# JML: testing hybridization from species trees

Simon Joly

Institut de recherche en biologie végétale, Université de Montréal and Jardin botanique de Montréal, 4101 Sherbrooke Est, Montréal (QC) H1X 2B2, Canada. Phone : +1 514-872-0344; Email : simon.joly@umontreal.ca

1    **Abstract.** I introduce the software JML that tests for the presence of hybridization in multi-species

2    sequence datasets by posterior predictive checking following Joly, McLenachan and Lockhart

3    (2009, American Naturalist 174:e54-e70). Although their method could potentially be applied on

4    any dataset, the lack of appropriate software made its application difficult. The software JML thus

5    fills a need for an easy application of the method, but also includes improvements such as the

6    possibility to incorporate uncertainty in the species tree topology. The JML software uses a

7    posterior distribution of species trees, population sizes and branch lengths to simulate replicate

8    sequence datasets using the coalescent with no migration. A test quantity, defined as the

9    minimum pairwise sequence distance between sequences of two species, is then evaluated on the

10   simulated datasets and compared to the one estimated from the original data. Because the test

11   quantity is a good predictor of hybridization events, departure from the bifurcating species tree

12   model could be interpreted as evidence of hybridization. Software performance in terms of

13   computing time is evaluated for several parameters. I also show an application example of the

14   software for detecting hybridization among native diploid North American roses.

## Introduction

16  Hybridization is an important evolutionary process (Arnold 1997; Barton 2001). Its role in

17  speciation (Mallet 2007; Rieseberg 1997; Rieseberg *et al.* 2003; Seehausen 2004) and adaptation

18  (Arnold 2004; Joly & Schoen 2011) is understood theoretically and has also been confirmed

19  experimentally. Yet, the role of hybridization is hard to confirm in many instances because it is

20  often difficult to find statistical evidence for hybridization. Here, the term hybridization is used in

21  the broad sense. That is, it refers both to the event, the successful mating between individuals

22  from two distinct species, and its outcomes: hybrid speciation and introgression, where

23  introgression is the transfer of genetic material between species via sexual reproduction.

24  Typically, hybridization is detected using measures of gene tree incongruence (Funk & Omland

25  2003), either among gene trees or between the gene tree and the species tree, although other

26  processes can be in cause. Thus, distinguishing between hybridization and other processes

27  resulting in gene tree incongruence is a critical issue in evolutionary biology. A specific question

28  that has received a lot of attention is that of distinguishing incongruence caused by introgression

29  from that caused by incomplete lineage sorting. Incomplete lineage sorting arises when ancestral

30  polymorphisms present in the ancestral species have not been completely sorted out by genetic

31  drift in the daughter species, resulting in non-monophyletic species. Even though several methods

32  have been described to address this problem, none provide a clear and general test for the

33  presence of hybridization (reviewed in Joly *et al.* 2009).

34  Joly *et al.* (2009) proposed a method based on the idea that incomplete lineage sorting imposes a

35  limit to the minimum expected distance between sequences of two species because the sequences

36  compared have been diverging since the speciation event. Such limit does not exist for

37  introgressed sequences. Consequently, it should be possible to statistically identify introgressed

38  sequences when the pairwise distance between sequences found in two distinct species is smaller

39 than that expected under a lineage sorting scenario. Simulations have confirmed that this statistic

40 is able to detect introgression, although the success rate depends on several parameters: the

41 relative timing of the hybridization and of speciation events, the population sizes and the

42 sequence length (Joly *et al.* 2009). The method of Joly *et al.* (2009) has the potential to be

43 applied on any dataset, but the lack of software implementing the method has limited its use.

44 Here, I introduce the software JML that implements the posterior predictive approach of Joly *et al.*

45 (2009). I also improve the original approach by accounting for the uncertainty in the species tree

46 topology.

