

## Accepted Manuscript

Registering the evolutionary history in individual-based models of speciation

Carolina L.N. Costa, Flavia M.D. Marquitti, S. Ivan Perez, David M. Schneider, Marlon F. Ramos, Marcus A.M. de Aguiar



PII: S0378-4371(18)30711-8  
DOI: <https://doi.org/10.1016/j.physa.2018.05.150>  
Reference: PHYSA 19690

To appear in: *Physica A*

Received date: 10 June 2017  
Revised date: 10 April 2018

Please cite this article as: C.L.N. Costa, F.M.D. Marquitti, S.I. Perez, D.M. Schneider, M.F. Ramos, M.A.M. de Aguiar, Registering the evolutionary history in individual-based models of speciation, *Physica A* (2018), <https://doi.org/10.1016/j.physa.2018.05.150>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Reply Ms. Ref. No.: PHYSA-171117

Title: Constructing phylogenetic trees in individual based models

Physica A

### Highlights

- We provide a link between individual-based models and macroevolutionary theory.
- We show how to track ancestral relationships and speciation/extinction events in IBMs.
- Genealogies of individuals and phylogeny of species are drawn from these algorithms.
- We illustrate these algorithms using a spatially-explicit model of speciation.
- We compare trees based on historical information with trees inferred from genetic data.

# 1 Registering the evolutionary history in individual-based 2 models of speciation

3 Carolina L. N. Costa<sup>a,\*</sup>, Flavia M. D. Marquitti<sup>b</sup>, S. Ivan Perez<sup>c,b</sup>, David M.  
4 Schneider<sup>b</sup>, Marlon F. Ramos<sup>b</sup>, Marcus A.M. de Aguiar<sup>a,b</sup>

5 <sup>a</sup>*Instituto de Biologia, Universidade Estadual de Campinas, Unicamp, 13083-859,*  
6 *Campinas, SP, Brazil*

7 <sup>b</sup>*Instituto de Física 'Gleb Wataghin', Universidade Estadual de Campinas, Unicamp,*  
8 *13083-859, Campinas, SP, Brazil*

9 <sup>c</sup>*División Antropología, Museo de La Plata, Universidad Nacional de La Plata, Paseo del*  
10 *Bosque s/n, 1900 La Plata, Argentina*

---

## 11 Abstract

12 Understanding the emergence of biodiversity patterns in nature is a central  
13 problem in biology. Theoretical models of speciation have addressed this  
14 question in the macroecological scale, but little has been done to connect  
15 microevolutionary processes with macroevolutionary patterns. Knowledge of the  
16 evolutionary history allows the study of patterns underlying the processes being  
17 modeled, revealing their signatures and the role of speciation and extinction  
18 in shaping macroevolutionary patterns. In this paper we introduce two  
19 algorithms to record the evolutionary history of populations and species in  
20 individual-based models of speciation, from which genealogies and phylogenies  
21 can be constructed. The first algorithm relies on saving ancestor-descendant  
22 relationships, generating a matrix that contains the times to the most recent  
23 common ancestor between all pairs of individuals at every generation (the Most  
24 Recent Common Ancestor Time matrix, MRCAT). The second algorithm directly  
25 records all speciation and extinction events throughout the evolutionary  
26 process, generating a matrix with the true phylogeny of species (the Sequential  
27 Speciation and Extinction Events, SSEE). We illustrate the use of these algorithms  
28 in a spatially explicit individual-based model of speciation. We compare  
29 the trees generated via MRCAT and SSEE algorithms with trees inferred by  
30 methods that use only genetic distance between individuals of extant species,  
31 commonly used in empirical studies and applied here to simulated genetic data.  
32 Comparisons between trees are performed with metrics describing the overall  
33 topology, branch length distribution and imbalance degree. We observe that  
34 both MRCAT and distance-based trees differ from the true phylogeny, with the  
35 first being closer to the true tree than the second.

36 **Keywords:** genealogies of individuals, phylogenies of species, macroevolutionary  
37 patterns, distance-based trees, tree statistics

---

\*Corresponding author. Phone: +55-19-35215466

Email address: [lemes.carol@gmail.com](mailto:lemes.carol@gmail.com) (Carolina L. N. Costa)

## 38 1. Introduction

39 The origin of the patterns of diversity at macroecological scale is a central  
40 problem in biology [1–3]. In the last decades patterns such as geographical  
41 variation in species richness, species abundance distributions and species-area  
42 relationships, have been studied from empirical and theoretical perspectives  
43 [4–8]. Neutral models of speciation – where differences between individuals are  
44 irrelevant for their birth, death, and dispersal rates [3, 9] – have played a central  
45 role in understanding the patterns of diversity at the macroecological scale.  
46 With the help of computers, it became possible to test different hypothesis about  
47 the mechanisms of speciation, such as sympatric versus allopatric processes,  
48 assortative mating and the effect of number of genes [10–12].

49 Among the different theoretical approaches designed to quantitatively study  
50 speciation [3, 13], models that explicitly incorporate space have allowed the  
51 study of major macroecological patterns that could be compared with those ob-  
52 served in nature [2, 7, 14, 15]. However, these models have given little attention  
53 to the historical or evolutionary dimension of the origin of diversity, which is  
54 reflected in the macroevolutionary patterns described by phylogenetic trees [16–  
55 19]. Because of the increased interest in the role of microevolutionary processes  
56 on the resulting macroecological patterns, the extension of these approaches to  
57 include algorithms that track the branching or phylogenetic divergence process  
58 is a next fundamental step to further explore models of speciation using simu-  
59 lations [16, 20, 21]. Individual-based models (IBM) widely used in biology [22]  
60 have the advantage that can be easily extended to include this historical per-  
61 spective and to provide a record of the ancestor-descendant relationships among  
62 the simulated individuals and/or species. These relationships can be stored in  
63 matrices from which individual genealogies and species trees (i.e. phylogenies)  
64 may be directly obtained.

65 In this article we describe two algorithms that save historical information in  
66 individual-based models of speciation. The first algorithm focuses on genealogies  
67 and the quantity saved is the parenthood of each individual. With parenthood  
68 registered, the *time to the most recent common ancestor*, i.e., the number of  
69 generations needed to go backward to find a common ancestor of one individual  
70 with another individual of the population, can be easily calculated in terms  
71 of the common ancestor of the parents. These times are computed at every  
72 generation between all pairs of individuals and, at the end of the simulation, are  
73 saved in a matrix (the Most Recent Common Ancestor Time matrix - MRCAT).  
74 The second algorithm focuses on phylogenies and consists of directly records  
75 all speciation and extinction events (the Sequential Speciation and Extinction  
76 Events - SSEE) and set a matrix analogous to MRCAT but whose entries are  
77 species rather than individuals. The SSEE matrix contains the exact branching  
78 times in the simulated clade or community, including all extinct species. The  
79 MRCAT and SSEE matrices can be used to draw the exact branching sequence  
80 of the simulated individuals and species, respectively. These procedures differ  
81 from the inference methods based on phenotypic and genetic traits used to  
82 estimate phylogenies in natural studies, because in our model we are looking

83 for the branching process forward in time, while in usual approaches the same  
 84 process is looked backwards in time. In addition to the presentation of the  
 85 MRCAT and SSEE algorithms, we compare the trees they generate with those  
 86 obtained by usual distance-based methods of phylogenetic inference using only  
 87 genetic data from simulated individuals of the final community. Comparing  
 88 these inferred phylogenies with those generated by MRCAT or SSEE algorithms  
 89 might offer a practical way to evaluate the reliability of the estimated trees to  
 90 recover natural macroevolutionary patterns.

