-	
_	

2	RH: MULTILOCUS DELIMITATION USING ROOTED TRIPLETS
3	
4	A rapid and scalable method for multilocus species delimitation using Bayesian model
5	comparison and rooted triplets
6	
7	Tomochika Fujisawa ¹ , Amr Aswad ^{2,*} , Timothy G. Barraclough ²
8	
9	¹ Department of Zoology, Kyoto University, Sakyo, Kyoto, 606-8502, Japan
10	² Department of Life Sciences, Imperial College London, Silwood Park Campus, Ascot,
11	Berkshire, SL5 7PY, United Kingdom.
12	* Current address: Department of Zoology, University of Oxford, South Parks Road, Oxford,
13	OX1 3PS, United Kingdom
14	
15	Corresponding author: Tomochika Fujisawa, Department of Zoology, Kyoto University,
16	Sakyo, Kyoto, 606-8502, Japan.
17	Email: <u>t.fujisawa05@gmail.com</u>
18	Telephone: +81(0)75-753-4092
19	
20	

21 Abstract

22 Multilocus sequence data provide far greater power to resolve species limits than the single locus data typically used for broad surveys of clades. However, current statistical methods 23 24 based on a multispecies coalescent framework are computationally demanding, because of the number of possible delimitations that must be compared and time-consuming likelihood 25 calculations. New methods are therefore needed to open up the power of multilocus 26 approaches to larger systematic surveys. Here, we present a rapid and scalable method that 27 introduces two new innovations. First, the method reduces the complexity of likelihood 28 calculations by decomposing the tree into rooted triplets. The distribution of topologies for a 29 triplet across multiple loci has a uniform trinomial distribution when the 3 individuals belong 30 to the same species, but a skewed distribution if they belong to separate species with a form 31 32 that is specified by the multispecies coalescent. A Bayesian model comparison framework was developed and the best delimitation found by comparing the product of posterior 33 probabilities of all triplets. The second innovation is a new dynamic programming algorithm 34 35 for finding the optimum delimitation from all those compatible with a guide tree by successively analyzing subtrees defined by each node. This algorithm removes the need for 36 heuristic searches used by current methods, and guarantees that the best solution is found and 37 potentially could be used in other systematic applications. We assessed the performance of the 38 method with simulated, published and newly generated data. Analyses of simulated data 39 demonstrate that the combined method has favourable statistical properties and scalability 40 with increasing sample sizes. Analyses of empirical data from both eukaryotes and 41 prokaryotes demonstrate its potential for delimiting species in real cases. 42

43 Keywords: Multilocus species delimitation, Bayesian model comparison, Dynamic
44 programming, Bacterial species

45 INTRODUCTION

Species constitute the basic taxonomic unit for exchanging information about biological 46 diversity. Defining species boundaries in a consistent manner is therefore of major importance 47 48 to a broad range of biological disciplines. DNA-based delimitation provides a universal method to detect the signature of species existence applicable to various organisms. 49 Consequently, methods to delimit species from DNA sequences alone have been actively 50 developed over the last decade. For early applications of DNA-based delimitation, available 51 markers were limited to a handful of barcoding loci customized for each type of organism 52 (such as *cox1* for animals, Hebert et al. 2003), and therefore delimitation methods were 53 designed to handle these single locus sequences (Pons et al. 2006; Puillandre et al. 2012; 54 Fujisawa and Barraclough 2013; Zhang et al. 2013). However, as the cost of sequencing large 55 56 amounts of DNA has dramatically decreased, and the ease of developing nuclear markers from genome data has increased, the focus has naturally shifted from single to multiple locus 57 approaches. 58

59 There has been huge progress recently in the development of statistical methods for 60 multilocus species delimitation, driven by theoretical advances in the multispecies coalescent model (Rannala and Yang 2003; Degnan and Rosenberg 2009). By comparing alternative 61 delimitation hypotheses and finding the best one based on probability distributions of gene 62 63 trees under the multispecies coalescent model, species can be delimited robustly even with incomplete lineage sorting. Several methods using Bayesian or information theoretic 64 frameworks have been published so far (O'Meara 2010; Yang and Rannala 2010; Ence and 65 Carstens 2011). Empirical studies have evaluated these methods using taxonomically difficult 66 groups (Carstens and Dewey 2010; Hambäck et al. 2013; Satler et al. 2013). Now, the 67 68 multispecies coalescent model is becoming a standard for multilocus DNA-based delimitation,

and there are attempts to integrate these methods with morphology and geography in order to
achieve integrative taxonomy (Fujita et al. 2012; Edwards and Knowles 2014).

One drawback of methods based on the multispecies coalescent model is their limited 71 72 scalability: they rely on the calculation of the probability of obtaining gene trees (or a sequence alignment) given a population tree under the coalescent model, which is relatively 73 74 time consuming. Also, the joint evaluation of species boundaries and species phylogeny requires searches through an enormous parameter space (Yang and Rannala 2014), and 75 computation becomes challenging even with small numbers of sampled individuals. Thus, 76 current procedures for multilocus delimitation often require prior assignments of samples to 77 populations, and they are therefore restricted to validation of candidate delimitations based on 78 79 the assignments. Delimitation without any a priori assignment (species discovery, Ence and 80 Carstens 2011) is feasible only with a limited number of samples, though techniques to reduce search space are being actively studied (Yang and Rannala 2010; Satler et al. 2013). With the 81 increasing ease of sequencing massive multiple nuclear markers (e.g. transcriptome, RAD; 82 83 Baird et al. 2008; anchored hybrid enrichment; Lemmon et al. 2012), the need for rapid and scalable delimitation methods is becoming more urgent. 84

An alternative strategy for potentially scalable multilocus species delimitation is to use 85 genealogical concordance. The congruence of between-species branching across gene trees 86 87 reconstructed for separate loci versus incongruence within species has been used as a signature of reproductive isolation and thereby species diversification (Barraclough et al. 88 2003). Early attempts that used topological congruence to detect species included the 89 delimitation of cryptic fungi using concordance of gene trees inferred from five loci 90 (Koufopanou et al. 1997). The "Genealogical Concordance Phylogenetic Species 91 92 Recognition" (GCPSR, Taylor et al. 2000) is now commonly used to delimit fungal species which often lack morphological or environmental information (Vialle et al. 2013; Millanes et 93

al. 2014). A disadvantage of using concordance measures between multiple gene trees is that 94 it is hard to treat them under statistical models of evolution. It has been known that a set of 95 multiple gene trees do not necessarily "concord" with each other even if they are generated 96 under the same species tree because of the stochastic nature of the coalescent process. 97 Moreover, the consensus topology of gene trees may not be congruent even with the species 98 tree that generated the gene trees (the anomalous gene tree problem, Degnan and Rosenberg 99 2006). Thus, the degree of concordance at which one can confidently infer species is not as 100 simple as first perceived. Modelling the distribution of congruence of trees is intrinsically 101 102 difficult as it must incorporate calculations of the probability of gene trees under a given species tree. Only one non-parametric method with a simulated null model has been devised 103 for statistical delimitation based on topological congruence (O'Meara 2010). 104

Here, we develop a new method for multilocus species delimitation using gene tree 105 congruence, which employs a likelihood model based on the distribution of triplets. We 106 define a triplet as a partial rooted tree consisting of three tips. Using the distribution of rooted 107 108 triplets is a promising approach to model congruence between gene trees under the coalescent framework for two reasons. First, the number of triplets with congruent topology is an 109 intuitive measure of topological similarity between trees. Second, the distribution of triplets is 110 readily tractable under the multispecies coalescent framework (Pamilo and Nei 1988) and has 111 been used successfully for rapid inference of phylogenetic trees (Liu et al. 2010). The 112 distribution model for triplet topology is simple and can be extended for intuitive and rapid 113 model-based delimitation. We tested the performance of the new method with various data 114 sets including simulated gene genealogies and both published and newly generated sequence 115 data from both eukaryotes and microbes. The method provides a tractable approach for 116 multilocus delimitation that is scalable to samples with hundreds of individuals across large 117 clades. 118

119

120 METHODS

121 Calculation of the likelihood of triplet distributions

We employ common assumptions of the multispecies coalescent model (Rannala and Yang 2003; Degnan and Rosenberg 2009): there is neutral random coalescence without structure within species (i.e. panmixia), no gene flow or horizontal transfer between species, and loci evolve independently without intra-locus recombination. In addition, to simplify, we assume initially that the topology of the gene tree is known without error. Under these assumptions, the distribution of triplet topologies is modelled by a simple trinomial distribution as follows.