## Formal description of the test

48 In JML, posterior predictive checking is used to test for the presence of hybridization. The

49 program uses as input a posterior distribution of species trees (*S*) with branch lengths (*l*) and

50 population sizes (*θ*). This posterior distribution is generally defined as

$$P(S, l, \theta | D) = \int_G \left( \prod_{i=1}^{n} P(d_i | g_i) P(g_i | S) \right) P(S) dG.$$

51 *D* is the data that consist of *n* multiple sequence alignments ($d_i$). The equation integrates over all

52 possible gene trees (*G*) for all alignments, and $g_i$ represents one specific gene tree. $P(d_i|g_i)$ is the

53 likelihood of the data given the gene tree (Felsenstein 1981), $P(g_i|S)$ is the multispecies

54 coalescent (Degnan & Rosenberg 2009; Rannala & Yang 2003), and *P(S)* is the prior on species

55 trees.

56 Replicated datasets are simulated from the posterior distribution $P(S,l,\theta|D)$. A test quantity is then

57 estimated on the observed data and on the simulated datasets to see how well the model is

58 consistent with the data. This approach of posterior predictive checking is commonly used in

59 Bayesian analyses to check the adequacy of a model (Gelman *et al.* 2004); if the test quantity

60    estimated on the observed data departs strongly from the quantities estimated from the simulated

61    data, then we can conclude that the model is inadequate. Here, the test quantity used is the

62    minimum pairwise distance between sequences of two species (*minDist*), which has been shown

63    to be a useful quantity for detecting hybridization (Joly *et al.* 2009). In the presence of

64    hybridization, *minDist* can sometimes be much smaller than that expected in a scenario without

65    hybridization. Suppose that *minDist(AB)* represents *minDist* between species *A* and *B* on the

66    observed data and that *minDist(AB)$^{sim}$* represents *minDist* between species *A* and *B* on simulated

67    data. The *p*-value for hybridization between species *A* and *B* is

$$p = \Pr\left(minDist(AB) < minDist(AB)^{sim}\right).$$

68    The probability is taken over the posterior distribution of parameters *S*, *l*, and *θ* (i.e., *P(S,l,θ|D)*)

69    and the posterior predictive distribution of *minDist(AB)$^{sim}$*. This probability can be approximated

70    by simulation. If we simulate *M* datasets from the posterior distribution *P(S,l,θ|D)*, we can

71    calculate *minDist*(AB)$^{sim(m)}$ on each simulated dataset *m* and the *p*-value is the proportion of these

72    *m* simulations for which *minDist(AB)* < *minDist(AB)$^{sim(m)}$*. If the model is good, then Pr(

73    *minDist(AB)* < *minDist(AB)$^{sim}$* ) ≈ 0.5. On the contrary, a small *p*-value will indicate that the

74    model doesn't fit the data well. Because a small value is characteristic of hybrid sequences in a

75    dataset, one can tentatively conclude that the inaccuracy of the model is due to the presence of

76    hybrid sequences.

77    **Implementation**

78    Incorporating species tree topology uncertainty in posterior predictive checking represents an

79    improvement compared to the original description of the method where the species tree topology

80    was fixed (Joly *et al.* 2009). This is done by using as input the posterior distribution obtained

81    from *BEAST analyses (Heled & Drummond 2010; Drummond & Rambaut 2007). *BEAST is a

82    Bayesian method that estimates the posterior distribution of species trees, branch lengths and

83    population sizes using sequence information from multiple genes. Note that posterior

84    distributions from other programs could also be used in JML as long as the tree file is in the same

85    format as the *BEAST nexus format. For the simulations, species trees (with branch lengths and

86    populations sizes) are sampled from the stationary phase of the Markov Chain Monte Carlo.

87    For each species tree, a gene tree is then simulated using the coalescent. The code for the gene

88    tree simulation routine was adapted from MCMCcoal (Yang 2007). The number of gene copies

89    simulated per species is the same as in the original dataset. The user can scale the species tree

90    population sizes using a heredity scalar to reflect the effective population size of the marker being

91    simulated. Similarly, the mutation rate of the species tree can also be scaled for the simulations to

92    allow the possibility that the mutation rate of the marker being simulated is not the same as the

93    mutation rate implied in the species tree.

