *c85m10*, more than 25% of assembled loci were identified as paralogs, yet
264 *c95m50* contained ~12% flagged paralogs (Fig. 1).

265 BAYESCAN and *pcadapt* showed a large discrepancy in detecting loci under selection:
266 the mean number of detected outliers of all datasets was about 21 for BAYESCAN and
267 1036 for *pcadapt*. In addition, the number of detected outliers increases with *maxSH*
268 value for *pcadapt* except for datasets with the highest clustering threshold, but
269 decreases for BAYESCAN in all datasets. It is also worth noting that BAYESCAN detected
270 no outliers on 4 of 6 parameter assemblies with *maxSH* = 0.25 and 0.50, and outliers
271 common to both programs were only detected when *maxSH* = 0.10 (Table 1).

272 After all filtering steps, the final number of neutral SNPs used for genetic clustering
273 and demographic modelling ranged from 5002 for *c85m10* to 22126 for *c95m50*
274 (Table 1).

275 3.2 Variation in population assignment

276 Most bayesian clustering results confirmed that samples labeled as hybrids exhibited a
277 genetic mixture of *P. alpicola* and *P. florindae*, especially for *maxSH* values of 0.25
278 and 0.50 (Fig. 2D-I). Additionally, differences between results of *maxSH* = 0.25 and
279 0.50 were relatively small when keeping clustering threshold constant. However,
280 hybrid ancestry proportion of *P. florindae* was relatively increased in datasets with the
281 lowest *maxSH* (Fig 2A-C). Particularly, five putative hybrids were genetically
282 clustered to *P. florindae* for *c95m10* (Fig. 2C). For all bayesian clusters, the optimal K
283 value is 2 (datafile S2). Interactive version of plots for all K values were all available
284 in supplemental files as datafile S3.

285 PCA results were roughly consistent with those obtained by the bayesian clustering
286 approach, especially the similarity between plots with *maxSH* = 0.25 and 0.50 (Fig.
287 3D-I). Besides, when *maxSH* > 0.10, PCA results supported *P. alpicola* and *P.*
288 *florindae* as genetically separated clusters, and samples of putative hybrids were
289 located between *P. alpicola* and *P. florindae*, indicating their genetically admixed
290 background. Unexpectedly, these datasets collectively segregated four putative
291 hybrids from other samples marked along PC1 or PC2. By contrast, part of putative
292 hybrids always showed their genetic similarity to *P. florindae* along at least one PC
293 when *maxSH* = 0.10 (Fig. 3A-C), particularly, plot of *c95m10* distinctively exhibited a
294 fusion of genetic clusters (Fig. 3C).

295 3.3 Variation in demographic modelling

296 Choice of *maxSH* showed dramatic impact on demographic modelling (Fig. 4).
297 Results of modelling were directed to two drastically different scenarios, therefore, all
298 results of parameter assemblies with *maxSH* = 0.25 and 0.50 inferred putative hybrids
299 were speciated from hybridization between *P. alpicola* and *P. florindae* accompanied
300 by interspecific gene flow (Model C1; Fig. 4D-I; Table s4-s9). While the best model
301 of three datasets with *maxSH* = 0.10 alternatively fit the scenario that putative hybrids
302 were diverged from *P. florindae* with post-divergence asymmetric gene flow (Model
303 A1; Fig. 4A-C; Table s1-3). Additionally, for datasets *c95m10*, the model indicating
304 putative hybrids were speciated from hybridization (model C1, Fig.S2) was
305 statistically equivalent to the best model to some extent ($\Delta AIC = 3.85$, Burnham &
306 Anderson, 1998).

307

308 Demographic estimations also varied dramatically between datasets, including split
309 time, effective population size, proportion of migrants, and migration rates (Fig. 4).
310 Seven of nine assemblies indicated *P. alpicola* and *P. florindae* diverged more than 20
311 Ma (Fig. 4B, D-I), while *c85m10* suggested they split about 13.5 Ma (Fig. 4A).. For
312 the speciation time of putative hybrids, three models with *maxSH* = 0.50 referring to
313 hybridization origin with gene flow all directed to near 0.18 Ma (Fig. 4G-I), while the
314 rest showed that speciation time of putative hybrids varied from 0.04 to 0.19 Ma
315 regardless of ancestral lineage. Unlike split time, effective population size inferences
316 were irregularly variable, yet all results suggested putative hybrids hold the smallest
317 effective population size. Besides, eight of nine models indicated expansion of

318 effective population size at first, then coming to the recent constriction. *P. florindae*
319 contributed a higher proportion of ancestry (from 0.62 to 0.75) to the hybrid lineage
320 than *P. alpicola* according to five of six models referring to hybridization origin with
321 gene flow (Fig. 4E-I), only *c85m25* supported both *P. florindae* and *P. alpicola*
322 contributed the same proportion (0.5) to the hybrid lineage (Fig. 4D).. Interestingly,
323 regardless of how putative hybrids originated, All models with gene flow got higher
324 AIC scores compared to models without gene flow in the same parameter assembly,
325 and eight of them shared a similar gene flow pattern: weak or moderate (only in
326 *c95m10*) continuous gene flow between *P. alpicola* and *P. florindae* ; asymmetric
327 gene flow between *P. alpicola* and putative hybrids, varying from moderate to strong;
328 moderate gene flow from putative hybrids to *P. florinade*, while the reverse was subtle
329 (Fig. 4A-E, 4G-I). Yet *c95m25* supported a different gene flow pattern: moderate gene
330 flow from *P. alpicola* to *P. florindae* was supported. Besides, contrary to previous
331 pattern, gene flow from putative hybrids to *P. florindae* was subtle, while the reverse
332 was strong (Fig. 4F). It is also worth mentioning that almost all results of parameter
333 estimation for *c95m25* were distinct from other datasets.