A bifurcating tree with K tips can be decomposed into $\binom{K}{3}$ rooted triplets. For a given 128 triplet of three individuals, a, b and c, there are three possible topologies, ab|c, ac|b and bc|a. 129 When genealogies from N independent loci are sampled, the numbers of gene trees that 130 conform to each topology - represented by n_1 , n_2 , and n_3 - are modelled by a trinomial 131 distribution for each triplet. When individuals a, b and c belong to a single species, then under 132 our assumption that the species is panmictic, there is an equal probability of observing each of 133 the three triplet topologies because coalescent events of any pair are equally likely in a 134 panmictic population. Therefore, the distribution of counts of the three topologies is 135 represented by an equiprobable trinomial distribution with likelihood: 136

137

$$P_W(n_1, n_2 n_3) = \frac{N}{n_1! n_2! n_3!} \left(\frac{1}{3}\right)^N$$

138 (Eq.1)

140 When individuals are sampled from 2 or 3 distinct species, under the assumptions of the 141 multispecies coalescent process above, the probability of observing triplets congruent with the 142 species tree topology is $1-2e^{-\lambda}/3$, where λ is the length of the internal branch measured by 143 coalescent time units on the species tree, and the probability of observing an incongruent 144 triplet is $e^{-\lambda}/3$ (Pamilo and Nei 1988; Degnan and Rosenberg 2006). Hence, the distribution of 145 triplet counts follows the skewed trinomial distribution.

146

$$P_B(n_1, n_2, n_3 | \nu = n_1, \lambda) = \frac{N}{n_1! n_2! n_3!} \left(1 - \frac{2}{3}e^{-\lambda}\right)^{n_1} \left(\frac{1}{3}e^{-\lambda}\right)^{n_2 + n_3}$$

147 (Eq. 2)

148

In the equation above, $v = n_1$ is the count of triplets congruent with the species tree topology (dominant topology) while n_2 and n_3 denote the counts of incongruent triplets (minority topologies). Note that this distribution does not distinguish the two-species case from the three-species case. Therefore, it is impossible to split a pair of species only represented by two samples but possible to split a species represented by a single sample from species with 2 or more samples.

155

156 A Bayesian model comparison framework

In the absence of prior knowledge of the species tree, the observer cannot know a priori which triplet is the triplet concordant with the species tree. Choosing the most frequently observed triplet and using its count as v in the above equation introduces a bias toward the three-species case and increases the rate of false positives (Supplementary figure

S1, available at Dryad: dryad. 3cb25). We therefore develop a Bayesian model comparisonframework to take the unknown species tree into account.

163 When the species tree is unknown, there are three models that conform to the three-164 species case described above. We call these three models, $T_B = \{\tau_{b1}, \tau_{b2}, \tau_{b3}\}$, each of which 165 is associated with one of three possible topologies of the underlying species tree. The 166 likelihood functions of the models in T_B are described by P_B in equation (2), with the 167 dominant triplet v matching n_1 , n_2 and n_3 for τ_{b1} , τ_{b2} and τ_{b3} , respectively.

We also consider three models for the case of a single species, following the scheme of Yang and Rannala (2014). We call the set of the three models, $T_W = \{\tau_{w1}, \tau_{w2}, \tau_{w3}\}$. Each

model in T_W is again associated with one of the three possible topologies and has its

171 counterpart in T_B (Yang and Rannala 2014). The likelihood functions are P_W in equation (1)

and they are identical across models.

173 With the six candidate models, the joint posterior probability of τ (model) and λ 174 (branch length) given triplet counts X = (n₁, n₂, n₃) is

175

$$P(\lambda, \tau | X) = \frac{P(X|\lambda, \tau)\pi(\lambda)\pi(\tau)}{\int \sum_{\tau \in \mathbf{T}_B \cup \mathbf{T}_W} P(X|\lambda, \tau)\pi(\lambda)\pi(\tau) \, d\lambda}$$

176 (Eq.3)

177 , where $\pi(\tau)$ and $\pi(\lambda)$ are posterior probabilities of τ and λ . We obtain the posterior probability 178 of τ by marginalizing the joint posterior by λ .

$$P(\tau|\mathbf{X}) = \frac{\int P(X|\lambda, \tau)\pi(\lambda)\pi(\tau)d\lambda}{\int \sum_{\tau \in \mathbf{T}_B \cup \mathbf{T}_W} P(X|\lambda, \tau)\pi(\lambda)\pi(\tau)\,d\lambda}$$

179 (Eq. 4)

To simplify the expression for the posterior, we now employ simple uniform priors, $\pi(\tau) = \frac{1}{6}$ and $\pi(\lambda) = \frac{1}{L} [0 \le \lambda \le L]$. We use a prior range up to *L*=5 throughout this study, which covers a realistic range of frequency of dominant triplets, $0.33 \le 1-2e^{-\lambda}/3 \le 0.996$. The posterior probability of the model with the uniform priors is,

$$P(\tau|\mathbf{X}) = \frac{\int_0^L P(X|\lambda,\tau)d\lambda}{\int_0^L \sum_{\tau \in \mathbf{T}_B \cup \mathbf{T}_W} P(X|\lambda,\tau)d\lambda}$$

184 (Eq.5)

185 The integration of $P(X|\lambda, \tau)$ over λ has a tractable analytical solution, therefore a reversible 186 jump MCMC is not required to characterize this posterior distribution. When τ is one of the 187 three models of T_W , the integration over λ is trivial.

$$\int_0^L P(X|\lambda,\tau\in T_W)d\lambda = L\cdot P_W(X)$$

188 When τ belongs to T_{*B*}, the integration of the likelihood function is represented by the incomplete beta 189 function. When the dominant triplet v is n₁,

$$\begin{split} \int_{0}^{L} P(X|\lambda,\tau\in T_{B})d\lambda &= \left(\frac{1}{2}\right)^{n_{2}+n_{3}} C \int_{\frac{2e^{-L}}{3}}^{\frac{2}{3}} (1-x)^{n_{1}} x^{n_{2}+n_{3}-1} dx \\ &= \left(\frac{1}{2}\right)^{n_{2}+n_{3}} C \left\{ \beta \left(\frac{2}{3}; n_{2}+n_{3}, n_{1}+1\right) - \beta \left(\frac{2e^{-L}}{3}; n_{2}+n_{3}, n_{1}+1\right) \right\} \end{split}$$

where $\beta(x; a, b)$ is the incomplete beta function and *C* is the multinomial coefficient in equation (2). Replacing n₁ with n₂ or n₃ gives solutions for v = n₂ or n₃.