91 The paper is organized as follows: in section 2 we describe the algorithms  
 92 to record ancestor-descendant relationships (MRCAT, subsection 2.1) and spe-  
 93 ciation/extinction events (SSEE, subsection 2.2). In subsection 2.3 we compare  
 94 the true phylogenetic tree obtained from the SSEE algorithm with genealogies  
 95 of individuals obtained from the MRCAT algorithm considering only one indi-  
 96 vidual per species. In section 3 we discuss the applications of the algorithms  
 97 proposed in section 2. First, we present an individual-based model of specia-  
 98 tion proposed in [2] in which the algorithms regarding the ancestor-descendant  
 99 relationships and the branching process were incorporated (subsection 3.1). We  
 100 emphasize that the algorithms are quite general and could be implemented in  
 101 most IBM's. Next, we briefly describe the Unweighted Paired Group Method  
 102 with Arithmetic mean (UPGMA) [23], the Neighbor Joining (NJ) [24] and the  
 103 Minimum Evolution (ME) [25] methods, which are based on genetic distances  
 104 calculated directly from one individual of each species present in the last gen-  
 105 eration of the simulation (subsection 3.2). While closer to what empiricists do,  
 106 the phylogenies derived from these methods are further from the true phylogeny  
 107 generated by the SSEE algorithm than is the phylogeny based on the MRCAT  
 108 algorithm presented here. We end this section presenting the statistical mea-  
 109 surements used to compare phylogenies obtained from algorithms proposed here  
 110 with those estimated by distance-based methods (subsection 3.3). The goal is  
 111 to show that the accuracy of some methods usually employed when the only  
 112 information available is the data of individuals collected from nature can be  
 113 evaluated with the help of models. In section 4 we present the results regarding  
 114 the output of simulations and the comparisons of phylogeny summary statistics.  
 115 Finally, section 5 was devoted to discussion and section 6 to conclusions.

## 116 2. Registering the history of individuals and species

117 In this section we describe two algorithms to record historical information  
 118 during the evolution of a population. The first algorithm records genealogical  
 119 relationships between all pairs of individuals at every generation. The second,  
 120 in turn, registers all the speciation and extinction events that occur along the  
 121 evolutionary history. These algorithms are general enough to be applied to most  
 122 individual-based models of speciation.

### 123 2.1. Ancestor-descendant relationships among individuals - MRCAT

124 In this subsection we show how the time to the most recent common ancestor  
 125 between all pairs of individuals can be obtained by keeping track of parental re-

Individuals at generation $t + 1$	Parent at generation $t$
1	$P(1) = 4$
2	$P(2) = 8$
3	$P(3) = 1$
4	$P(4) = 4$
$\dots$	$\dots$
$N_{t+1}$	$P(N_{t+1}) = 15$

Table 1: List of individuals ( $i$ ) at generation  $t + 1$  and their respective parents ( $P(i)$ ) at generation  $t$  in an asexual model. This information is necessary to construct the MRCAT matrix. Parents of each individual must be recorded to track the most recent common ancestor between individuals at the end of a simulation. Note that individuals at generation  $t$  are not the same individuals at generation  $t + 1$  (discrete generations).

126 relationships at every generation. We also show how this information can be used  
 127 to draw the genealogy of individuals of the last simulated generation. We dis-  
 128 tinguish between asexual and sexual models because of the technical differences  
 129 in tracking only one or two parents.

### 130 2.1.1. Asexual models

131 Consider a population of  $N_t$  asexual individuals at generation  $t$ . The pop-  
 132 ulation at the next generation,  $t + 1$ , will be comprised of offspring of these  
 133 individuals and the parent of individual  $i$  will be denoted  $P(i)$ .

134 An example is shown in Table 1, where  $P(1) = 4$ ,  $P(2) = 8$ ,  $P(3) = 1$ , etc.  
 135 The MRCAT between individuals  $i$  and  $j$  is

$$T_{t+1}(i, j) = T_t(P(i), P(j)) + 1. \quad (1)$$

136 which is simply the time to the most recent common ancestor between the  
 137 parents plus one, since a generation has passed [26]. As examples

$$T_{t+1}(1, 2) = T_t(4, 8) + 1$$

138 and

$$T_{t+1}(1, 4) = T_t(4, 4) + 1 = 1.$$

139 since in this last case they have the same parent. Starting from  $T_0(i, j) = 1$  if  
 140  $i \neq j$  and noting that  $T_t(i, i) = 0$  at all times the rule (1) allows one to compute  
 141 the MRCAT matrix for any number of generations. The matrix  $T$  is stored only  
 142 for two times, the past and the present generation, so that the memory cost  
 143 does not depend on time, only on the (square) size of population. A schematic  
 144 view of the algorithm is shown in Fig. 1, where the genealogical relationships  
 145 between 9 individuals originated from a single ancestor is represented. In this  
 146 example the total population size is kept fixed, so that the full MRCAT matrix  
 147 is always  $9 \times 9$ . The phylogeny of the community can be drawn by selecting one  
 148 individual per species at each moment in time. The corresponding matrices at  
 149  $t = 3$  and  $t = 6$  are given by

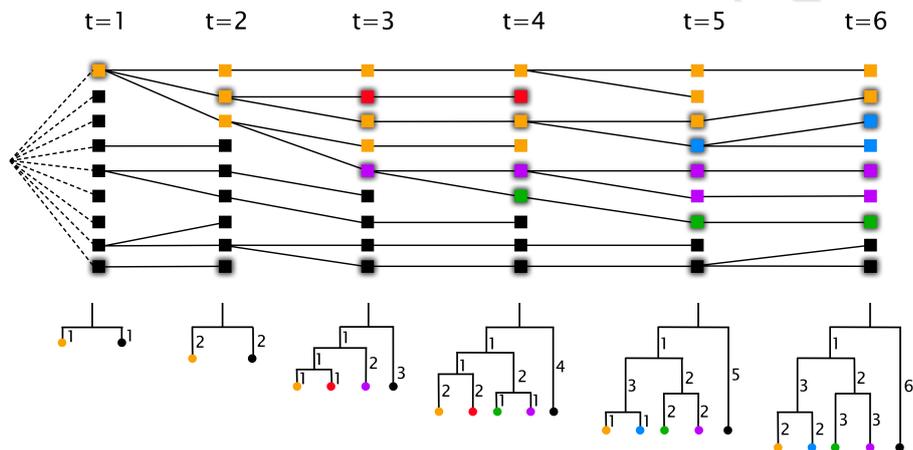


Figure 1: Illustration of ancestor-descendant relationships for an asexual population with constant size  $N = 9$  implemented with MRCAT algorithm. Each square is an individual and colors represent different species. Phylogenetic trees are constructed by selecting one individual per species (shaded squares).

$$T_3 = \begin{pmatrix} 0 & 1 & 2 & 3 \\ 1 & 0 & 2 & 3 \\ 2 & 2 & 0 & 3 \\ 3 & 3 & 3 & 0 \end{pmatrix}; \quad T_6 = \begin{pmatrix} 0 & 2 & 5 & 5 & 6 \\ 2 & 0 & 5 & 5 & 6 \\ 5 & 5 & 0 & 3 & 6 \\ 5 & 5 & 3 & 0 & 6 \\ 6 & 6 & 6 & 6 & 0 \end{pmatrix}. \quad (2)$$

150 where the selected individuals are shown in shaded colors (from top to bottom)  
151 at the corresponding times.