334 2.7 Verification from PYRAD datasets of Californian golden cup oaks

335 The best model for both PYRAD datasets with $maxSH = 0.1$ is Model B1 (Table s10,
336 s11), indicating the southern lineage of *Quercus chrysolepis* was diverged from *Q.*
337 *tomentella*. which is the second-best model for coalescent results of oak_Stacks.
338 Additionally, ModelC1 representing hybridization origin was even not the second-best
339 model for both oak_85 and oak_90. Additionally, for oak_90, modelA1 indicating

340 the southern lineage of *Q. chrysolepis* was diverged from the northern lineage of *Q.*
341 *chrysolepis* was statistically equivalent to the best model (Model B1) to some extent
342 ($\Delta AIC = 2.28$, Burnham & Anderson, 1998).

343 The advent of reduced-representation sequencing (RRS) has definitely facilitated
344 studies on ecology and evolution in depth (Twyford & Ennos, 2012; Andrews et al.,
345 2016). However, complex software and absence of standard analyses pipelines could
346 mislead the analyses process, requiring conclusions to be drawn with caution (Shafer
347 et al., 2017). Although various studies are dedicated to exploring biases from
348 bioinformatics analyses processes, more attention should be paid to each and every
349 detail, due to their potentially immeasurable influence (Gautier et al., 2013; Arnold et
350 al., 2013; Shafer et al., 2017). The comparison of different values of *maxSH* in this
351 study showed that an undesirable filtration of shared heterozygous sites can have a
352 large impact on downstream analyses in a population genomics framework, altering
353 the final biological inference. Since we carried out this research using empirical data,
354 parts of our results could merely reflect some unique characteristics of the two tested
355 datasets. However, the consistent influence of *maxSH* presented by two different
356 datasets, especially on demographic modelling, has bolstered our confidence to draw
357 a conclusion that a strict *maxSH* threshold is improper when conducting population
358 genomics analyses.

359 In this study, we filtered missing data using moderate thresholds. Changing missing
360 data rate could inevitably bring variations into downstream analyses, but for this first
361 approach, we avoided using too many variables, which may result in focus reduction.

362 In fact, although a handful of studies have discussed its influence, dealing with
363 missing data remains controversial (Huang and Knowles, 2016; Paris et al., 2017;
364 Shafer et al., 2017; Yi and Latch, 2021).

365 **4.1 Influence of *maxSH* and its interactions with clustering threshold**

366 Identification of paralogs has unendingly been a challenge we have to cope with,
367 because it can certainly act, as demonstrated here, on almost all downstream analyses,
368 for example, outlier detection (Table 1), phylogenetic reconstruction (Fitz-Gibbon et
369 al., 2017; McCartney-Melstad et al., 2019), demographic inferences (Fig. 4; Shafer et
370 al., 2017), and to some extent, population clustering (Fig. 2 & 3; Rodríguez-Ezpeleta
371 et al., 2016). For the increasingly prevalent PYRAD/IPYRAD, McCartney-Melstad et al.
372 (2019) have illustrated that clustering thresholds strongly affect paralogs filtering and
373 subsequent phylogenetic resolution in detail. In this study, we confirmed its
374 significance on filtering paralogs and additionally estimated its influence on typical
375 population genomic analyses. More importantly, we demonstrated that an
376 underestimated parameter of PYRAD/IPYRAD, *maxSH*, is as influential as clustering
377 threshold on handling potential paralogs. A low threshold of *maxSH* has remarkably
378 increased the proportion of flagged paralogs in the tested data set. What's more,
379 downstream analyses were all significantly impacted by *maxSH*, since variation of
380 population assignment and demographic modelling were far more closely associated
381 with *maxSH*, rather than clustering threshold. This could be mainly interpreted by
382 their totally different rules for filtering paralogs. For clustering threshold, paralogs are
383 identified via comparison of sequencing similarity, a lower threshold could lead to

384 underestimation of the number of loci and thus undersplitting (Rodríguez-Ezpeleta et
385 al., 2016). However, biases from undersplitting are somewhat unpredictable, as we
386 cannot know what such information represents. By contrast, a lower threshold of
387 *maxSH* can directly filter heterozygous sites across many samples, which should only
388 be originated from interspecific gene flows or sharing ancestral polymorphism if it is
389 a true heterozygous site (Fijarczyk and Babik, 2015). This could explain why
390 changing *maxSH* can strikingly alter the choice of the best model for coalescent
391 simulations, while clustering thresholds mainly impacted parameter estimates.
392 Besides, other downstream analyses should also clearly be more vulnerable once we
393 improperly filtered these informative sites as paralogs.