The models in T_B are supporting the three-species delimitation, *B*, (that is, samples are from three distinct species); therefore the posterior probability of the delimitation *B* is a sum of the three posterior probabilities of the models in T_B .

$$P(B|\mathbf{X}) = \frac{\sum_{\tau \in \mathbf{T}_B} \int_0^L P(X|\lambda, \tau) d\lambda}{\sum_{\tau \in \mathbf{T}_B} \int_0^L P(X|\lambda, \tau) d\lambda + 3L \cdot P_W(X)}$$

195 (Eq.6)

, and the posterior probability of the single-species case delimitation, *W*, (samples are from asingle species) is,

$$P(W|\mathbf{X}) = \frac{3L \cdot P_W(X)}{\sum_{\tau \in \mathbf{T}_B} \int_0^L P(X|\lambda, \tau) d\lambda + 3L \cdot P_W(X)}$$

198 (Eq.7)

With a given hypothesis of delimitation, each of the $\binom{K}{3}$ triplets is assigned to one of the two categories defined above, i.e. a, b and c either belong to the same species or to multiple species. The overall posterior probability of a given delimitation for all K taxa is the product of the posterior probabilities of all triplet counts of two categories. For a set of triplet counts, *w*, which is assigned to delimitation *W*, and a set *b*, which is assigned to delimitation *B*, the log-posterior probability of a delimitation *D* is as follows.

205

$$\log P(\mathbf{D}|\mathbf{X}) = \sum_{(n_1, n_2, n_3) \in \boldsymbol{w}} \log P(W|n_1, n_2, n_3) + \sum_{(n_1, n_2, n_3) \in \boldsymbol{b}} \log P(B|n_1, n_2, n_3)$$

206 (Eq.8)

We use this quantity as the posterior probability score. Note that it is not a true posterior probability of delimitation since it ignores the mutual dependence of the parameters of the triplet distribution caused by overlapping membership of some triplets. However, the similar approximation of likelihood functions has been used successfully in statistical phylogenetic inference (Liu et al. 2010) and we test its performance by simulation here. 212

213 Finding the best delimitation model

The posterior probability score described above is used to find the optimal 214 delimitation from a set of delimitations of samples. The number of all possible delimitations 215 of K samples is represented by the Bell number, $\sum_{i=1}^{K} {K \choose i}$ (Bell 1934), where ${K \choose i}$ is a Stirling 216 number of second kind, defined as the number of all possible ways to split K items into i 217 218 groups. This number is common to partitioning problems and intractably large. An approach taken to reduce the number of delimitations considered is using a guide tree (Yang and 219 220 Rannala 2010), which gives a hierarchical structure of multiple delimitations. Different combinations of splitting and lumping of lineages on a given guide tree are searched to find 221 the best delimitation. Conventional search methods with the guide tree approach use either 222 reversible-jump Markov chain Monte Carlo for characterizing posterior probabilities of 223 competing delimitations (Yang and Rannala 2010, 2014) or heuristic search algorithms to 224 225 find the optimal combinations of splits and lump of lineages (O'Meara 2010; Satler et al. 2013). 226

We now consider only the problem of finding the best delimitation on a fixed guide 227 tree without tree rearrangement. The number of all possible delimitations under a given guide 228 tree with S tips is approximately $|1.5^{S}|$ (floor of 1.5^{S}) in the worst case and S in the best case 229 (Fujisawa and Barraclough 2013). The size of the space is reduced compared with the Bell 230 number, but still grows exponentially with the number of species in the worst case. We 231 developed a new dynamic programming algorithm to rapidly find the best combination of 232 lineages on the guide tree, taking advantage of the optimal substructure of the likelihood 233 model and the guide tree structure. 234

Given a delimitation under a guide tree, D, its posterior probability score can be decomposed into a sum of the scores of the delimitations of two subtrees descending from the root, $logP(D_L)$ and $logP(D_R)$, and a constant factor since equation (8) is additive.

238

$$\log P(D) = \log P(D_L) + \log P(D_R) + c$$

$$c = \sum_{(n_1, n_2, n_3) \in b_{root}} \log P(B|n_1, n_2, n_3)$$

239 (Eq.9)

where c is a constant representing a score for triplets crossing over 2 subtrees descending 240 from the root node of the guide tree. Since triplets are not shared between subtrees and their 241 posterior probabilities are independently calculated, the optimal solutions for each subtree 242 must be included in the global solution. Therefore, finding the global optimal solution can be 243 244 reduced to finding solutions to subtrees' delimitations, and iteratively solving and combining 245 them yields the global solution. An exception is the case where logP(D) is represented by the root of the guide tree; that is, all samples are from the same species. In the case of root 246 delimitation, equation (9) does not hold because the constant of the third term must be 247 represented by P(W|X) not P(B|X). So, the dynamic programming algorithm must compare 248 the "root" delimitation with the aggregated solution of subtrees in each step. This leads to the 249 algorithm described in Figure 1 and supplementary text S1 (Online Appendix). This algorithm 250 calculates the global optimal posterior probability score from a guide tree, and the best 251 delimitation was obtained by keeping the set of nodes producing the best score. The algorithm 252 reduces the number of likelihood calculations to twice the number of the nodes on the guide 253 254 tree.

We implemented the method in a program called "tr2" (Trinomial distribution of Triplets) for calculation of posterior probability scores for delimitation hypotheses and the search algorithm for the best delimitation given a guide tree. The program is implemented in Python and can run on any operating system (Distributed at

259 <u>https://bitbucket.org/tfujisawa/tr2-delimitation</u>, archived version for paper is available at

260 Dryad: dryad.3cb25).

261

262 *Simulations and case studies*

We used simulated and real gene trees to test the performance of the method. First, we performed coalescent simulations with species trees with 3 and 10 tips. Then, we analyzed a published data set of rattlesnakes with 29 individuals and a newly sequenced data set of 144 *Bacillus cereus* isolates.

267

Three-species simulations.— In order to test the performance of the delimitation 268 model, we first conducted a simple 3 species simulation and assessed the error rates of the 269 model. In this simulation, we assume gene trees are known without error. Gene genealogies 270 were simulated within a species tree with three tips and fixed branch lengths, T_1 =4000 and 271 T_0 =8000 generations (Figure 2a). The number of samples per species was set to 10, totalling 272 30 individual samples. The effective population sizes were set to $1/2*T_1$ to $8*T_1$ for all 273 species ($T_1=1/8Ne - 2Ne$ generations). Coalescent trees within the species tree were 274 simulated using SIMCOAL (Excoffier et al. 2000) assuming that one species represents one 275 population and populations merge on speciation events. Custom scripts were used to generate 276 input files for SIMCOAL from species trees. Twenty-five independent loci were simulated 277

100 times, which resulted in 2500 gene trees in total. The posterior probability for a global delimitation *W* (all individuals are from a single species) and the three alternative models representing correct delimitation (a), over-splitting (b) and under-splitting (c) (Figure 2) were calculated with increasing numbers of loci between 5 and 25 with step 5. Error rates, i.e., the frequency of choosing an incorrect model as the best model, were recorded for each iteration.

283

Ten-species simulations.— The ten-species simulation considers more realistic 284 conditions. Species trees with 10 tips were simulated under the Yule model with a constant 285 speciation rate. The total depth of species trees, T, was rescaled to 20000 generations, and the 286 effective population size (Ne) of species was set as 1, 1/2, 1/4, 1/8 and 1/16 times of T 287 (Ne=1250-20000). These parameter settings cover speciation rates and effective population 288 sizes observed in various eukaryotic groups (Coyne and Orr 2004; Charlesworth 2009) 289 290 including extreme cases of rapid radiations. Gene trees with 10 samples per species (100 total samples) for 40 independent loci were simulated using SIMCOAL and the custom scripts. 291 Simulations were replicated 100 times. 292

In the first simulation, hereafter simulation A, we assume that the topology of the guide tree and assignment of terminals to species groups is known. This simulation tests whether the method can correctly find the positions of nodes which define species from multiple competing combinations on a guide tree. The tr2 program was run with the species tree as a guide and simulated gene trees as inputs. The number of loci used ranged from 5 to 40.