### 152 2.1.2. Sexual models

153 The generation of MRCAT matrices in sexual models is slightly different,  
154 since each individual  $i$  has two parents, a mother  $P_1(i)$  and a father  $P_2(i)$ . Con-  
155 sider as an example a population which has 4 females and 3 males in generation  
156  $t$  and gives rise to 5 females and 3 males in generation  $t + 1$  (Table 2). Notice  
157 that not only the total number of individuals but also the number of males and  
158 females may vary over generations. As the model is sexual, both maternal and  
159 paternal lineages can be followed in the simulations, allowing the generation of  
160 two different MRCAT matrices and their corresponding trees. A third option is  
161 not tracking lineages by sex, but record the most recent common ancestor tak-  
162 ing into account both parents, which is the only option if the model considers  
163 hermaphroditic individuals.

164  
165 – *Maternal and paternal lineages.* The maternal lineage of individuals is  
166 obtained by computing the time to the most recent common ancestor of their

Individuals at generation $t + 1$	Mother at generation $t$	Father at generation $t$
Females		
1	$P_1(1) = 4$	$P_2(1) = 6$
2	$P_1(2) = 3$	$P_2(2) = 7$
3	$P_1(3) = 1$	$P_2(3) = 7$
4	$P_1(4) = 4$	$P_2(4) = 5$
5	$P_1(5) = 2$	$P_2(5) = 6$
Males		
6	$P_1(6) = 1$	$P_2(6) = 5$
7	$P_1(7) = 3$	$P_2(7) = 5$
8	$P_1(8) = 3$	$P_2(8) = 7$

Table 2: List of individuals ( $i$ ) at generation  $t+1$  and their respective parents ( $P_1(i) = \text{mother}$  and  $P_2(i) = \text{father}$ ) at generation  $t$  in a sexual model. In this case each individual has two parents,  $P_1$  and  $P_2$ . Notice that the couple 3 and 7 at generation  $t$  had two offspring, the individuals 2 and 8 at generation  $t+1$ , while other couples had only one offspring. Additionally, notice that there were 4 females and 3 males at generation  $t$ , while there are 5 females and 3 males at generation  $t + 1$ .

167 corresponding mothers:

$$T_{t+1}^M(i, j) = T_t^M(P_1(i), P_1(j)) + 1 \quad (3)$$

168 with  $T_0^M(i, j) = 1$  if  $i \neq j$  and  $T_t^M(i, i) = 0$ . Similarly, the paternal lineage is  
169 computed with

$$T_{t+1}^F(i, j) = T_t^F(P_2(i), P_2(j)) + 1 \quad (4)$$

170 with  $T_0^F(i, j) = 1$  if  $i \neq j$  and  $T_t^F(i, i) = 0$ . Both  $T^M$  and  $T^F$  are computed for  
171 all individuals, females and males.

172

173 – *Lineages of hermaphroditic individuals.* Many simulations consider, for  
174 simplicity, hermaphroditic individuals. In this case, the separation into maternal  
175 and paternal lineages does not make sense and the definition of the MRCAT  
176 matrix is

$$T_{t+1}(i, j) = \min_{\{k,l\}} \{T_t(P_k(i), P_l(j))\} + 1 \quad (5)$$

177 with  $k, l = \{1, 2\}$ ,  $T_0(i, j) = 1$  and  $T_t(i, i) = 0$ . This considers, literally, the most  
178 recent common ancestor of  $i$  and  $j$ , taking all parental combinations into ac-  
179 count. The same definition is applied to sexual models with sex separation when  
180 the recorded genealogy does not separate the maternal and paternal lineages. In  
181 the case of hermaphroditic model the MRCAT matrix does not determine the  
182 tree uniquely. A detailed example of this situation is described in Supporting  
183 Information, section I.

### 184 2.1.3. Drawing genealogies from MRCAT matrices

185 At the end of the simulated evolutionary process the MRCAT matrix con-  
186 tains the time to the most recent common ancestor between every pair of in-  
187 dividuals of the extant population and this information can be used to draw

188 genealogical trees. Drawing the tree from the MRCAT matrix consists in join-  
 189 ing individuals into groups according to their most recent common ancestral  
 190 (Fig. 1). The tree starts with  $N$  units (the extant individuals) and at each  
 191 step of the process two of these units are joined together to form a group, so  
 192 that the number of units decreases by 1. Next, the time to the most recent  
 193 common ancestral between the newly formed group and the other units of the  
 194 tree (previously formed groups or extant individuals) are recalculated with a  
 195 so called *clustering method*. Once the times have been recalculated, the pair of  
 196 units with the least time is joined into a new group. The process ends when  
 197 a single unit is left, the root of the tree. As discussed in the SI, section I, a  
 198 unique tree is generated independently of the clustering method for asexual,  
 199 maternal or paternal lineages. For hermaphroditic populations or for sex sep-  
 200 aration but with the MRCA taking into account both parents that is not the  
 201 case. In these situations more than one tree can be constructed from the same  
 202 MRCAT matrix using different clustering procedures. In all cases the tips (or  
 203 leaves) of the tree represent extant individuals whereas internal nodes represent  
 204 the most recent common ancestor between a pair of individuals. Branch length  
 205 denote the time in generations between an ancestor and its descendants (see,  
 206 for instance, Fig. S1 in the SI). More information about the drawing of trees is  
 207 available in Supporting Information, section II.

## 208 2.2. Recording all speciation and extinction events - SSEE

209 The algorithm described in subsection 2.1 records the ancestor-descendant  
 210 relationships between all pairs of individuals in the population at a given point  
 211 in time. This allows the drawing of entire genealogies. However, information  
 212 about individuals that died without leaving descendants or species that went  
 213 extinct is totally lost. In this subsection we describe an algorithm that allows  
 214 the construction of the true phylogenetic tree, retaining information about all  
 species that ever existed during the evolution (Fig. 2).

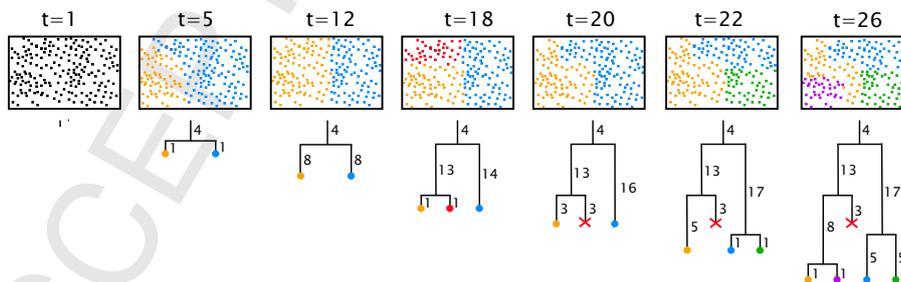


Figure 2: Illustration of speciation and extinction events implemented with SSEE algorithm and the corresponding phylogenetic trees exhibiting the complete history. Colored squares represent individuals of different species, and colored circles in phylogenies represent each species, with numbers denoting the time to speciation and extinction events.