394 Besides the different expected behaviors of clustering threshold and *maxSH*, there
395 may also be some unexpected interactions between them. Oversplitting due to
396 extremely high clustering thresholds has been demonstrated to cause a split between
397 true allelic variants of orthologous loci into putatively separate loci (McCartney-
398 Melsted et al., 2019). On the one hand, exorbitant request of sequence similarity by
399 high clustering threshold can directly limit the proportion of heterozygosity for
400 assembled loci. On the other hand, a low *maxSH* value can further exclude loci
401 regarding dissimilarity components. Thus interactions between clustering threshold
402 and *maxSH* can lead to considerable but unaccounted decrease of genetic distance
403 among lineages, which could explain the fusion of genetic clustering intensively
404 represented by PCA result for *c95m10*. It is also strongly supported by STRUCTURE
405 results of golden cup oaks. in the STACKS datasets, part of individuals of *Q. tomentella*

406 in CAT population and most individuals of *Q.chrysolepis* in MOJ, LAG, BER, GAB,
407 FIG, and HAS populations exhibited some extent genetic admixture. While both
408 proportion of individuals exhibiting genetic admixture and the extent of admixture has
409 decreased in two PYRAD datasets, and higher clustering threshold has resulted in
410 heavier decrease (Ortego et al., 2017, Fig. S4). Extending this study to other datasets
411 would be helpful for confirming how prevalent the issue really is. Yet it has already
412 implied an optimal clustering threshold is urgent, as it is the first and great influential
413 filtering step in PYRAD/IPYRAD.

414 Handling heterozygous sites has frequently been neglected for RRS data
415 processing. On the one hand, most studies do not clearly exhibit information on
416 whether heterozygous sites were dropped or retained. On the other hand, phasing
417 between loci is almost impossible when a reference genome is unavailable (Garrick et
418 al., 2010; Lischer et al., 2014). As a consequence, heterozygous sites are improperly
419 filtered or totally excluded. However, multidimensional information of introgression
420 and incomplete lineage sorting from heterozygous positions is undoubtedly precious
421 for population genomic studies tackling closely related lineages. Our study revealed
422 the fathomless influence from processing heterozygous sites at first, and illustrated
423 filtering paralogs according to shared heterozygosity is risky and unreliable.

424 **4.2 Interpretation of divergent biological inference**

425 Clustering threshold and *maxSH* biological inference were determined by tuning said
426 parameters, with special emphasis on *maxSH*. Individuals labeled as hybrids were

427 genetically admixed in most cases according to STRUCRURE and PCA. These results
428 are in agreement with their intermediate morphology, indicating potential
429 hybridization origin of putative hybrids. Demographic models with $maxSH = 0.25$ and
430 0.50 supported a hybridization origin of putative hybrids scenario. However, a low
431 $maxSH$ value shifts this conclusion to a putative divergent origin, especially when
432 combined with high clustering threshold values. Given that patterns of gene flow
433 across eight modelling results collectively indicated gene flow from putative hybrids
434 to *P. florindae* is always several times higher than to *P. alpicola*, putative hybrids are
435 less likely to share heterozygous sites with *P. alpicola*. When we set a small $maxSH$
436 value, those limited shared heterozygous sites have to be preferentially filtered, while
437 on the other hand, shared heterozygous sites between hybrids and *P. florindae* can be
438 more likely retained. These results were also partially verified by the coalescent
439 results for golden cup oaks. The best model of STACKS datasets without filtering
440 sharing heterozygous sites indicated hybrid origin of the southern lineage of
441 *Q. chrysolepis*, while two PYRAD datasets collectively tend to the model describing
442 *Q. chrysolepis* was diverged from *Q. tomentella* with post-divergence gene flow. As
443 such, the mutual confirmation between population assignment and demographic
444 modelling illustrated that extreme $maxSH$ has brought excessive and asymmetrical
445 removal of truly sharing heterozygous sites as paralogs in this study. We thus tend to
446 infer putative hybrids to be originated by hybridization between *P. alpicola* and *P.*
447 *florindae*, with an asymmetrical pattern of gene flow.

448 Hybridization could be widespread and play vital role on the diversification of
449 *Primula* (Schimidt-Lebuhn et al., 2012; Boucher et al., 2016; Keller et al., 2021).
450 Besides, increasing evidences have proofed the existence of multiple gene/genome
451 duplication and their significance on the origin of heterostyly, the most famous feature
452 of *Primula* (Li et al., 2016; Huu et al., 2020; Potente et al., 2022). On this basis,
453 sharing heterozygous sites among close related *Primula* species could contain plenty
454 of both paralogs and truly heterozygous sites. It could be the reason why a harsh
455 *maxSH* threshold can sharply reduce the number of retained SNPs. While it also
456 suggested that we should depend on other specific ways to filter paralogs when
457 performing population genomics studies on *Primula* and other similar taxa.