In the second simulation, B, delimitation was conducted solely from sets of gene trees (species discovery approach). A consensus tree was built from gene trees from multiple loci using the rooted triple consensus (Ewing et al. 2008). Then, the consensus tree was used as

the guide tree in the delimitation step. This guide tree contains all possible hierarchical 302 delimitations, from each individual representing a separate species to all individuals 303 representing a single species. Polytomies on consensus trees were randomly resolved by the 304 "multi2di" function in the "ape" package (Paradis et al. 2004). In addition, we performed a 305 set of simulations to assess the effect of increasing numbers of loci and individual samples. 306 Gene trees were simulated within the same species trees as above with Ne=T/4, but the total 307 number of samples was reduced to 50 (5 per species) and the number of loci was doubled, 308 keeping the total sample size (number of loci X number of samples) constant. Delimitation 309 with tr2 was conducted in the same procedure as simulation B. 310

The third simulation, C, considers conditions where gene trees and species trees are 311 estimated from DNA sequences. Sequences were simulated along the branches of the gene 312 313 trees simulated above using Seq-Gen (Rambaut and Grassly 1997) assuming HKY+G model (Ts/Tv = 2.5 and α = 0.1) and 3% of overall genetic variations. These parameters were chosen 314 to be comparable to the case studies described in the next sections. Sequence length was set to 315 316 a constant length of 750bp. Gene trees were reconstructed from the simulated sequences using RAxML with a GTR+G model (Stamatakis 2014) and rooted by the "-I f" option of RAxML. 317 Guide trees were estimated from the reconstructed gene trees with the rooted triple consensus, 318 and delimitation was conducted with tr2. Under the parameter settings above, within-species 319 genetic variation of simulated sequences ranged from 0.3% to 1.4% depending on Ne, and 320 between-species variation was 3.0%. 321

The number of estimated species and the number of exact matches between estimated and true species were measured as the accuracy of delimitation. The elapsed time for each trial was also recorded. The numbers of non-monophyletic species were counted to measure the degree of incomplete lineage sorting. The effects of Ne, the number of loci, simulation type (A, B and C) and their second interaction terms on the proportion of exact matches were tested using GLM. For simulation B, the effect of the two sampling strategies was also tested.
Simulations and delimitations were run on a Linux personal computer with a 2.3 GHz Intel i5
quad-core processor and 4GB memory.

330

Case study one: Sistrurus Rattlesnakes.— Kubatko et al. (2011) sampled 18 nuclear 331 loci and one mtDNA locus of Sistrurus rattlesnakes. The data set of the nuclear loci included 332 58 phased sequences from 29 individuals of six known subspecies of S. catenatus/S. miliarius 333 and two outgroups. Kubatko et al. (2011) reported that one subspecies, S. catenatus catenatus, 334 exhibited signatures of a distinct species status while the other five subspecies did not show 335 significant evidence of independent species based on the monophyly-based test described by 336 Rosenberg (2007). We reanalyzed this data set. The gene trees and an alignment matrix of 18 337 nuclear loci were downloaded from TreeBase (accession:TB2:S11174). The trees were 338 339 randomly resolved with "multi2di". Then, a consensus tree was built using the rooted triple consensus from them, and the best delimitation was determined with the consensus as the 340 guide tree. A re-sampling procedure of loci was conducted by progressively adding single loci 341 342 in random order. Polytomies were randomly resolved in each iteration. The re-sampling was repeated 50 times to characterize the effect of increasing number of loci on the delimitation. 343 Genetic variation within subspecies was 0.2% and between species 2.2%. 344

345

346 *Case study two: Bacillus Multilocus Sequence Typing.*— We tested the applicability of 347 the tr2 method to bacterial species using a multilocus sequence typing (MLST) data set of the 348 *Bacillus cereus* complex. Multilocus sequence typing (MLST) is a typing scheme for bacterial 349 species/subspecies using a few (typically seven) loci (Maiden et al. 1998; Maiden 2006). It is 350 widely used in clinically relevant bacteria and occasionally in environmental prokaryotes to

delimit species (e.g. Papke et al. 2007). Although bacterial reproduction is largely clonal, in 351 many bacteria including Bacillus cereus, genetic exchange also occurs (Vos and Didelot 352 353 2009). If there was frequent gene exchange within a group of closely related individuals, but none between distantly related groups, this could lead to units equivalent to reproductively 354 355 isolated species in sexual eukaryotes (Didelot et al. 2011; Barraclough et al. 2012). The tr2 method should be able to delimit such a group as a putative species. However, in clonal 356 bacteria without any recombination, the delimitation method based on gene tree congruence 357 would delineate all individuals as separate species because the true genealogy of each locus 358 would be identical. Another complication is that horizontal transfer might occur rarely 359 between otherwise distinct species. This could introduce additional incongruence among loci 360 between otherwise separate species. We were interested to see how the method coped with a 361 prokaryotic clade that might display these complications. 362

Our sample comprised 144 isolates originally collected from evenly spaced quadrats 363 in the walled garden at Silwood Park for the study by Collier et al. (2005). In brief, freezer 364 365 isolates were regrown on *B. cereus* selective agar and DNA extracted using Chelex Instagene matrix method. The 7 house-keeping genes used for standard B. cereus MLST (Jolley et al. 366 2004) were PCR amplified and Sanger sequenced using primers and conditions at the MLST 367 database (http://pubmlst.org/bcereus/info/primers.shtml). Sequences were edited in Geneious 368 and trimmed to the lengths used at the MLST database. Full details are provided elsewhere 369 (Collier et al. 2005; Barraclough et al. in preparation), and sequences are available at Genbank 370 (Accession: KT806485-KT807462). 371

Alignment lengths of the MLST sequences ranged from 348 to 504bp, and there were 29 to 55 unique haplotypes at each locus (maximum of 55 for purH and minimum of 29 for glpF and gmk). The complete data matrix excluding missing loci contained 2806 bp from each of 114 isolates, which included 99 unique multilocus sequence types. Overall genetic

variation was 4.0%. Sequences from the seven loci were separately aligned with MUSCLE 376 3.8 (Edgar 2004). Gene genealogies of the seven loci were estimated using BEAST 1.80 377 378 (Drummond et al. 2012). Ten million generations of MCMC sampling were run with a GTR+G substitution model and the log-normal relax clock model (Drummond et al. 2006). 379 380 Twenty percent of the MCMC samples were discarded as burn-in. The convergence of the parameters was checked by effective sampling size using Tracer (Rambaut and Drummond 381 2007), and the maximum clade credibility trees (MCC trees) were extracted from the MCMC 382 runs using TreeAnnotator. 383

Two methods were used to obtain guide trees for delimitation of the Bacillus group. 384 First, a consensus tree was constructed using the rooted triple consensus from the MCC trees 385 of seven loci. Second, in order to account for the effects of horizontal transfer on the guide 386 387 tree estimation, we ran ClonalFrame (Didelot and Falush 2007) on the concatenated alignment. ClonalFrame estimates the most likely clonal genealogy by removing putative 388 horizontally transferred regions. An MCMC of ClonalFrame was run with 800 thousand 389 generations, and 50% of the chain was discarded as burn-in. Convergence of parameters was 390 examined by checking effective sample size using Tracer. The 50%-majority consensus from 391 the ClonalFrame MCMC was used as a second guide tree. Re-sampling of loci was 392 conducted 50 times using these two guide trees. To further account for the uncertainty of tree 393 building, 100 trees were sampled from the MCMC chain from BEAST for each locus and 394 from the chain of ClonalFrame, and delimitation was repeated with these 100 sets. The 395 frequency for each pair of samples to be grouped in the same species was recorded. (Sequence 396 alignments and trees are available at Dryad: dryad. 3cb25) 397

398

399 RESULTS

400 *Three-species simulations*

The overall false positive rate (FPR, rate of over-splitting) in the three-species simulations is 401 0.0 in all iterations with all numbers of loci between 5 and 25. False negative rates (FNR, rate 402 403 of under-splitting) decrease as the number of loci used increases (Figure 3). FNR of less than 30% were attained with only 5 loci when Ne was 2000 and 4000 (equivalent to T=2Ne and 404 Ne) whereas the FNRs reached 30% with 20 loci when Ne was 8000 (T=1/2Ne). With larger 405 Ne values, the decrease of FNRs was much slower, and the method was not able to correctly 406 delimit species within the range of loci used in the simulations when Ne was 16000 and 407 32000 (T=1/4Ne-1/8Ne). The average time required for one trial was 0.5 seconds. 408

409 *Ten-species simulations*

When true species trees are given as the guide tree, the method appeared to delimit species
consistently. The proportion of exact matches increased with the number of loci used (Figure
4, A), and the number of estimated species approached the true number of species, 10 (Figure
5, A). With low Ne value (Ne=1250), the median number of exact matches reached 10 when
25 or more loci were used. The increase in the number of exact matches slowed down with
larger Ne values, for example, when Ne≥5000, 40 loci were not enough to attain 100% exact
matches.