94    Sequences are then simulated on the gene tree. This was implemented by adapting the code of the

95    software seq-gen 1.3.2 (Rambaut & Grassly 1997), which allows any nucleotide substitution

96    model to be used. This procedure is repeated for all species tree of the posterior distribution (or a

97    subset of them). Finally, JML outputs the posterior predictive distribution of the smallest distances

98    between sequences of any two species of the dataset, from which $p$-values could be estimated.

99    JML can also output the exact $p$-value for each pairwise species comparison if the empirical

100   sequence dataset is given.

101   **Interpretation and multiple comparisons**

102   Different approaches can be used for interpreting results from posterior predictive checking. An

103   intuitive one is to interpret the $p$-value(s) directly. The $p$-values estimated by JML are posterior

104   probabilities (Gelman *et al.* 2004) and can be interpreted as the probability that the model will

105    generate a minimum distance between sequences of two species smaller than that observed from

106    the data, given the data. However appealing is this interpretation, it could lead to statistical issues

107    when multiple tests are performed. Indeed, the need to correct for multiple statistical testing (Rice

108    1989) diminishes the likelihood of finding statistically significant results. This is especially

109    problematic for the present application because the large variance in mutation rate for short

110    sequences (Edwards & Beerli 2000), combined with the difficulty to get long nucleotide

111    sequence stretches that lack evidence of recombination in practice, result in power issues (Joly *et*

112    *al.* 2009). The problem is even more acute when the approach is used in an explorative way, that

113    is if there are no a priori hypotheses of hybridization to test and if JML is only used to investigate

114    the presence of hybridization in the dataset. In such cases, all pairwise species comparisons can

115    be tested simultaneously and the statistical power will be greatly affected. To minimize power

116    issues it could thus be important to specify hybridization hypotheses a priori without reference to

117    the sequence data.

118    There is an alternative interpretation of posterior predictive checking, which is to see "how

119    particular aspects of the data would be expected to appear in replications" (Gelman *et al.* 2004).

120    For instance, we could evaluate the overall adequacy of a model by assessing if there are some

121    aspects of the data that are not well predicted by the model. To do this, it would be of interest to

122    report all observed distances that have a low probability of being observed, e.g. distances with $p$

123    < 0.1 (this value is arbitrary and can be fixed by the user). This could indicate species

124    comparisons where the model cannot adequately predict the observed minimum distances. If

125    there were several of those instances, one could thus conclude that a strictly bifurcating species

126    tree model is not adequate, probably because of the presence of hybridization. Note, however,

127    that this is not the same as concluding that there has been hybridization between two given

128    species. With such interpretations of posterior predictive distributions, the type I error is less of a

129    concern because we use posterior predictive checking to evaluate the fit of the model rather than

130    to test a specific hypothesis (Gelman *et al.* 2004).

131    Regardless of the multiple comparison issues associated with posterior predictive checking, there

132    are two points that should always be kept in mind when interpreting results from JML. First,

133    posterior predictive checking is a test of the model and not of hybridization. If one rejects the

134    model (bifurcating species tree without gene flow), this may well be because of the presence of

135    hybridization, although it could also be due to other properties of the data such as undetected

136    gene duplication (Maddison 1997), population substructure along the branches of the phylogeny

137    (Machado *et al.* 2002), and parallel evolution (Joly *et al.* 2010). The second point to take into

138    account is that a lack of evidence for hybridization with JML should not be interpreted as an

139    absolute absence of hybridization in the dataset because (1) a lack of statistical significance can

140    also be caused by a lack of data and that (2) not all hybridization events leave a detectable

141    molecular signature (Joly *et al.* 2009, 2006).