458 ise e

459 **4.3 Differentiated behavior of similar pipelines and approaches**

460 The flourishing of sequencing technology has prompted software development around
461 all aspects of downstream analyses. Yet differences between underlying algorithms
462 and logic of different software can lead the same analysis to divergent inference. For
463 example, Chen et al. (2021) demonstrated incredibly different behaviors of two
464 mainstream lines—McDonald-Kreitman (MK) test and PAML test—for positive
465 selection detection. Considering outlier removal is not only prior for inferring
466 unbiased population structure and estimating demographic history, but also crucial for
467 tackling adaptive divergence, numerous approaches have been developed for

468 detecting outliers, mostly lending F_{st} -related statistics as criteria, such as BAYESCAN,
469 OUTFLANK (Whitlock & Lotterhos 2015), SELESTIM (Vitalis et al., 2014) and so on.
470 Unlike them, *pcadapt* identifies outliers regarding their relationship with population
471 structure ascertained with principal component analysis. In this study, F_{st} based
472 BAYESCAN and PCA based *pcadapt* showed distinct efficiency on outlier identification.
473 *pcadapt* could always flag a large number of loci under selection, while BAYESCAN
474 conservatively detected quite a few number of outliers, furthermore, for only part of
475 parameter assemblies. Likewise, large discrepancy and limited intersection of outliers
476 detected by BAYESCAN and *pcadapt* have been elaborated by several studies (e.g.
477 Kotsakiozi et al., 2017, Bekkevold et al., 2019). Nevertheless, quite few studies took
478 interpretation of their discrepancy into consideration. Luu et al. (2017) pointed out the
479 power of BAYESCAN decreased sharply when admixed individuals are included, which
480 has been verified by our results and partially explains the distinct behavior of
481 BAYESCAN and *pcadapt* in this study. Except for the influence of admixed individuals,
482 *maxSH* has resulted in extra difference for these two approaches. Turning up *maxSH*
483 can strikingly increase the number of total SNPs, while it can also decrease the
484 genetic distance (F_{st} value) among populations. Thereby those F_{st} methods would be
485 more vulnerable to tuning *maxSH*. This could mainly explain the decrease of detected
486 outliers when increasing *maxSH* and the absence of outliers in some parameter
487 assemblies when *maxSH* = 0.25 and 0.50, while non-sensitive PCADAPT will detect
488 more outliers along with the increasing total SNPs.

489 STRUCTURE and PCA are some of the most popular approaches for tackling
490 population assignment in population genetics/genomics. In this study, results of these
491 two approaches are mostly concordant. Yet, we also find a number of differences on
492 how they respond to variations of parameter assemblies. Tuning *maxSH* has brought
493 greater impact on PCA, because when *maxSH* =0.10, part of putative hybrids
494 constantly can not separate from or *P. florindae* or *P. alpicola* in any PCs. Yet when
495 clustering threshold = 0.85 and 0.90, putative hybrids still kept their genetic
496 admixture in STRUCTURE. According to the identity of putative hybrids. Genetic
497 clustering by STRUCTURE could have offered robust results in this study. In *c95m10*,
498 both STRUCTURE and PCA have showed the fusion of genetic clusters, This could be a
499 result of oversplitting effect on reducing genetic distances among populations due to
500 extreme clustering threshold values (Harvey et al., 2015; Rodríguez-Ezpeleta et al.,
501 2016). At last, four putative hybrids were repetitively drifted to others. According to
502 their higher missing data rate than other putative hybrids, we inferred their irregular
503 drift should stem from the impute limitations for missing data in PCA (Yi and Latch,
504 2021).

505 Demographic modelling was most vulnerable to *maxSH* variation among all
506 population genomics analyses in this study. As we performed demographic modelling
507 using FASTSIMCOAL2, the filtration of sharing heterozygous sites can intensively alter
508 site frequency spectrum, the most important impute information for modelling. Thus,
509 estimation of the best model would be close related to tuning *maxSH*. Taking
510 consideration of increasing popularity and significance of coalescent simulation in

511 population genomics, only a handful of studies shed light on how bioinformatics
512 processes affect demographic modelling or inference (Harvey et al., 2015; Shafer et
513 al., 2017). Based on our results, we would like to highlight again the extreme
514 importance of precise establishment of orthologs before performing demographic
515 modelling for reliable biological inference.

516 **Conclusions**

517 Overall, this study highlights the feasibility but risk of tuning *maxSH* values on
518 filtering paralogs. Our results illustrate its remarkable effect on almost all downstream
519 analyses within a population genomics framework. According to the mutual
520 confirmation between population assignment and demographic modelling, we inferred
521 that $maxSH = 0.10$ has brought excessive and asymmetrical removal of truly sharing
522 heterozygous sites as paralogs into this study. On this basis, we tend to approve the
523 hybrid origin of putative hybrids between *P. alpicola* and *P. florindae* with an
524 asymmetrical gene flow pattern deserving further investigation.

525 Here we give some suggestions on how to minimize the influence of *maxSH* from
526 excessive and asymmetrical removal of heterozygous sites. Foremost, no single value
527 could be expected for *maxSH* to be universal for all studies. Setting optimal clustering
528 threshold following McCartney-Melstad et al (2019) would be beneficial as we have
529 demonstrated the amplified biases from interactions between clustering threshold and
530 *maxSH*. Then no matter what kind of analyses are arranged for closely related
531 lineages, especially those with potential hybridization or introgression, one should not
532 rely on *maxSH* for filtering paralogs, when we use PYRAD to generate datasets for

533 population genomics study, do not forget to turn up this parameter as the number of
534 half samples (same as the default value of IPYRAD). And if we use IPYRAD, just keep it
535 as default value.