The accuracy was slightly reduced when the guide trees were estimated by the consensus method (Figure 4 and 5, B). However, the effect of simulation type was not significant (z = -0.33, p = 0.74 for simulation type, GLM with binomial errors) while Ne and the number of loci were highly significant (p << 0.001 for both Ne and the number of loci). In addition to the under-split observed in the simulation A, a few oversplits occurred especially when the number of loci was small. In 0.9% of trials, the method estimated more than 10 species. Overall accuracy still increased when more loci were added. When the gene trees

were estimated from the simulated sequences, the accuracy further decreased, especially when the number of loci was small (Figure 4, C). The accuracy was significantly lower than other simulation types (z = -4.42, p << 0.001, GLM with binomial error). Even more frequent oversplits were observed: the number of trials with >10 estimated species reached 2.0%. (Figure 5, C).

The time required for a delimitation process increased nearly linearly with the number of loci (Supplementary Figure S2). Median time ranged from 23 to 47 seconds for 10 tip guide trees and from 135 sec to 162 seconds for guide trees with 100 tips. Average proportions of non-monophyletic species were between 0.34 for Ne=1250 and 0.97 for Ne=20000 (Supplementary Figure S3), indicating non-monophyly is prevalent even for small Ne values. The accuracy of delimitation was significantly lower when fewer loci and more samples were used (*z*=-3.27, p=0.001, GLM with binomial error, Supplementary Figure S7).

436

437 *Rattlesnakes*

The method delimited 4 putative species of the Sistrurus rattlesnakes, including two 438 ingroup and two outgroup species (Figure 6 Left, Supplementary Figure S4). The two ingroup 439 species matched with the known taxonomic species, S. catenatus and S. miliarius. Random 440 resampling of loci indicated the number of estimated species does not saturate within the 441 range of loci used in this study (Figure 6 Left). With 18 loci, 28% of repeated delimitations 442 split S.catenatus into two groups: one group exclusively consisted of a subspecies 443 444 S.c.catenatus and another group consisted of S.c.edwardsii and S.c.tergeminus. Three subspecies of S.miliarius were always grouped together into a single species. 445

Delimitation using the seven MCC trees and rooted triple consensus tree resulted in 7 putative 448 species while the delimitation with ClonalFrame consensus resulted in 11 species. The 449 majority of nodes on the rooted triple consensus were unresolved (Supplementary Figure S5). 450 ClonalFrame robustly recovered three clades, two of which were unresolved in the rooted 451 triple consensus (Clade A, B and C in Figure 7). The difference between the two approaches 452 is consistent with horizontal transfer affecting topologies deeper in the tree; we mainly focus 453 on the result of delimitation using the ClonalFrame guide tree. Re-sampling of loci showed 454 that there was substantial variation in the number of estimated species (Figure 6 Right): the 455 sample of 7 loci might be too few for robust delimitation in this case. Repeated delimitations 456 run on 100 sets of MCMC tree samples exhibited 18 species that were consistently delimited 457 458 (Figure 7). While clades A and B were grouped into three or four large clusters, clade C was more frequently separated into small singleton species. Frequencies for isolates to be grouped 459 in species with other isolates within these clades were on average 61% and 40% for clade A 460 and B and 35% for clade C. 461

462 We estimated linkage disequilibrium (LD) within subsets of these groups to test for variation in recombination rate. Samples were taken from within the largest clusters in clade 463 A and B respectively, and randomly from within clade C, and LD of variable sites was 464 calculated for each group by the "LD" function of an R package "pegas" (Paradis 2010). The 465 test calculates the correlation between pairs of variable sites (Zaykin et al. 2008). There are 466 distinctive linkage patterns between and within the seven loci in the three groups 467 (Supplementary Figure S6). In clade A and B, strong to moderate LD within each locus and 468 LD between a few pairs of loci were observed, but LDs between loci were small (Median 469 within-locus $R^2=0.49$ and 0.16 and Median between-locus $R^2=0.07$ and 0.08 for clades A 470 and B, respectively). This is consistent with recombination among separate loci, but linkage 471

within loci. On the other hand, in clade C, there were moderate or high levels of LD between most loci (Median within-locus $R^2=0.38$ and between-locus $R^2=0.29$ for clade C), consistent with low rates of recombination even between loci.

475

476 DISCUSSION

477 Congruence between gene trees provides intuitive and readily tractable statistical models for multilocus species delimitation. In this paper, we developed a method to delimit species based 478 on topological congruence or incongruence of triplets quantified by two types of trinomial 479 distribution models. These models were derived from the multispecies coalescent framework 480 and can be used for robust delimitation of species from gene trees with incomplete lineage 481 sorting. The simulation studies confirmed that the method can consistently delimit species 482 without monophyly, and its performance increased with the number of loci and decreased 483 with larger effective population size relative to divergence time. 484

The accuracy of the method is slightly lower than the reported performance of 485 conventional multilocus delimitation methods (Camargo et al. 2012); more than 25 loci were 486 required to delimit with 95% success rate under the condition T=0.5*Ne (Figure 2) while 487 Camargo et al. (2012) reported 60 - 100 % success with 10 loci by conventional methods. The 488 advantage of tr2 appears to be its speed and applicability to large data. According to 489 Camargo et al. (2012), with a four species guide tree, SpedeSTEM (Ence and Carstens 2011) 490 ran in 30 seconds with 20 samples and BP&P (Yang and Rannala 2010) with 80 samples in 491 492 6.5hrs. The order of speed of the tr2 (30 seconds with 100 samples and 10 species with known gene trees) is comparable to the fastest conventional method. When a sequence alignment is 493 used, additional time for tree reconstruction is required (e.g. approximately 1 minute per locus 494 by RAxML), but the reconstruction - delimitation procedure can still scale to large data sets. 495

In addition, the dynamic programming algorithm finds the global solution on a given guide
tree while most heuristic optimizations do not guarantee it. The method was sufficiently
conservative to over-splitting, which is a favorable property for DNA-based species
delimitation methods (Carstens et al. 2013).