215 We will use a new matrix  $S_t$  (the SSEE matrix) such that  $S_t(i, j)$  is the time  
 216 when species  $i$  and  $j$  branched off a common ancestral species. Species that  
 217

218 go extinct will be kept in the matrix but will be assigned a label to distinguish  
 219 them from living (extant) species. This label will be stored in a *extinction vector*  
 220  $E_t$  such that  $E_t(i) = 0$  indicates a living species at time  $t$  and  $E_t(i) = \tau \neq 0$   
 221 indicates the moment  $\tau$  when the species disappeared.

222 The algorithm is as follows: consider the hypothetical sequence of speciation  
 223 and extinction events displayed in Fig. 2. At time  $t=18$  there are three species  
 224 that we denote as Orange(18), Red(18) and Blue(18) and the corresponding S  
 225 matrix and E vector are

$$S_{18} = \begin{pmatrix} 0 & 1 & 14 \\ 1 & 0 & 14 \\ 14 & 14 & 0 \end{pmatrix}; \quad E_{18} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}. \quad (6)$$

226 Two generations later, at  $t = 20$ , one finds only two species, Orange(20) and  
 227 Blue(20). Notice that names (and colours) are arbitrary and to determine the  
 228 relation between these species and the ones at the previous time step we need  
 229 to look at the parents of individuals in each species. Suppose, as illustrated in  
 230 the figure, that we find that the parents of individuals in Orange(20) belonged  
 231 to species Orange(18). In this case we draw a link between Orange(18) and  
 232 Orange(20) and mark Orange(18) as a species that survived that time step, i.e.,  
 233 we set  $E_{20}(1) = 0$ . Similarly Blue(20) links with Blue(18) and  $E_{20}(2) = 0$ .  
 234 Looking at the previous generation we notice that species Red(18) did not leave  
 235 any descendant species, i.e., it went extinct. In order to keep track of it we  
 236 create a virtual species Red(20) and set  $E_{20}(3) = 20$  as a mark that it is no  
 237 longer a living species and went extinct at time 20. The SSEE and E vector at  
 238 time 20 become

$$S_{20} = \begin{pmatrix} 0 & 16 & 3 \\ 16 & 0 & 16 \\ 3 & 16 & 0 \end{pmatrix}; \quad E_{20} = \begin{pmatrix} 0 \\ 0 \\ 20 \end{pmatrix}. \quad (7)$$

239 Extinct species are, therefore, treated as species that will never again spe-  
 240 ciate, but will be kept in the matrix. When drawing the corresponding tree  
 241 its branch will stop at the value  $E(i)$ . Proceeding in this way, with the living  
 242 species always filling the first part of the matrix, followed by copies of extinct  
 243 species, we can draw the complete phylogeny and study extinction dynamics as  
 244 well. At time  $t = 26$  the SSEE matrix and extinction vector E are

$$S_{26} = \begin{pmatrix} 0 & 1 & 22 & 22 & 9 \\ 1 & 0 & 22 & 22 & 9 \\ 22 & 22 & 0 & 5 & 22 \\ 22 & 22 & 5 & 0 & 22 \\ 9 & 9 & 22 & 22 & 0 \end{pmatrix}; \quad E_{26} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 20 \end{pmatrix}. \quad (8)$$

245 One important case occurs when two species merge into a single species  
 246 (speciation reversal). This might happen, for instance, when two species that  
 247 have just become reproductively isolated are able to breed again because of a  
 248 mutation. The resulting merged species will have individuals with parents in

249 both ancestral species and we need to define which one “survived” and which  
 250 went extinct. Although this is just a matter of labeling the species, we call the  
 251 surviving species the one with most parents in the previous generation.

252 The drawing of species phylogenies for SSEE matrices is almost identical  
 253 to that for MRCAT matrices. The only differences are that internal nodes  
 254 represent speciation events, not the time to MRCA, and branches associated  
 255 to extinct species should not be drawn all the way down to present time, but  
 256 should stop at the extinction time recorded in the vector  $E$ . As in the MRCAT  
 257 case of separation of lineages by sex, a unique tree is generated independently  
 258 of the clustering procedure chosen, due to the exact times of speciation and  
 259 extinction recorded in simulations based on this algorithm.

### 260 *2.3. Phylogenies generated by ancestor-descendant relationships (MRCAT) ver-* 261 *sus trees from speciation and extinction events (SSEE)*

262 At the end of a simulation the MRCAT matrix contains the exact time to  
 263 the most recent common ancestor between every pair of individuals in the pop-  
 264 ulation. The SSEE matrix contains the equivalent information at the species  
 265 level, including extinct species. Both matrices can be used to draw phylogenetic  
 266 trees. To draw a phylogeny of species considering the ancestor-descendant rela-  
 267 tionships between individuals we can use the MRCAT matrix with the following  
 268 reasoning: if  $N_S$  species exist at time  $t$  and  $ind(i, j)$  is the  $j$ -th individual of the  
 269  $i$ -th species, a  $N_S \times N_S$  sub-matrix of the full MRCAT matrix can be generated  
 270 considering only one individual per species (Fig. 1); a simple choice is to take  
 271  $ind(i, 1)$  for  $i = 1, 2, \dots, N_S$  so that  $T_{i,j}^{phy} \equiv T_{ind(1,i),ind(1,j)}$ .

272 The tree drawn from the SSEE algorithm is the true phylogeny of species,  
 273 because it records the exact speciation and extinction events, representing the  
 274 actual branching process. On the other hand, the phylogeny of species drawn  
 275 from the MRCAT algorithm is different, although similar, from the true phy-  
 276 logeny, because the time to the most recent common ancestor between individ-  
 277 uals of different species is only an approximation to the speciation time, since  
 278 speciation can happen several generations later. Figure 3 illustrates this situ-  
 279 ation: if a population splits into three species in two closely spaced speciation  
 280 events, it might happen that the first group to speciate, species A in the figure,  
 281 has a more recent common ancestor with the subgroup B than B with C. During  
 282 the time when B and C still form a single species reproduction between their in-  
 283 dividuals might not happen for a while until they split, preserving the long time  
 284 ancestry. This is more likely to happen in populations with a spatial structure  
 285 when individuals belonging to the two subpopulations occupy different areas.

## 286 **3. Applications of MRCAT and SSEE algorithms to an individual-** 287 **based model**

### 288 *3.1. The speciation model*

289 The model considered here to exemplify the MRCAT and SSEE algorithms  
 290 is an extension of the speciation model introduced in [2] and adapted in [27] to

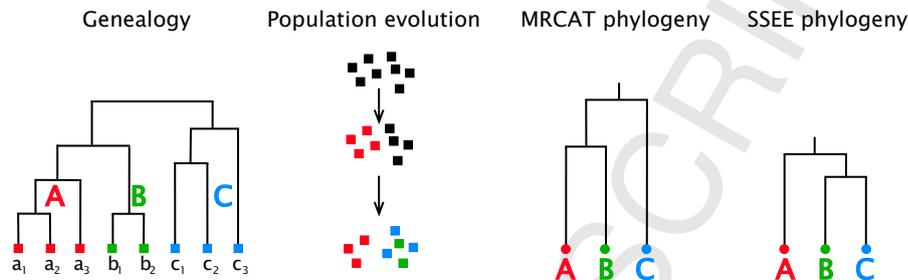


Figure 3: Illustration of a genealogy recorded with MRCAT and the corresponding population evolution. The phylogenies constructed via MRCAT and SSEE differ in this case because, although individuals from species A and B have a more recent common ancestor than with individuals in C, species A split first, followed by the separation of B and C.