## Performance

143    Thorough simulations regarding the performance of the test statistic have already been conducted

144    for several parameters such as sequence length, population size, speciation time and time of the

145    hybridization event (Joly *et al.* 2009). Here I report results on the impact of different parameters

146    values on computing time. The parameters investigated were the number of species (5, 10, 15),

147    the number of sequences per species (5,10,15), the number of simulations (1000, 2000, 4000),

148    and the sequence length (500, 1000, 1500). Random species trees were simulated under a birth

149    and death model with the R package 'geiger' (Harmon *et al.* 2008); the birth and death

150    parameters were set to 0.00025 and 0.000125, respectively, and the phylogeny was evolved for

151    0.01 units of time. These settings resulted in phylogenies with a tree depth (time × mutation rate)

152   similar to that of empirical datasets (Joly *et al.* 2009). The first phylogenies obtained with five,

153   ten and fifteen extant species were retained for the simulations (extinct species were pruned from

154   the tree). Mutational population sizes ($\theta = 4N_e\mu$) of the tree were generated randomly by

155   sampling from a truncated normal distribution with mean and standard deviation of 0.005, with a

156   lower cut-off of 0.0001. Again, this is similar to empirical observations. These phylogenies were

157   treated as "fixed" and JML generated simulated datasets (using the GTR + I + Γ substitution

158   model) using combinations of the parameters mentioned above. Because repeated runs had very

159   small coefficients of variation (0.5%), only one full run was performed for each combination of

160   parameters. Simulations were performed on a HP desktop computer with an Intel core2 duo CPU

161   at 2.33 GHz with 2 Gb of RAM.

162   The results show that the computing time for a complete run grows linearly with the number of

163   datasets simulated (data not shown) and with the sequence length (Fig 1a). In contrast, the

164   computing time increases according to a power function relative to the number of species and

165   relative to the number of sequences per species (Fig. 1b).

166   **An application example—North American roses**

167   To give an application example of the software, I reanalyse here sequence data from three nuclear

168   genes for the native diploid roses of North America. Three nuclear genes (*GAPDH*, *TPI*, *MS*)

169   have been sequenced for 46 individuals from eight species and have been analysed with distances

170   and gene tree parsimony approaches (Joly & Bruneau 2006, 2009). Alleles within individuals

171   were obtained through direct sequencing or via cloning when an individual was heterozygous for

172   a gene (Joly & Bruneau, 2006). Previous studies showed that there might be introgressed

173   sequences in the dataset; i.e. some sequences in one species are often either identical or one

174    mutation away from a sequence of another species (Joly & Bruneau 2006). Yet, no formal tests of

175    hybridization have been conducted to date.

176    Previous studies could not find evidence of recombination in these datasets (Joly & Bruneau

177    2006) and thus the three genes could be analysed integrally. Species tree analyses were

178    performed in *BEAST. The nucleotide substitution model used was the one that received the

179    highest Akaike Information Criteria (AIC) score in Modeltest 3.7 (Posada & Crandall 1998)

180    when fitted on a maximum likelihood tree obtained from five independent searches in Garli 1.0

181    (Zwickl 2006) with a GTR $+ I + \Gamma$ substitution model.  A strict clock was used for all genes; the

182    rate of the *GAPDH* gene was set to 1 and the rate of the other genes were estimated relative to

183    *GAPDH*. Population sizes were modelled as constant along branches. More details on parameters

184    and priors can be found in the .xml file given as supplementary information. The analysis was run

185    for $10^7$ generations, recording the trees and parameters every $10^4$ generations, and the first

186    million generations was discarded as burnin. Independent runs converged on the same parameters

187    values and species tree topologies.

188    The species tree obtained with *BEAST (Fig. 2) was identical to one of the two most

189    parsimonious species trees obtained by gene tree parsimony (Joly & Bruneau 2009). The branch

190    support was relatively high for most nodes, but there is nevertheless clearly some uncertainty in

191    the tree topology which was clearly worth accounting for in the hybridization tests. The wide

192    branches along the backbone of the tree are likely the results of gene tree incongruence, which

193    could be caused by either incomplete lineage sorting or hybridization.

194    The species trees (with branch length and population sizes) estimated by *BEAST were then

195    input into JML and posterior predictive distributions generated for *minDist* between all species for

196    all genes. For each gene, sequences of the same length as the original ones were simulated

197    according to the best substitution model and parameter values as determined by the AIC in

198    ModelTest (see above). The relative mutation rate used in the simulations for each gene was set

199    to the median posterior value obtained from the *BEAST analyses. The species tree from the first

200    million generations were discarded as burnin in JML and the remaining 9000 trees were used for

201    the simulations. Because I did not have a specific hypothesis of hybridization to test, I decided to

202    investigate the overall fit of the model and report all observed distances that had a probability <

203    0.1 of being generated by the posterior distribution.