The simulation studies also showed that the accuracy of the guide tree is crucial for 500 accurate multilocus delimitation. It has been known that the incorrect assignment of samples 501 on guide trees results in oversplits of species in multilocus delimitation (Leaché and Fujita 502 2010; Zhang et al. 2014). The oversplits observed in the discovery approach (simulation B 503 and C) are likely to have resulted from the incorrect placement of samples on guide trees. 504 However, except for the excess of oversplits, the effect of unknown guide tree was minimal. 505 506 The number of exact matches was not significantly different between known and unknown 507 guide tree simulations, and even when DNA sequences were used, accuracy was comparable to the other simulations with a sufficient number of loci. It appears that, when the consensus 508 species tree estimation can resolve a particular node on a guide tree, tr2 does not erroneously 509 merge or split species on the node. This is a useful property since there are discrepancies 510 between the number of loci required for correct delimitation, guide tree estimation and initial 511 population assignment in the conventional delimitation procedures (Zhang et al. 2014). The 512 inaccurate estimate of guide tree and delimitation may be mediated simply by adding more 513 loci as the number of loci is not a major computational obstacle. 514

The delimitation results for *Sistrurus* rattlesnakes were partially consistent with the reported results in Kubatko et al. (2011). Though the two known taxonomic species, *S. catenatus* and *S. miliarius*, were consistently delimited as putative species, only about 30% of iterations supported the distinctiveness of the subspecies *S. c. catenatus*. Considering the number of loci necessary to delimit species in the simulation studies, 18 nuclear markers appear to be insufficient to fully delimit this group with the present method. The resampling

also indicates that polytomies are an important source of uncertainty on delimitations. The 521 two alternative outcomes with 18 loci resulted solely from the different resolutions of 522 polytomies. The lack of mutations and resulting polytomies do not positively mislead the 523 delimitation when the identical sequences are randomly inserted or polytomies are randomly 524 resolved. Nevertheless, simulations and case studies show the use of uninformative loci 525 compromise the power of species detection and introduce uncertainty. We used repeated 526 527 delimitations with randomly resolved gene trees and guide trees, and this approach was able to capture the level of uncertainty of gene tree reconstruction in the rattlesnakes. Resampling 528 trees from bootstrap trees or MCMC runs, as done in the Bacillus data set, is an alternative 529 way to handle the uncertainty. 530

The results of re-sampling analysis of *Bacillus* complex indicate more uncertainty in 531 532 their delimitation than the rattlesnakes. The reduced number of species observed on the rooted triple consensus may partly result from the unresolved guide tree due to horizontal transfer 533 between distantly related groups. However, distinctive patterns of bacterial diversification 534 535 were still observed. Clade A and B were consistently delimited into large groups while clade C mainly consisted of weakly connected singletons. Samples from these two categories 536 exhibited a contrasting pattern of linkage disequilibrium patterns. Especially, low LD between 537 loci observed in samples from the largest clusters detected in clades A and B indicates that 538 there is frequent gene exchange between members of those groups. Homologous 539 recombination creates local topological discordance on bacterial genomes (Didelot et al. 540 2010), and if the recombination events are localized only within closely related groups, the 541 mutually recombining groups can be detected by tr2 through genealogical discordance. The 542 clusters delimited in the clades A and B are likely to be such groups. Clade C has low 543 recombination rates and methods based on recombination and gene congruence are 544 inappropriate. It may still be possible to identify independently evolving groups in such clades 545

using alternative concepts and methods developed for clonal bacteria and asexuals (Cohan
2001; Barraclough et al. 2003). Clearly, the mixture of high and low-recombining lineages in
the *Bacillus* data adds complexities to species delimitation (which we will address in detail
elsewhere) and the number of loci may not be large enough to fully elucidate diversification
patterns. However, the result demonstrates the potential for detecting 'recombinationally
isolated' groups in prokaryotes.

The parameters to be considered for the computational complexity of delimitation are 552 the number of samples (K), the number of species (S) and the number of loci (N). The 553 dynamic programming algorithm introduced in this paper finds the best delimitation and 554 reduces the complexity of search through a guide tree to time scale linear to S, O(S), which 555 allows a thorough search of a guide tree. For example, using a guide tree that assigns every 556 557 individual into a distinct species has often been prohibitive with large samples, but, in our simulations, tr2 was able to process guide trees with 100 tips within 150 seconds. Combined 558 with good performance with respect to other parameters - cubic dependency of time on 559 overall sample size, $O(K^3)$ and linear for loci, O(N) - the method could be used to provide a 560 rapid search method through candidate delimitation hypotheses before applying more 561 statistically rigorous methods to large datasets. Current next generation sequencing projects 562 often target a large number of loci from relatively few individuals. The tr2 method is suitable 563 for this type of sampling design since the impact of increasing loci on computations is smaller 564 than increasing individuals. A simulation shows that higher accuracy is achieved with more 565 loci than with more individuals when total sample size (loci X individuals) is fixed 566 (Supplementary Figure S7). This demonstrates a potential use of current sequence 567 technologies for species delimitation though the optimal sampling strategy is yet to be 568 investigated. A final point is that the dynamic programming algorithm introduced in this 569 study may be applied for other optimal partitioning problems using hierarchical structure, 570

such as finding optimal partitioning of sequence alignments for phylogenetic inference (Li etal. 2008; Lanfear et al. 2012).

In this study, we did not consider possible violations of the assumption of the 573 574 multispecies coalescent model including gene flow between populations. Gene flow between sister species reduces the number of dominant triplets and increases two minority triplets 575 equally, and may compromise the accuracy of the method. Incorporating branch lengths into 576 the model may be required to tease apart the effects of gene flow and incomplete linage 577 sorting. Introgression events from distantly related groups may be detected as an increase in 578 one of two minority triplet counts. Indeed, deviation from equal counts of minority triplets is 579 used for tests of introgressive gene flow (Durand et al. 2011; Zwickl et al. 2014). Violation of 580 the model assumption of panmixia could be detected in a similar manner by extending the 581 582 trinomial distribution model used in this study to a three-rate model.

583 The method now uses estimated gene trees as inputs. In addition, it uses a guide tree estimated from the given gene trees or other independent methods. This procedure does not 584 take the uncertainty of gene tree and species tree inference into account. Also, most 585 586 computational time of the delimitation procedure was spent on the tree-building steps (a BEAST run on 1 locus of *Bacillus* took 2.5hrs while the tr2 ran in 2 minutes with > 100 587 samples). For guide tree inference, one possible solution would be to incorporate joint 588 589 inference of species tree and delimitation using triplets. Triplet- or quartet-based phylogenetic inference methods using known gene trees under the multispecies coalescent framework have 590 been developed and implementations to handle large datasets already exist (Liu et al. 2010; 591 Mirarab et al. 2014). The delimitation step based on the trinomial distributions could be easily 592 integrated into these procedures. Also, gene tree inference could be bypassed by directly 593 594 counting triplets estimated from 3 corresponding sequences and an outgroup as done in some

phylogenetic inference programs (DeGiorgio and Degnan 2010). Combining these methodscould potentially lead to a highly scalable joint estimation of species tree and delimitation.

In conclusion, we present a method for species delimitation from multilocus data that 597 can potentially scale to the kind of sample sizes that are currently only feasible for single-598 locus approaches. The method uses exact methods derived from the multispecies coalescent, 599 but by splitting the problem into triplets it circumvents the computational challenges. As it 600 becomes easier to sequence non-model genomes, and consequently to assay variable nuclear 601 markers across clades, we envisage a growth in the number of studies using standardized 602 multiple unlinked markers across entire clades, equivalent to current DNA barcoding sample 603 regimes. Our method is designed with these scenarios in mind to complement more intensive 604 605 methods.

606

607 AUTHOR CONTRIBUTIONS

TF devised the new methods, wrote the software, ran the analyses and wrote the manuscript.
AA generated the *Bacillus* sequence data. TGB provided advice on the methods, helped with
some analyses and wrote the manuscript.

611

612 ACKNOWLEDGEMENTS

We thank Richard Ellis for providing the frozen *Bacillus* samples and Yumi Moon and Kevin
Balbi for help with that project. TGB and AA were supported by BBSRC grant
BB/G004250/1.

617 REFERENCES

- Baird N.A., Etter P.D., Atwood T.S., Currey M.C., Shiver A.L., Lewis Z.A., Selker E.U.,
 Cresko W.A., Johnson E.A. 2008. Rapid SNP discovery and genetic mapping using
 sequenced RAD markers. PLoS One. 3:e3376.
- Barraclough T.G., Balbi K.J., Ellis R.J. 2012. Evolving Concepts of Bacterial Species. Evol.
 Biol. 39:148–157.
- Barraclough T.G., Birky JR. C.W., Burt A. 2003. Diversification in sexual and asexual
 organisms. Evolution (N. Y). 57:2166–2172.
- Bell E.T. 1934. Exponential Numbers. Am. Math. Mon. 41:411–419.