536

537 **Acknowledgements**

538 This work was supported by the National Natural Science Foundation of China
539 (U1202261, U1602263), We especially thanks the South-East Tibetan Plateau Station
540 for Integrated Observation and Research of Alpine Environment for fieldwork
541 assistance. We also thanks Joaquín Ortego, Evan McCartney-Melsted, Julio Rozas for
542 their help with data analysis. The authors declare that there are no conflict of interests.

543

544 **Data Accessibility**

545 Sequencing Data for this study will be available as soon as acceptance at NCBI
546 Sequence Read Archives as PRJNA669915.

547

548 **References**

- 549 Andrews KR, Good JM, Miller MR, Luikart G, Hohenlohe PA. 2016. Harnessing the power of
550 RADseq for ecological and evolutionary genomics. *Nature Reviews Genetics*, 17: 81-92.
- 551 Arnold B, Corbett-Detig RB, Hartl D, Bomblies K. 2013. RADseq underestimates diversity and
552 introduces genealogical biases due to nonrandom haplotype sampling. *Molecular Ecology*,
553 22(11): 3179-3190.
- 554 Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA, Johnson,
555 EA. 2008. Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers. *Plos*
556 *One*, 3(10): e3376.
- 557 Bekkevold D, Höjesjö J, Nielsen EE, Aldvén D, Als TD, Sodeland M, Kent MP, Lien S, Hansen MM.
558 2020. Northern European *Salmo trutta* (L.) populations are genetically divergent across
559 geographical regions and environmental gradients. *Evolutionary Applications*, 13(2): 400-416.

- 560 Burnham KP, Anderson DR. 1998. *Model selection and inference: a practical information-theoretic*
561 *approach*. New York: Springer.
- 562 Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA. 2013. Stacks: An analysis tool set for
563 population genomics. *Molecular Ecology*, 22(11): 3124-3140.
- 564 Cavender-Bares J, Gonzalez-Rodriguez A, Eaton DAR, Hipp AAL, Beulke A, Manos PS. 2015.
565 Phylogeny and biogeography of the American live oaks (*Quercus* subsection *Virentes*): a
566 genomic and population genetics approach. *Molecular Ecology*, 24(14): 3688-3687.
- 567 Chen, QP, Yang H, Feng X, Chen QJ, Shi SH, Wu CI, He ZW. 2021. Two decades of suspect evidence
568 for adaptive molecular evolution-Negative selection confounding positive selection signals.
569 *National Science Review*. doi: 10.1093/nsr/nwab217
- 570 Chong Z, Ruan J, Wu CI. 2012. Rainbow: an integrated tool for efficient clustering and assembling
571 RAD-seq reads. *Bioinformatics*, 28(21): 2732-2737.
- 572 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, 1000 Genomes
573 Project Analysis Group. 2011. The variant call format and VCFtools. *Bioinformatics*, 27(15):
574 2156–2158.
- 575 Darwin C. 1877. *The Different Forms of Flowers on Plants of the Same Species*. London: John Murray.
- 576 Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML. 2011. Genome-wide genetic
577 marker discovery and genotyping using next-generation sequencing. *Nature reviews. Genetics*,
578 12: 499-510.
- 579
- 580 de Vos, JM, Wüest RO, Conti E. 2014. Small and ugly? Phylogenetic analyses of the "selfing
581 syndrome" reveal complex evolutionary fates of monomorphic primrose flowers. *Evolution*,
582 68(4): 1042-1057.
- 583 Doyle JJ, Doyle JL. 1987. A Rapid DNA Isolation Procedure from Small Quantities of Fresh Leaf
584 Tissues. *Phytochem Bull*, 19: 11-15.
- 585 Eaton DAR. 2014. PYRAD: assembly of de novo RADseq loci for phylogenetic analyses.
586 *Bioinformatics*, 30(13): 1844-1849.
- 587 Eaton DAR, Hipp AL, González-Rodríguez A, Cavender-Bares J. 2015. Historical introgression among
588 the American live oaks and the comparative nature of tests for introgression. *Evolution*,
589 69(10): 2587-2601.
- 590 Eaton DAR, Overcast I. 2020. IPYRAD: Interactive assembly and analysis of RADseq datasets.
591 *Bioinformatics*, 36(8): 2592-2594.
- 592 Escudero M, Eaton DAR, Hahn M, Hipp AL. 2014. Genotyping-by-sequencing as a tool to infer
593 phylogeny and ancestral hybridization: A case study in *Carex* (Cyperaceae). *Molecular*
594 *Phylogenetics and Evolution*, 79: 359–367.