291 characterize individuals with separated sexes (males and females). The model  
 292 has already been studied in terms of speciation rates, species-area relationships  
 293 and species abundance distributions. Here we are adding the historical informa-  
 294 tion generated by MRCAT and SSEE algorithms, i.e., recording the parenthood  
 295 of individuals from one generation to another (genealogy) as well as the pattern  
 296 and time of the speciation and extinction events (phylogeny or time tree).

297 The model describes a population of  $N$  haploid individuals that are gene-  
 298 tically identical at the beginning of the simulation and are randomly distributed  
 299 in a  $L \times L$  spatial lattice with periodic boundary conditions. More than one  
 300 individual is allowed in each site of the lattice, but because the density of the  
 301 population is low, this seldom occurs. The genome of each individual is repre-  
 302 sented by a sequence of  $B$  binary *loci*, with state 0 or 1, where each *locus* plays  
 303 the role of an independent biallelic gene. Individuals also carry one separate  
 304 label that specify their sex, male or female. The evolution of the population  
 305 involves the combined influence of sexual reproduction, mutation and dispersal  
 306 [2].

307 The reproduction trial starts with individual 1 and goes to individual  $N$ , so  
 308 that all individuals of the population have a chance to reproduce. The individ-  
 309 ual selected for reproduction, the *focal individual*, searches for potential mates  
 310 in its *mating range*, a circular area of radius  $S$  centered on its spatial location.  
 311 The focal individual can only reproduce with those within its mating range and  
 312 if they are genetically compatible, i.e., if the genetic distance between them is  
 313 below a particular threshold  $G$ . Among the compatible individuals within its  
 314 mating range one of the opposite sex is randomly chosen as mating partner.  
 315 Individuals whose genetic distance is larger than  $G$  are considered reproduc-  
 316 tively isolated (threshold effect [3]). Genetic distances between individuals are  
 317 calculated as the Hamming distance [28] between their genetic sequences, i.e.,  
 318 the number of *loci* at which the corresponding alleles are different.

319 Once the focal individual finds a compatible mate of the opposite sex, repro-  
 320 duction proceeds with the combination of their genetic materials to produce the  
 321 offspring genome, with each *locus* having an equal probability of being transmit-

322 ted from mother or father. After combination of parental genomes, each *locus*  
323 in the offspring genome can mutate with probability  $\mu$ . Finally, the offspring  
324 replaces the focal reproducing individual. In each reproductive event only one  
325 descendant is generated. The offspring is then dispersed with probability  $D$  to  
326 one of the 20 nearest sites (radius approximately equal to  $\sqrt{5} \approx 2.24$ ) around  
327 the expiring focal parent. Conversely, with probability  $1 - D$  the offspring will  
328 be placed exactly in the same site of its focal expiring parent. Hence, close to  
329 the location of every individual of the previous generation there will be an indi-  
330 vidual in the present generation, keeping the spatial distribution homogeneous.  
331 There is a probability  $Q$  that the focal individual will die without reproducing.  
332 In this case a neighbor is randomly selected from its mating range to reproduce  
333 in its place, so that the population size remains constant.

334 Evolution proceeds in non-overlapping discrete generations such that the  
335 entire population is replaced by offspring. Species are defined as groups of in-  
336 dividuals connected by gene flow, so that any pair of individuals belonging to  
337 different species are reproductively isolated (genetic distance greater than  $G$ ).  
338 However, two individuals belonging to the same species can also be reproduc-  
339 tively isolated, as long as they can exchange genes indirectly through other  
340 individuals of the species. This model is considered neutral because individuals  
341 choose their mates randomly from a mating range, independent of their genetic  
342 composition except for the genetic threshold of reproductive compatibility, so  
343 differences between individuals are irrelevant for their birth, death, and dispersal  
344 rates [3, 9].

### 345 3.2. Phylogenies based on genetic distances

346 As we have described in the previous subsection, the genome of all individ-  
347 uals are identical at the beginning of the simulation but mutations introduce  
348 differences and after many generations the population will display a distribu-  
349 tion of genomes. Genetic distances can, therefore, be calculated between pairs  
350 of individuals and be used as a proxy for ancestry, such that the larger the ge-  
351 netic distance between two individuals the farther back should be their common  
352 ancestor. In order to estimate phylogenies by genetic distance, we selected the  
353 same individuals per species that were used to draw the phylogeny via MRCAT  
354 and computed a matrix of genetic distances. This process mimics the sampling  
355 of individuals from a real population and the comparison of their DNA's as a  
356 measure of ancestry.

357 From the genetic distance matrix, we estimated trees from three distance-  
358 based methods. Firstly, we used the UPGMA hierarchical clustering method  
359 [23]. In this algorithm two groups of species are clustered based on the average  
360 distance between all members of the groups. This method assumes a constant  
361 rate of change, generating ultrametric trees in which distances from the root  
362 to all tips are equal. Secondly, we used the NJ method [24] of phylogenetic  
363 inference. In this method the procedure is to find pairs of neighbors in which  
364 the total branch length at each stage of the clustering is minimal, starting with a  
365 starlike tree. Finally, we used the ME method [25], which assumes that the true  
366 phylogeny is probably the one with the smallest sum of branch lengths, as in the

367 NJ method. The difference is that in the ME method a NJ tree is constructed  
 368 first and next tree topologies close to this NJ tree are estimated by certain  
 369 criteria, with all these trees being examined and the tree with the small sum of  
 370 branch lengths being chosen. We used the function `hclust` of the `stats` package  
 371 in R [29] to estimate ultrametric trees from the UPGMA method. To estimate  
 372 trees from the NJ method, we used the `nj` function of the `ape` package in R [30].  
 373 In this case, the estimated trees are not ultrametric, so we transform them in  
 374 ultrametric trees using the `chronMPL` and `multi2di` functions in `ape` package  
 375 [30, 31]. We used the `Rkitsch` function of the `Rphylip` package in R [32, 33] to  
 376 estimate ultrametric trees from the ME method assuming an evolutionary clock.  
 377 The NJ and ME methods are generally considered superior to UPGMA because  
 378 they optimize a tree according to minimum evolution criteria. Similarly to the  
 379 UPGMA, the NJ and ME methods are fast and efficient computationally.

### 380 3.3. Statistical indexes to compare phylogenies

381 To evaluate the accuracy of the phylogenies generated by the MRCAT algo-  
 382 rithm and by the genetic distance methods (UPGMA, FM and ME) in relation  
 383 to the true phylogeny generated by SSEE we use three statistics: the Robinson  
 384 and Foulds (RF [34]) metric, the gamma statistic ( $\gamma$  [35]) and the Sackin's index  
 385 ( $I_s$  [36, 37]).

386 The RF metric measures the distance between phylogenetic trees, providing  
 387 the overall topological resemblance of the phylogenies. Specifically, the RF  
 388 metric calculates the number of internal branches present in only one of the  
 389 trees being compared. Given two trees,  $T_1$  and  $T_2$ , we define

$$RF(T_1, T_2) = \frac{L_1}{L'_1} + \frac{L_2}{L'_2} \quad (9)$$

390 where  $L_1$  and  $L_2$  are the number of branches on  $T_1$  and  $T_2$ , respectively. The  
 391 number of branches shared by  $T_1$  and  $T_2$  are represented by  $L'_1$  and  $L'_2$ . The  
 392 RF metric was calculated using the `RF.dis` function of the `phangorn` package  
 393 in R [38].