204    Six distances between alleles were smaller than the $10^{th}$ quantile in the posterior predictive

205    distributions (Table 1). These involved one individual of *Rosa blanda* (incl. *R. woodsii*) and one

206    of *R. pisocarpa*, each with three individuals of *R. gymnocarpa* for the *TPI* gene. Although the

207    observed distances are not statistically significant at the 5% level, they are small enough to

208    suggest that the model does not explain these observations very well. In other words, although

209    there is not statistical evidence for a hybridization event between *R. gymnocarpa* and *R. blanda* /

210    *R. pisocarpa*, the data suggest this could be the case. Hybridization could have occurred in

211    different ways, but most likely towards *R. gymnocarpa* given that *R. gymnocarpa* sequences are

212    nested with a *R. blanda* / *R. pisocarpa* clade (see supplementary Figures), whereas the species

213    tree suggest *R. gymnocarpa* is basal to the other species (Fig. 2). Because both *R. blanda* and *R.*

214    *pisocarpa* share the introgressed allele, the hybridization event could have occurred between

215    either of these species and *R. gymnocarpa* or between the ancestor of *R. blanda* and *R. pisocarpa*

216    and *R. gymnocarpa*. More data are needed to confirm these hypotheses. For instance, the addition

217    of genes might help to narrow down the confidence intervals of the species tree and perhaps

218    provide stronger statistical results in the future.

219    One interesting observation from this example is that although there were several cases of shared

220    alleles between species (*R. nitida* and *R. palustris* (*TPI, MS, GAPDH*); *R. pisocarpa* and *R.*

221    *blanda* (*TPI, MS, GAPDH*), *R. blanda* and *R. foliolosa* (MS), *R. blanda* and *R. nitida* (*TPI*); see

222    supplementary figures), none of these were found to be significant. In other words, even

223    relatively good evidence for the presence of hybridization such as identical sequences between

224    non-sister species does not mean that it is necessarily caused by hybridization. Due to

225    stochasticity in the coalescent process and in the mutation rates for short sequences, it is

226    relatively difficult to statistically infer hybridization events from empirical data. In the present

227    example, only one possible instance of hybridization was confirmed. In this case identical

228    sequences were found in a putative hybrid formed between two of the most diverged species in

229    the group.

230    This application example shows why it is important to test hybridization hypotheses. Lack of

231    significance could mean that hybridization is not responsible for the observed pattern, but it could

232    also stimulate the gathering of additional data to eventually obtain statistical support for

233    hybridization hypotheses. The statistical approach implemented in JML should thus help

234    researchers to attain a better knowledge regarding the presence of hybridization in their study

235    groups and hopefully contribute to better understand the contribution of hybridization to

236    evolution.

237    **Availability**

238    JML is written in C++ and is released under the GNU General Public License 3+. Source code and

239    precompiled binaries can be downloaded from www.plantevolution.org/jml.html. The manual of

240    JML version 1.0 is available as supplementary material.

## Acknowledgements

## Literature Cited

246   Arnold ML (1997) *Natural hybridization and evolution*. Oxford University Press, New York.

247   Arnold ML (2004) Transfer and origin of adaptations through natural hybridization: were
248       Anderson and Stebbins right? *The Plant Cell*, **16**, 562-570.

249   Barton NH (2001) The role of hybridization in evolution. *Molecular Ecology*, **10**, 551-568.

250   Degnan JH, Rosenberg NA (2009) Gene tree discordance, phylogenetic inference and the
251       multispecies coalescent. *Trends in Ecology & Evolution*, **24**, 332-340.

252   Drummond AJ, Rambaut A (2007) BEAST: Bayesian evolutionary analysis by sampling trees.
253       *BMC Evolutionary Biology*, **7**, 214.