- 595 Evanno G, Regnaut S, Goudet J. 2005. Detecting the number of clusters of individuals using the
596 software STRUCTURE: A simulation study. *Molecular Ecology*, 14(8): 2611-2620.
- 597 Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. 2013. Robust Demographic Inference
598 from Genomic and SNP Data. *Plos Genetics*, 9(10): e1003905
- 599 Fijarczyk A, Babik W. 2015. Detecting balancing selection in genomes: limits and prospects. *Molecular
600 ecology*, 24(14): 3529-3545.
- 601 Fitz-Gibbon S, Hipp AL, Pham KK, Manos PS, Sork VL. 2017. Phylogenomic inferences from
602 reference-mapped and de novo assembled short-read sequence data using RADseq sequencing
603 of California white oaks (*Quercus* section *Quercus*). *Genome*, 60(9): 743-755.
- 604 Foll M, Gaggiotti O. 2008. A genome-scan method to identify selected loci appropriate for both
605 dominant and codominant markers: a Bayesian perspective. *Genetics*, 180(2): 977-993.
- 606
- 607 Garrick RC, Sunnucks P, Dyer RJ. 2010. Nuclear gene phylogeography using PHASE: dealing with
608 unresolved genotypes, lost alleles, and systematic bias in parameter estimation. *BMC
609 evolutionary biology*, 10(1): 1-17.
- 610 Gautier M, Gharbi K, Cezard T, Foucaud J, Kerdelhué C, Pudlo P, Cornuet J, Estoup A. 2013. The
611 effect of RAD allele dropout on the estimation of genetic variation within and between
612 populations. *Molecular Ecology*, 22(11): 3165-3178.
- 613 Gossmann TI, Keightley PD, Eyre-Walker A. 2012. The effect of variation in the effective population
614 size on the rate of adaptive molecular evolution in eukaryotes. *Genome Biology and
615 Evolution*, 4(5): 658-667.
- 616 Harvey MG, Judy CD, Seeholzer GF, Maley JM, Graves GR, Brumfield RT. 2015. Similarity
617 thresholds used in DNA sequence assembly from short reads can reduce the comparability of
618 population histories across species. *PeerJ*, 3: e895.
- 619 Huang H, Knowles LL. 2016. Unforeseen consequences of excluding missing data from next-
620 generation sequences: simulation study of RAD sequences. *Systematic Biology*, 65(3): 357-
621 365.
- 622 Ilut DC, Nydam ML, Hare MP. 2014. Defining loci in restriction-based reduced representation genomic
623 data from nonmodel species: sources of bias and diagnostics for optimal clustering. *BioMed
624 Research International*, 2014: 675158.
- 625 Innan H, Kondrashov F. 2010. The evolution of gene duplications: classifying and distinguishing
626 between models. *Nature Reviews Genetics*, 11(2): 97-108.
- 627 Kotsakiozi P, Richardson JB, Pichler V, Favia G, Martins AJ, Urbanelli S, Armbruster PA, Caccone A.
628 2017. Population genomics of the Asian tiger mosquito, *Aedes albopictus*: insights into the
629 recent worldwide invasion. *Ecology and evolution*, 7(23): 10143-10157.

- 630 Lischer HEL, Excoffier L, Heckel G. 2014. Ignoring heterozygous sites biases phylogenomic estimates
631 of divergence times: implications for the evolutionary history of *Microtus* voles. *Molecular*
632 *Biology and Evolution*, 31(4): 817-831.
- 633 Luu K, Bazin E, Blum MG. 2017. *pcadapt*: an R package to perform genome scans for selection based
634 on principal component analysis. *Molecular Ecology Resources*, 17(1): 67-77.
- 635 Ma YP, Xie WJ, Tian XL, Sun WB, Wu ZK, Richard M. 2014. Unidirectional hybridization and
636 reproductive barriers between two heterostylous primrose species in north-west Yunnan,
637 China. *Annals of Botany*, 113(5): 763-775.
- 638 Mao XG, Zhang JP, Zhang SY, Rossiter SJ. 2010. Historical male-mediated introgression in horseshoe
639 bats revealed by multilocus DNA sequence data. *Molecular Ecology*, 19(7): 1352-1366.
- 640 Mastretta-Yanes A, Arrigo N, Alvarez N, Jorgensen TH, Piñero D, Emerson, BC. 2015. Restriction site-
641 associated DNA sequencing, genotyping error estimation and de novo assembly optimization
642 for population genetic inference. *Molecular Ecology Resources*, 15(1): 28-41.
- 643 McCartney-Melstad E, Gidis M, Shaffer HB. 2019. An empirical pipeline for choosing the optimal
644 clustering threshold in RADseq studies. *Molecular Ecology Resources*, 19(5): 1195-1204.
- 645 McKinney GJ, Waples RK, Seeb LW, Seeb JE. 2017. Paralogs are revealed by proportion of
646 heterozygotes and deviations in read ratios in genotyping-by-sequencing data from natural
647 populations. *Molecular Ecology Resources*, 17(4): 656-669.
- 648 Nordborg M, Hu TT, Ishino Y, Jhaveri J, Toomajian C, Zheng HG, Bakker E, Calabrese P, Gladstone J,
649 Goyal R, Jakobsson M, Kim S, Morozov Y, Padhukasahasram B, Plagnol V, Rosenberg NA,
650 Shah C, Wall JD, Wang J, Zhao KY, Kalbfleisch T, Schulz V, Kreitman M, Bergelson, J.
651 2005. The Pattern of Polymorphism in *Arabidopsis thaliana*. *PLoS biology*, 3(7): e196.
- 652 Ortego J, Gugger PF, Sork VL. 2018. Genomic data reveal cryptic lineage diversification and
653 introgression in Californian golden cup oaks (section *Protobalanus*). *New Phytologist*, 218(2):
654 804-818.
- 655 O'Leary SJ, Puritz JB, Willis SC, Hollenbeck CM, Portnoy DS. 2018. These aren't the loci you'e
656 looking for: Principles of effective SNP filtering for molecular ecologists. *Molecular Ecology*,
657 27(16): 3193-3206.
- 658
- 659 Papadopoulou A, Knowles LL. 2015. Species-specific responses to island connectivity cycles: Refined
660 models for testing phylogeographic concordance across a Mediterranean Pleistocene
661 Aggregate Island Complex. *Molecular Ecology*, 24(16): 4252-4268.
- 662 Paris JR, Stevens JR, Catchen JM. 2017. Lost in parameter space: A road map for stacks. *Methods in*
663 *Ecology and Evolution*, 8(10): 1360-1373.