394 The  $\gamma$ -statistic measures the distribution of branch lengths of a tree and is  
 395 defined as [35]:

$$\gamma = \frac{1}{D} \left[ \frac{1}{N_S - 2} \sum_{k=2}^{N_S-1} T(k) - T(N_S)/2 \right] \quad (10)$$

396 with

$$T(k) = \sum_{j=2}^k j g_j; \quad (11)$$

$$D = T(N_S) / \sqrt{12(N_S - 2)} \quad (12)$$

398 where  $N_S$  is the number of leaves and  $g_k$  is the time interval between speciation  
 399 events as represented by the nodes of the tree (see Fig. S4 in section III of  
 400 the SI). The  $\gamma$ -statistic was calculated using the `gammaStat` function of the `ape`  
 401 package in R [30].

402 The Sackin index measures the degree of imbalance, or asymmetry, of a tree  
403 [36, 37]. It is defined as

$$I_s = \sum_j d_j \quad (13)$$

404 in which  $d_j$  is the number of nodes to be traversed between each leaf  $j$  and the  
405 root, including the root [39]. The expected Sackin index under a pure birth  
406 process (the Yule model [40]) is

$$E(I_s(N_S)) = 2N_S \sum_{k=2}^{N_S} \frac{1}{k} \approx 2N_S \log N_S \quad (14)$$

407 where the approximation holds for  $N_S$  large [37]. Since the expected value of  
408 the Sackin index increases with the tree size, a normalized index is defined to  
409 compare trees of different sizes:

$$I_s^n = \frac{I_s(N_S) - E(I_s(N_S))}{N_S} \quad (15)$$

410 Here we used the normalized Sackin index to compare the phylogenies and  
411 calculated it using the `sackin` function of the `apTreeshape` package in R [41].

#### 412 4. Results

413 We ran simulations of the speciation model described in section 3.1 with  
414 parameters  $N = 1500$ ,  $L = 100$ ,  $B = 150$ ,  $S = 5$ ,  $G = 7$ ,  $\mu = 0.001$ ,  $D = 0.05$ ,  
415  $Q = 0.05$ . We start with the results of a single simulation to show examples of  
416 phylogenies. Figure 4 shows the population after 1000 generations, with squares  
417 representing individuals and colors indicating the 36 species generated. Species  
418 form spatial clusters, a consequence of the small  $S$  value used the simulation.

419 The true phylogenetic tree of the population, generated using the SSEE  
420 algorithm, is shown in Fig. 5. Figure 5(a) shows the full phylogeny, which  
421 includes all speciation and extinction events. The large number of events seen  
422 near the root of the tree correspond mostly to unsuccessful or incomplete speci-  
423 ation events, in which a group of individuals momentarily splits in two species  
424 but quickly recombines into a single species due to mutations. We distinguish  
425 these events from *true extinctions*, which are characterized by the collapse of a  
426 long living species by a sharp decline in population size. This phenomenon is  
427 very common at the beginning of the speciation process in the model described  
428 in section 3.1. In Fig. 5(b),(c),(d) the full phylogeny was filtered in order to  
429 remove speciation reversals and keep only *true* extinction events. In the model,  
430 extinctions occur by stochastic fluctuations in the number of individuals of a  
431 species, which might become very small and go to zero. Figure 5(b) shows the  
432 phylogeny filtered by the criterion of population size at the moment of van-  
433 ishing: species that disappear with more than 20 individuals were considered  
434 speciation reversals and removed from the tree. Figures 5(c) and (d) display  
435 the same phylogenies but filtered also by the criterion of persistence in time:

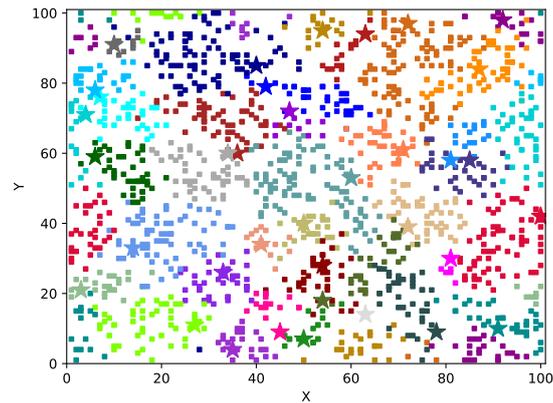


Figure 4: Spatial distribution of individuals from one simulation based on the model described in section 3.1. Individuals are represented by circles, and each color represents a different species. Stars indicate the individuals used to draw the phylogenies shown in figure 6.

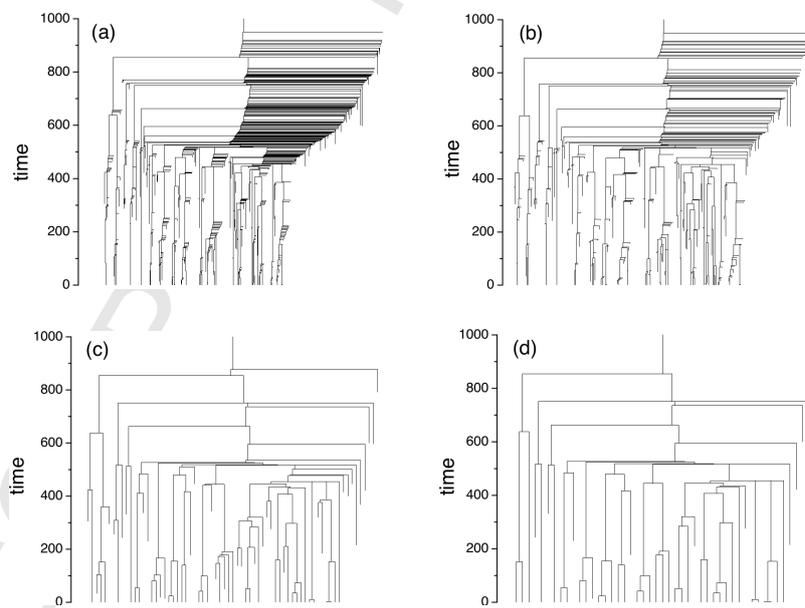


Figure 5: True phylogenies obtained with the SSEE method. (a) full phylogeny, including all speciation and extinction events; (b) filtered phylogeny, excluding branches (species) which had more than 20 individuals at the moment of extinction; (c) filtered phylogeny, excluding also branches that lasted less than 50 generations and (d) 100 generations.

436 branches of species that lasted less than 50 generations (c) or 100 generations  
437 (d) were also removed.

438 Phylogenies computed from the SSEE, MRCAT and genetic distance ma-  
439 trices are shown in Fig. 6. Panel (a) shows the true SSEE phylogeny, filtered  
440 to exhibit only the extant species. Panel (b) was obtained from the MRCAT  
441 algorithm, with one individual from each species being selected to represent  
442 the species. We showed in section II of the SI (Fig. S2) that the choice of  
443 the individual for constructing the phylogenetic tree with MRCAT can matter.  
444 However, the final structure of the tree will barely vary. Finally, panel (c) shows  
445 the phylogeny estimated from the genetic distance matrix of the same individu-  
446 als used in Fig. 6(b) by the UPGMA clustering method. Differences in topology  
447 and branch lengths are qualitatively visible between these trees. Maternal and  
448 paternal genealogies obtained from the MRCAT algorithm are shown in Fig. S3  
449 in the SI.