Camargo A., Morando M., Avila L.J., Sites J.W. 2012. Species delimitation with ABC and
other coalescent-based methods: a test of accuracy with simulations and an empirical
example with lizards of the liolaemus darwinii complex (Squamata: Liolaemidae).
Evolution (N. Y). 66:2834–2849.

- Evolution (N. Y). 60.2834 2849.
- Carstens B.C., Dewey T.A. 2010. Species delimitation using a combined coalescent and
 information-theoretic approach: an example from north american myotis bats. Syst. Biol.
 59:400–14.
- Carstens B.C., Pelletier T.A., Reid N.M., Satler J.D. 2013. How to fail at species delimitation.
 Mol. Ecol. 22:4369–83.
- Charlesworth B. 2009. Effective population size and patterns of molecular evolution and
 variation. Nat. Rev. Genet. 10:195–205.
- 637 Cohan F.M. 2001. Bacterial species and speciation. Syst. Biol. 50:513–524.

Collier F.A., Elliot S.L., Ellis R.J. 2005. Spatial variation in Bacillus thuringiensis/cereus
populations within the phyllosphere of broad-leaved dock (Rumex obtusifolius) and
surrounding habitats. FEMS Microbiol. Ecol. 54:417–425.

- 641 Coyne J.A., Orr H.A. 2004. Speciation. Sunderland, MA.: Sinauer Associates. Inc.
- DeGiorgio M., Degnan J.H. 2010. Fast and consistent estimation of species trees using
 supermatrix rooted triples. Mol. Biol. Evol. 27:552–69.
- Degnan J.H., Rosenberg N.A. 2006. Discordance of species trees with their most likely gene
 trees. PLoS Genet. 2:e68.

646 647	Degnan J.H., Rosenberg N.A. 2009. Gene tree discordance, phylogenetic inference and the multispecies coalescent. Trends Ecol. Evol. 24:332–40.
648 649 650	Didelot X., Bowden R., Street T., Golubchik T., Spencer C., McVean G., Sangal V., Anjum M.F., Achtman M., Falush D., Donnelly P. 2011. Recombination and population structure in salmonella enterica. PLoS Genet. 7.
651 652	Didelot X., Falush D. 2007. Inference of Bacterial Microevolution Using Multilocus Sequence Data. Genetics. 175:1251–1266.
653 654	Didelot X., Lawson D., Darling A., Falush D. 2010. Inference of homologous recombination in bacteria using whole-genome sequences. Genetics. 186:1435–49.
655 656	Drummond A.J., Suchard M.A., Xie D., Rambaut A. 2012. Bayesian phylogenetics with BEAUti and the BEAST 1.7. Mol. Biol. Evol. 29:1969–73.
657 658	Durand E.Y., Patterson N., Reich D., Slatkin M. 2011. Testing for Ancient Admixture between Closely Related Populations. Mol. Biol. Evol. 28:2239–2252.
659 660	Edgar R.C. 2004. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. BMC Bioinformatics. 5:113.
661 662	Edwards D.L., Knowles L.L. 2014. Species detection and individual assignment in species delimitation: can integrative data increase efficacy? Proc. Biol. Sci. 281:20132765.
663 664	Ence D.D., Carstens B.C. 2011. SpedeSTEM: a rapid and accurate method for species delimitation. Mol. Ecol. Resour. 11:473–80.
665 666	Ewing G.B., Ebersberger I., Schmidt H.A., von Haeseler A. 2008. Rooted triple consensus and anomalous gene trees. BMC Evol. Biol. 8:118.
667 668	Excoffier L., Novembre J., Schneider S. 2000. SIMCOAL : A General Coalescent Program for the Simulation of Molecular Data. J. Hered. 91:506–510.
669 670 671	Fujisawa T., Barraclough T.G. 2013. Delimiting species using single-locus data and the Generalized Mixed Yule Coalescent approach: a revised method and evaluation on simulated data sets. Syst. Biol. 62:707–24.
672 673	Fujita M.K., Leaché A.D., Burbrink F.T., McGuire J.A., Moritz C. 2012. Coalescent-based species delimitation in an integrative taxonomy. Trends Ecol. Evol. 27:480–8.
674 675	Hambäck P.A., Weingartner E., Ericson L., Fors L., Cassel-Lundhagen A., Stenberg J.A., Bergsten J. 2013. Bayesian species delimitation reveals generalist and specialist parasitic

676 677	wasps on Galerucella beetles (Chrysomelidae): sorting by herbivore or plant host. BMC Evol. Biol. 13:92.
678 679	Hebert P.D.N., Cywinska A., Ball S.L., DeWaard J.R. 2003. Biological identifications through DNA barcodes. Proc. R. Soc. B Biol. Sci. 270:313–321.
680 681	Jolley K.A., Chan MS., Maiden M.C.J. 2004. mlstdbNet - distributed multi-locus sequence typing (MLST) databases. BMC Bioinformatics. 5:86.
682 683 684	Koufopanou V., Burt A., Taylor J.W. 1997. Concordance of gene genealogies reveals reproductive isolation in the pathogenic fungus Coccidioides immitis. Proc. Natl. Acad. Sci. U. S. A. 94:5478–82.
685 686 687	Kubatko L.S., Gibbs H.L., Bloomquist E.W. 2011. Inferring species-level phylogenies and taxonomic distinctiveness using multilocus data in Sistrurus rattlesnakes. Syst. Biol. 60:393–409.
688 689 690	Lanfear R., Calcott B., Ho S.Y.W., Guindon S. 2012. PartitionFinder: Combined Selection of Partitioning Schemes and Substitution Models for Phylogenetic Analyses. Mol. Biol. Evol. 29:1695–1701.
691 692	Leaché A.D., Fujita M.K. 2010. Bayesian species delimitation in West African forest geckos (Hemidactylus fasciatus). Proc. Biol. Sci. 277:3071–7.
693 694	Lemmon A.R., Emme S. a, Lemmon E.M. 2012. Anchored hybrid enrichment for massively high-throughput phylogenomics. Syst. Biol. 61:727–44.
695 696	Li C., Lu G., Ortí G. 2008. Optimal data partitioning and a test case for ray-finned fishes (Actinopterygii) based on ten nuclear loci. Syst. Biol. 57:519–39.
697 698	Liu L., Yu L., Edwards S. V. 2010. A maximum pseudo-likelihood approach for estimating species trees under the coalescent model. BMC Evol. Biol. 10:302.
699 700 701 702	Maiden M.C.J., Bygraves J.A., Feil E., Morelli G., Russell J.E., Urwin R., Zhang Q., Zhou J., Zurth K., Caugant D.A., Feavers I.M., Achtman M., Spratt B.G. 1998. Multilocus sequence typing: A portable approach to the identification of clones within populations of pathogenic microorganisms. Proc. Natl. Acad. Sci. 95:3140–3145.
703 704	Maiden M.C.J. 2006. Multilocus sequence typing of bacteria. Annu. Rev. Microbiol. 60:561– 88.
705 706 707	Millanes A.M., Truong C., Westberg M., Diederich P., Wedin M. 2014. Host Switching Promotes Diversity in Host-Specialized Mycoparasitic Fungi: Uncoupled Evolution in the Biatoropsis-Usnea System. Evolution (N. Y).:1–18.