254   Edwards SV, Beerli P (2000) Gene divergence, population divergence, and the variance in
255       coalescence time in phylogeographic studies. *Evolution*, **54**, 1839-1854.

256   Felsenstein J (1981) Evolutionary trees from DNA sequences: a maximum likelihood approach.
257       *Journal of Molecular Evolution*, **17**, 368-376.

258   Funk DJ, Omland KE (2003) Species-level paraphyly and polyphyly: frequency, causes, and
259       consequences, with insights from animal mitochondrial DNA. *Annual Reviews in
260       Ecology, Evolution, and Systematics*, **34**, 397-423.

261   Gelman A, Carlin JB, Stern HS, Rubin DB (2004) *Bayesian data analysis*. Chapman & Hall,
262       Boca Raton, FL.

263   Harmon LJ, Weir JT, Brock CD, Glor RE, Challenger W (2008) GEIGER: investigating
264       evolutionary radiations. *Bioinformatics*, **24**, 129 -131.

265   Heled J, Drummond AJ (2010) Bayesian inference of species trees from multilocus data. *Mol
266       Biol Evol*, **27**, 570-580.

267   Joly S, Bruneau A (2006) Incorporating allelic variation for reconstructing the evolutionary
268       history of organisms from multiple genes: an example from Rosa in North America.
269       *Systematic Biology*, **55**, 623-636.

270   Joly S, Bruneau A (2009) Measuring Branch Support in Species Trees Obtained by Gene Tree
271       Parsimony. *Systematic Biology*, **58**, 100-113.

272 Joly S, Schoen DJ (2011) Migration rates, frequency-dependent selection and the self-
273      incompatibility locus in Leavenworthia (Brassicaceae). *Evolution*, **65**, 2357-2369.

274 Joly S, McLenachan PA, Lockhart PJ (2009) A statistical approach for distinguishing
275      hybridization and incomplete lineage sorting. *The American Naturalist*, **174**, e54-e70.

276 Joly S, Pfeil BE, Oxelman B, McLenachan PA, Lockhart PJ (2010) Correction. *The American*
277      *Naturalist*, **175**, 621-622.

278 Joly S, Starr JR, Lewis WH, Bruneau A (2006) Polyploid and hybrid evolution in roses east of
279      the Rocky Mountains. *American Journal of Botany*, **93**, 412-425.

280 Machado CA, Kliman RM, Markert JA, Hey J (2002) Inferring the history of speciation from
281      multilocus DNA sequence data: the case of Drosophila pseudoobscura and close relatives.
282      *Molecular Biology and Evolution*, **19**, 472-488.

283 Maddison WP (1997) Gene trees in species trees. *Systematic Biology*, **46**, 523-536.

284 Mallet J (2007) Hybrid speciation. *Nature*, **446**, 279-283.

285 Posada D, Crandall KA (1998) Modeltest: testing the model of DNA substitution. *Bioinformatics*,
286      **14**, 817-818.

287 Rambaut A, Grassly NC (1997) Seq-Gen: an application for the Monte Carlo simulation of DNA
288      sequence evolution along phylogenetic trees. *Computer Applications in the Biosciences*,
289      **13**, 235-238.

290 Rannala B, Yang Z (2003) Bayes estimation of species divergence times and ancestral population
291      sizes using DNA sequences from multiple loci. *Genetics*, **164**, 1645-1656.

292 Rice WR (1989) Analysing tables of statistical tests. *Evolution*, **43**, 223-225.

293 Rieseberg LH (1997) Hybrid origins of plant species. *Annual Reviews in Ecology and*
294      *Systematics*, **28**, 359-389.

295 Rieseberg LH, Raymond O, Rosenthal DM, Lai Z, Livingstone K, Nakazato T, Durphy JL,
296      Schwarzbach AE, Donovan LA, Lexer C (2003) Major ecological transitions in wild
297      sunflowers facilitated by hybridization. *Science*, **301**, 1211-1216.

298 Seehausen O (2004) Hybridization and adaptive radiation. *Trends in Ecology and Evolution*, **19**,
299      198-207.