- 664 Pina-Martins F, Silva DN, Fino J, Paulo OS. 2017. *Structure_threader*: An improved method for
665 automation and parallelization of programs STRUCTURE, FASTSTRUCTURE and *Maverick* on
666 multicore CPU systems. *Molecular Ecology Resources*, 17(6): e268–e274.
- 667 Pina-Martins F, Baptista J, Pappas Jr. G, Paulo OS. 2019. New insights into adaptation and population
668 structure of cork oak using genotyping by sequencing. *Global Change Biology*, 25(1): 337-
669 350.
- 670 Pritchard JK, Stephens M, Donnelly P. 2000. Inference of population structure using multilocus
671 genotype data. *Genetics*, 155(2): 945–959.
- 672 Puritz JB, Hollenbeck CM, Gold JR. 2014. *dDocent*: a RADseq, variant-calling pipeline designed for
673 population genomics of non-model organisms. *PeerJ*, 2: e431.
- 674 R Core Team. 2018. *R: A language and environment for statistical computing*. Vienna: R Foundation
675 for Statistical Computing.
- 676 Richards J. 2003. *Primula*. Oregon: Timber Press, Inc.
- 677 Rodríguez-Ezpeleta N, Bradbury IR, Mendibil I, Álvarez P, Cotano U, Irigoien X. 2016, Population
678 structure of Atlantic mackerel inferred from RAD-seq-derived SNP markers: effects of
679 sequence clustering parameters and hierarchical SNP selection. *Molecular Ecology Resources*,
680 16(4): 991-1001.
- 681 Rozas J, Ferrer-Mata A, Sánchez-DelBarrio JC, Guirao-Rico S, Librado P, Ramos-Onsins SE, Sánchez-
682 Gracia A. 2017. DnaSP 6: DNA Sequence Polymorphism analysis of Large
683 Datasets. *Molecular Biology and Evolution*, 34(12): 3299-3302.
- 684 Ren T, Yang Y, Zhou T, Liu ZL. 2018. Comparative plastid genomes of *Primula* species: Sequence
685 divergence and phylogenetic relationships. *International journal of molecular sciences*, 19(4):
686 1050.
- 687 Shafer ABA, Peart CR, Tusso S, Maayan I, Brelsford A, Wheat CW, Wolf JBW. 2017. Bioinformatic
688 processing of RAD-seq data dramatically impacts downstream population genetic inference.
689 *Methods in Ecology and Evolution*, 8(8): 907-917.
- 690 Sota T, Vogler AP. 2003. Reconstructing species phylogeny of the carabid beetles *Ohomopterus* using
691 multiple nuclear DNA sequences: heterogeneous information content and the performance of
692 simultaneous analyses. *Molecular Phylogenetics and Evolution*, 26(1): 139-154.
- 693 Spofford JB. 1969. Heterosis and the evolution of duplications. *The American Naturalist*, 103(932):
694 407-432.
- 695 Tonzo V, Papadopoulou A, Ortego, J. 2020. Genomic footprints of an old affair: Single nucleotide
696 polymorphism data reveal historical hybridization and the subsequent evolution of
697 reproductive barriers in two recently diverged grasshoppers with partly overlapping
698 distributions. *Molecular Ecology*, 29(12): 2254-2268.

- 699 Twyford AD, Ennos RA. 2012. Next-generation hybridization and introgression. *Heredity*, 108(3): 179-
700 189.
- 701 Vitalis R, Gautier M, Dawson KJ, Beaumont MA. 2014. Detecting and measuring selection from gene
702 frequency data. *Genetics*, 196 (3): 799–817.
- 703 Whitlock MC, Lotterhos KE. 2015. Reliable Detection of Loci Responsible for Local Adaptation:
704 Inference of a Null Model through Trimming the Distribution of F_{ST} . *The American Naturalist*,
705 186(S1): S24-S36.
- 706 Xie YP, Zhu XF, Ma YP, Zhao JL, Li L, Li QJ. 2017. Natural hybridization and reproductive isolation
707 between two *Primula* species. *Journal of Integrative Plant Biology*, 59(8): 526-530.
- 708 Yi X, Latch EK. 2022. Nonrandom missing data can bias Principal Component Analysis inference of
709 population genetic structure. *Molecular Ecology Resources*, 22(2): 602-611.
- 710 Zhu XF, Li Y, Wu GL, Fang ZD, Li QJ, Liu JQ. 2009. Molecular and morphological evidence for
711 natural hybridization between *Primula secundiflora* franchet and *P. poissonii* franchet
712 (Primulaceae). *Acta Biologica Cracoviensia*, 51(2): 29-36.