- Mirarab S., Reaz R., Bayzid M.S., Zimmermann T., Swenson M.S., Warnow T. 2014.
 ASTRAL: genome-scale coalescent-based species tree estimation. Bioinformatics.
 30:i541–8.
- O'Meara B.C. 2010. New Heuristic Methods for Joint Species Delimitation and Species Tree
 Inference. Syst. Biol. 59:59–73.
- Pamilo P., Nei M. 1988. Relationships between Gene Trees and Species Trees. Mol. Biol.
 Evol. 5:568–583.
- Papke R.T., Zhaxybayeva O., Feil E.J., Sommerfeld K., Muise D., Doolittle W.F. 2007.
 Searching for species in haloarchaea. Proc. Natl. Acad. Sci. U. S. A. 104:14092–14097.
- Paradis E., Claude J., Strimmer K. 2004. APE: Analyses of Phylogenetics and Evolution in R
 language. Bioinformatics. 20:289–290.
- Paradis E. 2010. pegas: an R package for population genetics with an integrated-modular
 approach. Bioinformatics. 26:419–420.
- Pons J., Barraclough T.G., Gomez-Zurita J., Cardoso A., Duran D.P., Hazell S., Kamoun S.,
 Sumlin W.D., Vogler A.P. 2006. Sequence-based species delimitation for the dna
 taxonomy of undescribed insects. Syst. Biol. 55:595–609.
- Puillandre N., Lambert A., Brouillet S., Achaz G. 2012. ABGD, Automatic Barcode Gap
 Discovery for primary species delimitation. Mol. Ecol. 21:1864–77.
- Rambaut A., Drummond A.J. 2007. Tracer v1.4. .
- Rambaut A., Grassly N.C. 1997. Seq-Gen: an application for the Monte Carlo simulation of
 DNA sequence evolution along phylogenetic trees. Comput. Appl. Biosci. 13:235–8.
- Rannala B., Yang Z. 2003. Bayes estimation of species divergence times and ancestral
 population sizes using DNA sequences from multiple loci. Genetics. 164:1645–1656.
- Rosenberg N.A. 2007. Statistical tests for taxonomic distinctiveness from observations of
 monophyly. Evolution (N. Y). 61:317–23.
- Satler J.D., Carstens B.C., Hedin M. 2013. Multilocus species delimitation in a complex of
 morphologically conserved trapdoor spiders (mygalomorphae, antrodiaetidae, aliatypus).
 Syst. Biol. 62:805–23.
- Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
 large phylogenies. Bioinformatics. 30:1312–3.

- Taylor J.W., Jacobson D.J., Kroken S., Kasuga T., Geiser D.M., Hibbett D.S., Fisher M.C.
 2000. Phylogenetic species recognition and species concepts in fungi. Fungal Genet. Biol.
 31:21–32.
- Vialle A., Feau N., Frey P., Bernier L., Hamelin R.C. 2013. Phylogenetic species recognition
 reveals host-specific lineages among poplar rust fungi. Mol. Phylogenet. Evol. 66:628–
 44.
- Vos M., Didelot X. 2009. A comparison of homologous recombination rates in bacteria and
 archaea. ISME J. 3:199–208.
- Yang Z., Rannala B. 2010. Bayesian species delimitation using multilocus sequence data.
 Proc. Natl. Acad. Sci. U. S. A. 107:9264–9269.
- Yang Z., Rannala B. 2014. Unguided Species Delimitation Using DNA Sequence Data from
 Multiple Loci. Mol. Biol. Evol. 31:3125–3135.
- Zaykin D. V, Pudovkin A., Weir B.S. 2008. Correlation-based inference for linkage
 disequilibrium with multiple alleles. Genetics. 180:533–45.
- Zhang C., Rannala B., Yang Z. 2014. Bayesian Species Delimitation Can Be Robust to
 Guide-Tree Inference Errors. Syst. Biol. 63:993–1004.
- Zhang J., Kapli P., Pavlidis P., Stamatakis A. 2013. A general species delimitation method
 with applications to phylogenetic placements. Bioinformatics. 29:2869–2876.
- Zwickl D.J., Stein J.C., Wing R.A., Ware D., Sanderson M.J. 2014. Disentangling
- methodological and biological sources of gene tree discordance on oryza (poaceae)
 chromosome 3. Syst. Biol. 63:645–59.

760 CAPTIONS OF FIGURES

Figure 1. An illustration of how the dynamic programming algorithm finds the optimal
delimitation. Below each node in the guide tree, two alternative delimitations are compared
(horizontal arrows) and the better one is chosen (dotted squares). The best delimitation below
one node is inserted into the comparison at the higher level successively to yield the final
optimal delimitation.

766

Figure 2. Schematic representations of alternative hypotheses of delimitation for the threespecies simulations: a) Correct hypothesis, b) Under-split and c) Over-split.

769

Figure 3. Relationships between false negative rate and the number of loci used for

delimitation in the three-species simulations that simulated different effective population sizes

within species relative to the divergence time between species.

773

Figure 4. Relationships between the number of exact matches and the number of loci used in
the 10-species simulations. A) Both guide trees and gene trees are known, B) Guide trees are
estimated but gene trees are known and C) Guide trees and gene trees are estimated from
DNA sequences.

778

Figure 5. Relationships between the number of estimated species and the number of loci usedin the 10-species simulations. A) Both guide trees and gene trees are known, B) Guide trees

are estimated but gene trees are known and C) Guide trees and gene trees are estimated fromDNA sequences.

783

Figure 6. The number of species estimated when randomly re-sampling loci in the empirical
data sets. Left) Rattlesnakes. Right) *Bacillus* complex.

786

Figure 7. Results of delimitation with 100 sets of gene trees and guide trees sampled from
MCMC runs. Trees from ClonalFrame MCMC were used as guide trees. Left) The 50%
majority consensus tree built with ClonalFrame . Right) The frequency that each pair of
isolates was grouped by tr2.

791

792 CAPTIONS OF SUPPLEMENTARY FIGURES

Figure S1. a) Frequency of false positives for three types of model comparison procedures. 793 The count of each topology was drawn from a trinomial distribution with equal rate, which 794 795 simulates samples from a single species, and model comparisons using AIC and the Bayesian modelling were conducted. AICs were calculated using the likelihood of eq.(1) for the single-796 species case and eq.(2) for the three-species case. Posterior probabilities were calculated with 797 798 eq.(6) and eq.(7). The numbers of trials where the three-species case had larger AICs or smaller posterior probabilities than the single-species case were recorded as false positives. 799 Abbreviations of models are: AIC MF: AIC with the most frequent topology as dominant 800 topology. AIC Random: AIC with randomly chosen triplet as dominant topology. Bayes: 801 Bayesian model comparison with six models. 802

803

804

Figure S2. Relationships between the required time for delimitation and the number of loci.

Figure S3. Average number of non-monophyletic species for each setting in the 10-speciessimulations.

809

Figure S4. A rooted triple consensus tree inferred from 18 nuclear loci of *Sistrurus* snakes.
Polytomies were randomly resolved. Bars in dark grey indicate species delimited by the tr2.
Bars in light grey show an alternative delimitation that appeared in the repeated delimitations.
The numbers on nodes represent the average difference of log posterior probability scores
between the null and alternative models defined by the nodes. A positive value indicates that
the delimitation B (separate species) is preferred over the delimitation W (same species).
Nodes indicated by asterisks are the most recent common ancestor nodes for species groups.

Figure S5. Results of delimitation of *Bacillus* with 100 sets of gene trees sampled from
MCMC runs using rooted triple consensus as guide trees. Left) The rooted triple consensus
built with seven MCC trees . Right) Frequency that each isolate is grouped by the tr2.

821

Figure S6. Distribution of linkage disequilibrium measured by R². Dotted grey lines indicate
the border of loci and each grid cell represents the positions of each locus, from left hand side,
glpF, gmk, ilv, pta, purH, pycA and tipD.

Figure S7. Relationship between the number of exact matches and total sample size (loci X
individual samples). Dark grey: trials with sample size=50. Light grey: trials with sample size
= 100.









type A B C









False posivite rates





elapsed time







Clade A



Clade B



Clade C