300 Yang Z (2007) *MCMCcoal: Markov chain monte carlo coalescent program*. London.

301 Zwickl DJ (2006) Genetic algorithm approaches for the phylogenetic analysis of large biological
302      sequence datasets under the maximum likelihood criterion.

303

304

305     **Table 1**. List of distances with p-values < 0.1 according to the posterior predictive distributions.

| Gene | individual 1 | individual 2 | Obs. Distance | *p*-value |
|------|--------------|--------------|---------------|-----------|
| *TPI* | *R. pisocarpa* 847 | *R. gymnocarpa* 543 | 0 | 0.0529 |
| *TPI* | *R. pisocarpa* 847 | *R. gymnocarpa* 751 | 0 | 0.0529 |
| *TPI* | *R. pisocarpa* 847 | *R. gymnocarpa* 767 | 0 | 0.0529 |
| *TPI* | *R. blanda* 741 | *R. gymnocarpa* 543 | 0 | 0.0812 |
| *TPI* | *R. blanda* 741 | *R. gymnocarp*a 751 | 0 | 0.0812 |
| *TPI* | *R. blanda* 741 | *R. gymnocarpa* 767 | 0 | 0.0812 |

306     *Note: the number designing the individual is the accession number. See Joly et al. (2006) for*
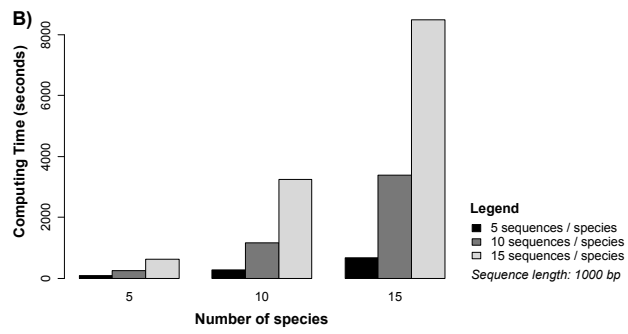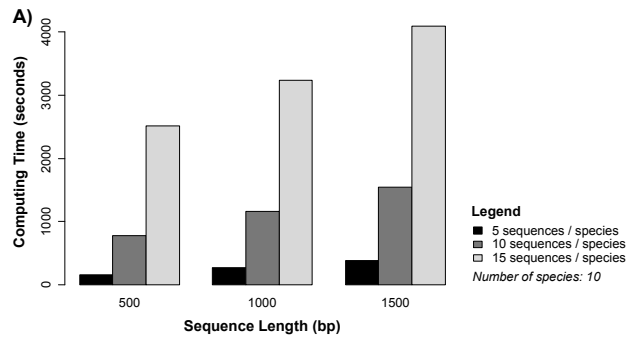307     *more details on accessions.*
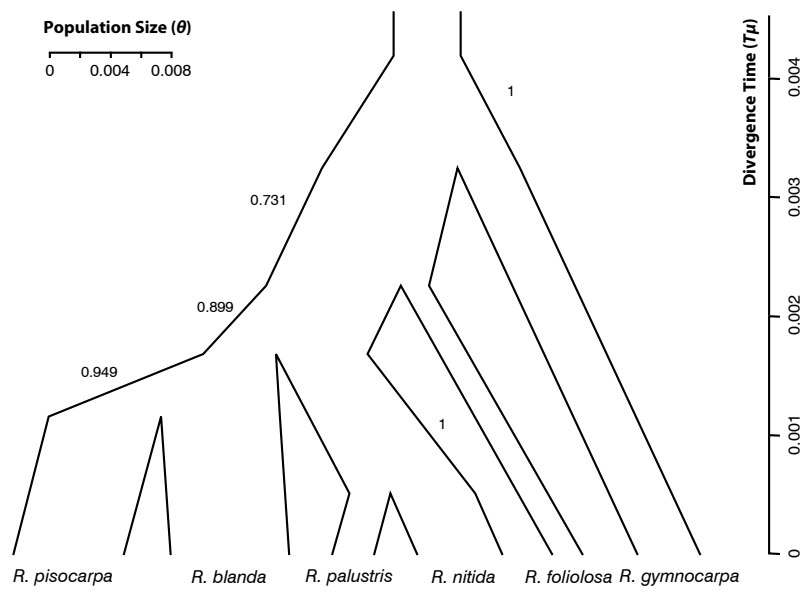308

309 **Figure Legends**

310

311 **Figure 1**.  Performance of the JML software in terms of computing time for (A) different

312 sequence lengths and number of sequences per species, keeping the number of species to 10, and

313 for (B) different number of species and sequences per species, keeping the sequence length to

314 1000 bp.