713 **Tables**714 **Table 1** Summary of IPYRAD output and outlier detection of different parameter assemblies

Parameter assemblies	Total loci	Total SNPs	Filtered SNPs	BAYESCAN outliers	<i>pcadapt</i> outliers	Outliers by both	Neutral SNPs
c85m10	33328	226965	5276	25	267	18	5002
c85m25	45671	351147	14374	0	816	0	13558
c85m50	47652	371317	16346	0	836	0	15510
c90m10	44783	307539	7076	33	342	21	6722
c90m25	58940	453269	17362	0	876	0	16486
c90m50	61046	475112	19456	0	1130	0	18326
c95m10	70321	427733	10870	73	1948	60	8909
c95m25	86620	586421	22187	36	1090	0	21061
c95m50	88630	605829	24166	25	2015	0	22126

715 Notes: Total loci and Total SNPs were generated by IPYRAD with at least 20% individuals
716 containing data at a given locus; Filtered SNPs were generated by total SNPs further filtered by
717 missing data, minimum allele frequency and keeping only the center one SNP per locus; BAYESCAN
718 outliers, outliers detected by the software BAYESCAN; *pcadapt* outliers, outliers detected by the R
719 package *pcadapt*; Neutral SNPs, filtered SNPs with detected outlier removed.

720 Figure Legends

721 **Fig. 1.** Proportion of loci flagged as paralogs and filtered by IPYRAD. Results are
722 grouped according to clustering threshold in order to highlight differences of flagged
723 paralogs resulted from *maxSH*.

724

725 **Fig. 2.** Comparison of genetic clustering by the bayesian clustering approach for *P.*
726 *alpicola*, *P.florindae* and their putative hybrids implemented in the program
727 STRUCTURE, Only $K = 2$ is shown here for being the Best K value. Each column
728 shared the same clustering threshold and each row shared the same *maxSH* value. In
729 every plot, each individual is represented by a vertical bar for every independent plot
730 and color composition of each bar is referred to the individual's ancestry.

731

732 **Fig. 3.** Comparison of genetic clustering by PCA approach for *P. alpicola*, *P. florindae*
733 and their putative hybrids. Each column shared the same clustering threshold and each
734 row shared the same *maxSH* value. In every plot, each individual is represented by a
735 dot and dots are colored according to sampling classification.

736

737 **Fig. 4.** The optimal demographic model for nine tested parameter assembly indicated
738 by AIC and ΔAIC computation. Value of estimated parameters for the best model are
739 showed in each plot, including divergence time (T_{DIV}), admixture time (T_{ADMIX}) for
740 admixture model, effective population size (θ), rates of gene flow (m), and proportion
741 of lineages transfer (α).

742 Support Information

743 **Fig. S1.** Alternative demographic models for exploring the origin of putative hybrids.

744 The only difference between upper three models and lower three is the existence of
745 interspecific gene flow. Parameter estimation include divergence time (T_{DIV}),
746 admixture time (T_{ADMIX}) for admixture model, effective population size (θ), rates of
747 gene flow (m), and proportion of lineages transfer (α).

748 **Fig. S2.** Results of genetic clustering by the BAYESCAN clustering approach for
749 p95_60_10 after removing six samples totally similar to *P. florindae* in genetics. Only
750 $K = 2$ is shown here for being the Best K value. Each individual is represented by a
751 vertical bar for every independent plot and color composition of each bar is referred to
752 the individual's ancestry.

753 **Fig. S3.** Demographic model statistically equivalent to the best model for p95_60_10.
754 Value of estimated parameters includes divergence time (T_{DIV}), admixture time
755 (T_{ADMIX}) for admixture model, effective population size (θ), rates of gene flow (m),
756 and proportion of lineages transfer (α).

757 **Table S1.** Comparison of demographic models for *c85m10*.

758 **Table S2.** Comparison of demographic models for *c90m10*.

759 **Table S3.** Comparison of demographic models for *c95m10*.

760 **Table S4.** Comparison of demographic models for *c85m25*.

761 **Table S5.** Comparison of demographic models for *c90m25*.

762 **Table S6.** Comparison of demographic models for *c95m25*.

763 **Table S7.** Comparison of demographic models for *c85m50*.

764 **Table S8.** Comparison of demographic models for *c90m50*.

765 **Table S9.** Comparison of demographic models for *c95m50*.

766 **datafile S1.** Complete information of each parameter assembly for running IPYRAD.

- 767 **datafile S2.** The best K estimation for each datasets according to ΔK method.
- 768 **datafile S3.** Interactive version of plots of STRUCTURE results for all K values.