450 Statistical comparisons between phylogenies generated by the MRCAT algo-  
451 rithm and by the genetic distance methods (UPGMA, NJ and ME) in relation  
452 to the true phylogeny (SSEE) are shown in Fig. 7. The first line shows compar-  
453 isons of topology (RF metric), branch length distribution ( $\gamma$ -statistic) and  
454 degree of imbalance (Sackin index) among phylogenies after 500 generations in  
455 50 simulations. The second line shows the same comparisons after 1000 gen-  
456 erations for the same 50 simulations. Colors represent the different methods  
457 utilized to generate the trees. In the RF scatterplots (Fig. 7(a)(b)) the coor-  
458 dinates of each point refer to the normalized topological distance between the  
459 tree calculated with the MRCAT matrix ( $y$ -axis) or by genetic distance matrix  
460 ( $x$ -axis) from the true phylogenies generated by the SSEE algorithm. Small  
461 values of RF indicate that phylogenies are closer to the true phylogeny (SSEE).  
462 The diagonal dotted line defines the condition in which the topology of the  
463 phylogenies (RF-value) was equal in trees generated by genealogical relation-  
464 ships (MRCAT trees) and that estimated by genetic distance (UPGMA, NJ  
465 and ME methods). The scatterplot for  $T = 500$  (Fig. 7(a)) shows that phyloge-  
466 nies generated by MRCAT and genetic distance using UPGMA method (orange  
467 points) were similar in their RF-values, while trees estimated from NJ and ME  
468 methods (yellow and pink) had more different RF-values. For  $T = 1000$  (Fig.  
469 7(b)) all phylogenies estimated by genetic distance-based methods differ from  
470 those obtained by MRCAT. The density distribution of RF values shown above  
471 the scatterplots indicates that MRCAT is always closer to SSEE, especially for  
472  $T = 1000$ .

473 Regarding the branch length distribution, the scatterplots (Fig. 7(c),(d))  
474 show the difference between  $\gamma$ -values in SSEE phylogenies ( $y$ -axis) and MRCAT  
475 or genetic distance (UPGMA, NJ or ME) phylogenies ( $x$ -axis). The diagonal  
476 dotted line defines the condition in which the  $\gamma$ -values of trees generated by  
477 genealogical relationships (MRCAT trees) or by genetic distance (by UPGMA,  
478 NJ and ME methods) were equal to values of true phylogenies. We observe  
479 that for both times (Fig. 7(c),(d)) MRCAT trees had  $\gamma$  distributions closer  
480 to true phylogenies (SSEE) than all genetic distance-based trees, with a good  
481 match for  $T = 1000$ . Finally, the normalized Sackin index is presented in Fig.

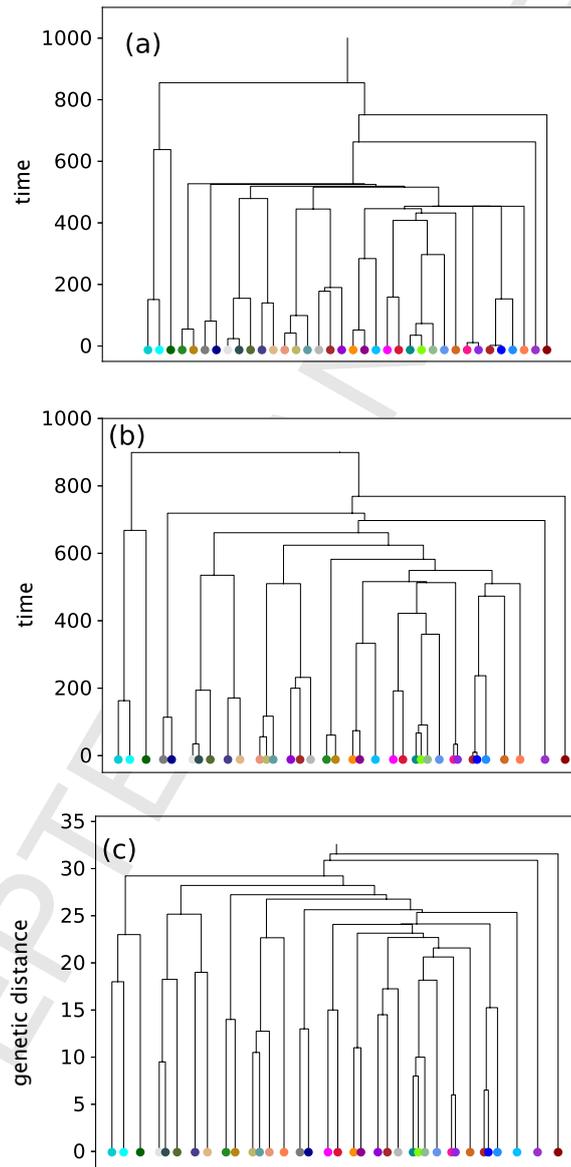


Figure 6: (a) Extant phylogeny obtained via SSEE (species are separated by one unit on x-axis); (b) via MRCAT; (c) via genetic distance matrix using UPGMA (neighbor species are separated by genetic distances). Colors correspond to species in Fig. 4.

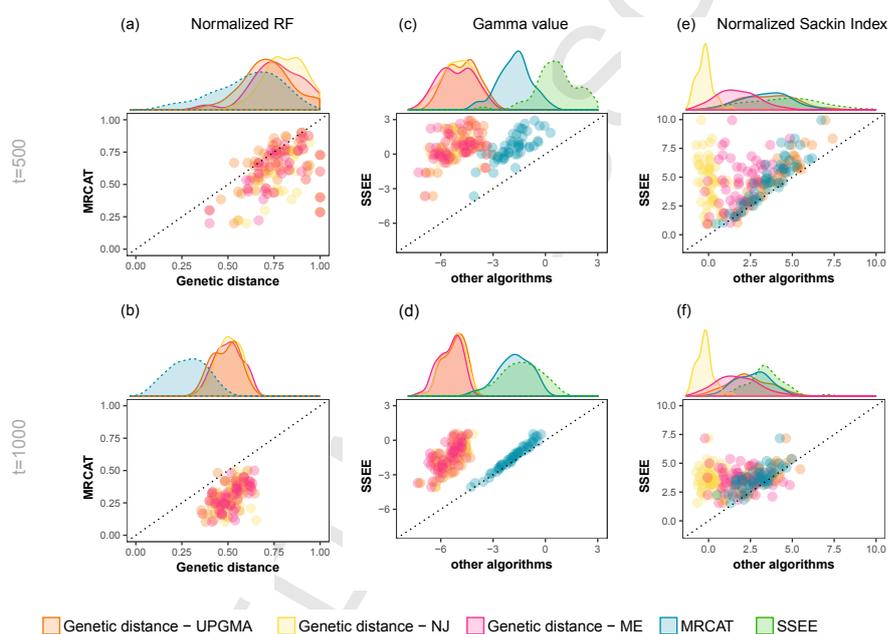


Figure 7: Comparisons among phylogenies generated by the algorithms proposed here (MRCAT and SSEE) and phylogenies estimated from genetic distance by UPGMA, NJ and ME methods. Lines exhibit the comparisons of RF, gamma and Sackin's metrics of 50 simulations at times 500 (first line) and 1000 (second line) generations. Colors represent the different methods utilized to generate the trees. (a) and (b): difference between RF-values of phylogenies obtained by MRCAT ( $y$ -axis) and by genetic distance-based methods ( $x$ -axis). Small values of RF indicate that phylogenies are closer to the true phylogeny (SSEE). (c) and (d): difference between branch length distributions ( $\gamma$ ) of phylogenies generated by SSEE ( $y$ -axis, green distribution) and MRCAT algorithm (blue) or genetic distance-based methods (orange, yellow and pink) ( $x$ -axis). (e) and (f): the same as (c) and (d), but considering now the degree of imbalance (Sackin index). Distributions above all scatterplots illustrate qualitatively the differences in topology (a,b), branch length distribution (c,d) and degree of imbalance (e,f) of phylogenies generated from each algorithm or method in the 50 simulations.