315

316 **Figure 2**.  Species tree of diploid North American roses obtained with *BEAST. The branch

317 widths are proportional to the estimated population sizes and the branch lengths are proportional

318 to their divergence times (both median estimates). The variations in population sizes along the

319 branches are a consequence of the graphical representation; population sizes were constant along

320 branches and the correct population sizes are those at the beginning of the branches. Numbers

321 besides branches represent the posterior probabilities of the groups. The outgroup (*R. setigera*

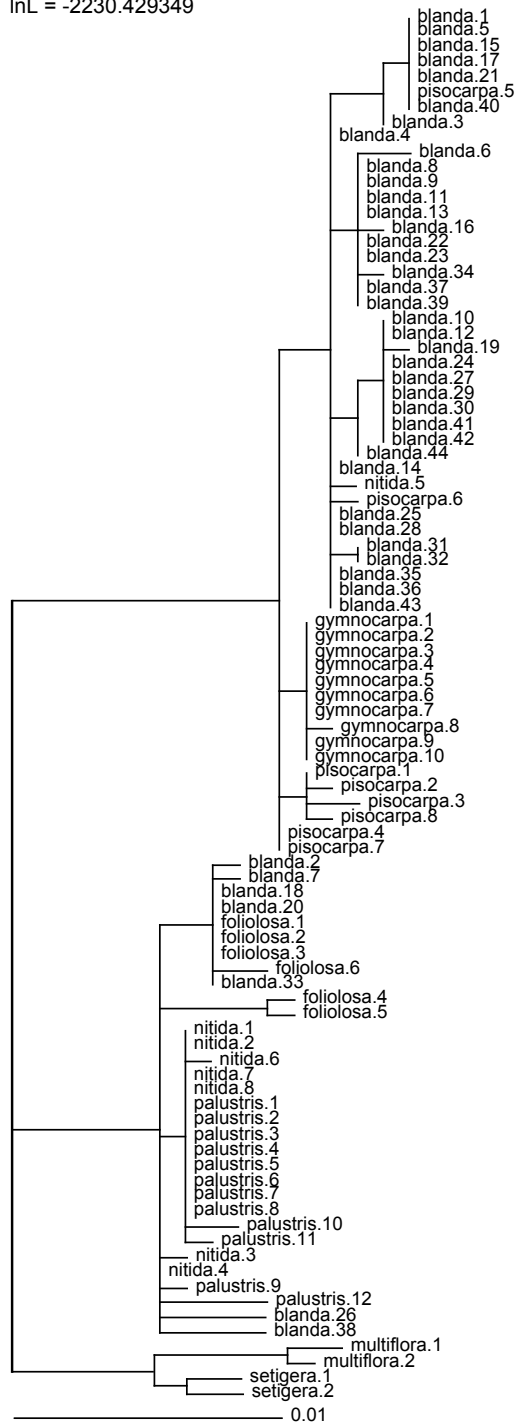322 and *R. multiflora*) is not shown.

**A)**

Computing Time (seconds) vs Sequence Length (bp)

**Legend**
- 5 sequences / species
- 10 sequences / species
- 15 sequences / species

*Number of species: 10*

**B)**

Computing Time (seconds) vs Number of species

**Legend**
- 5 sequences / species
- 10 sequences / species
- 15 sequences / species

*Sequence length: 1000 bp*

**MS**
GTR+I+G
lnL = -2230.429349

**GAPDH**
GTR+I+G
lnL = -1752.771898

**TPI**
GTR+I+G
lnL = -1886.737453

0.01