(Fig. 7(e),(f)). The imbalance of MRCAT phylogenies was closer to the true phylogenies for  $T = 500$  (Fig. 7(e)). On the other hand, for  $T = 1000$  the imbalance was similar for MRCAT and all distance-based methods, except for the NJ. The NJ trees exhibited the most incorrect Sackin index (Fig. 7(e),(f)), possibly because NJ trees are not rooted, a necessary condition to compute this index. The rooting procedure chosen can be quite arbitrary, affecting the balance of the trees and consequently the Sackin index. The distributions above all scatterplots show qualitatively the differences in topology (Fig. 7(a),(b)), branch length distribution (Fig. 7(c),(d)) and degree of imbalance (Fig. 7(e),(f)) of phylogenies generated from each algorithm or method in the 50 simulations performed in each time ( $t = 500$  or  $t = 1000$ ).

## 5. Discussion

Understanding all the mechanisms that promote speciation is still an open problem in evolutionary biology [3, 42]. Even more challenging is to identify which of these mechanisms were important in a particular case. A large number of mathematical and computational models were developed in the past years to understand different speciation processes, such as neutral [43–46], sexual [47–49] and ecological selection [12, 50]. Models have also considered the role of geography in speciation, such as allopatric [51–54], parapatric [10, 55] and sympatric [12, 49, 56, 57] scenarios. The results of models, however, can seldom be compared with real data [58, 59]. In these cases comparisons are often made in a macroecological scale, including qualitative species abundance and spatial distributions, species-area relationships and genetic or phenotypic distances [2, 6, 7, 14, 15]. Nevertheless, little attention has been given to the evolutionary history of individuals and species, neglecting the macroevolutionary scale underlying the speciation process [16, 21].

In this paper we have described two procedures to register the history of individuals (MRCAT) and species (SSEE) in individual-based models. With the ancestor-descendant relationships or speciation events saved in MRCAT and SSEE matrices we have constructed trees using a clustering algorithm. These trees have properties demonstrated in section I of Supporting Information. In the MRCAT algorithm, genealogies of individuals and phylogenies of species were obtained, whereas in the SSEE algorithm only phylogenies of species can be accessed. In the SSEE algorithm speciation events are precisely recorded and the resulting phylogenetic tree is the *true* tree of the community, whereas in the MRCAT algorithm the relations among species are recovered from genealogical relationships between individuals of each species. The MRCAT algorithm allows the construction of maternal, paternal and general lineages, the last being analogous to cases with hermaphroditic individuals. We have applied these algorithms to a spatially explicit IBM where individuals are separated into males and females and sexual reproduction is restricted by genetic difference below a threshold and by spatial proximity. We showed that maternal, paternal and general genealogies generated from the MRCAT algorithm are different even if the same individuals are chosen to draw the trees (Supporting Information,

526 section II). Maternal and paternal genealogies (Fig. S3(a),(b)) are different be-  
527 cause they were obtained from different MRCAT matrices. In the first case, the  
528 MRCAT matrix contains the time to the most recent common *female* ancestor  
529 between each pair of individuals, while in the second case the MRCAT matrix  
530 has the time to the most recent common *male* ancestor between the same in-  
531 dividuals, which lead to different ancestor times and genealogical relationships.  
532 In addition, for the general genealogy - taking the most recent common ancestor  
533 among females and males (*i.e.*, disregarding sex) - the resulting MRCAT matrix  
534 does not uniquely specify the genealogy (Fig. S3(c)). Regarding the phyloge-  
535 netic trees, we showed that they may be different if obtained by MRCAT or  
536 SSEE algorithm (Fig. 6(a),(b), Fig. 7). As discussed in subsection 2.3, this  
537 mismatch happens because the time to the most recent common ancestor be-  
538 tween individuals of different species is only an approximation to the speciation  
539 time, since speciation can happen several generations later (Fig. 3).

540 Structural properties of phylogenies, such as the Sackin index and the gamma  
541 distribution, obtained from SSEE and MRCAT trees were compared to values  
542 calculated in phylogenies estimated from the genetic distance between individ-  
543 uals of extant species by distance-based methods (UPGMA, NJ and ME). The  
544 aim of this comparison was to show that the validity of these methods commonly  
545 used in empirical studies, where the complete past history is inaccessible, can  
546 be assessed with the help of models. Differences in topology and branch length  
547 distribution measured by the RF metric and  $\gamma$ -statistic, respectively, revealed  
548 that MRCAT trees were closer to the true phylogenies (SSEE) than genetic  
549 distance-based trees. The difference between the results of these two methods  
550 possibly lies in back mutations that can happen in the genome of individuals,  
551 erasing the information needed to uncover the real history among species [60].  
552 This phenomenon is more likely to happen at long times and for small genome  
553 size. Indeed, we observed that in 500 generations (Fig. 7(a)(c)) the phylogenies  
554 estimated from genetic distance were closer to the ones generated from MR-  
555 CAT algorithm than in 1000 generations (Fig. 7(b)(d)), because in the first  
556 case the number of back mutations were probably smaller. Another factor that  
557 might explain the difference between genetic distance-based and true phyloge-  
558 nies is the sampling of only one individual to estimate the trees in the first  
559 case [61]. However, phylogenies generated with MRCAT algorithm also used  
560 only one individual per species - the same individuals used to compute genetic  
561 distance indeed - which suggests that this is not a very important factor (Fig.  
562 7(a),(b),(c),(d)). The degree of imbalance showed a different picture, with less  
563 differences between MRCAT trees and genetic distance trees. Still, MRCAT  
564 trees were closer to the true phylogenies than the others. Trees estimated from  
565 genetic information in IBMs should be closer to the true phylogenies for larger  
566 genome sizes, where the probability of back mutations is smaller. Individual-  
567 based models with large or infinite genome sizes already available [26, 62] would  
568 provide good tests for measuring the accuracy of trees obtained by distance-  
569 based methods.

570 The better performance of MRCAT algorithm in recover the topology and  
571 balance of phylogenetic trees is not surprising, since matrices generated from

572 this algorithm hold the exact times to the most recent common ancestors. How-  
573 ever, this type of exact information cannot be recovered from empirical data of  
574 contemporary samples. On the other hand, distance-based methods are com-  
575 monly used for inference of phylogenetic trees from empirical data [61]. The  
576 advantage of these methods, especially the NJ method, is their computational  
577 efficiency. Indeed, cluster algorithms are faster than optimality criteria used  
578 in character-based methods, like maximum parsimony and maximum likelihood  
579 [61, 63]. Distance methods are particularly useful for analysis of data sets con-  
580 taining sequences with low levels of divergence [61]. However, methods based  
581 on genetic distances can perform poorly when the data set contains sequences  
582 with high levels of divergence due to greater sampling error in larger genetic dis-  
583 tances. As most distance-based methods do not account for the high variances  
584 of large distance estimates, the inference of phylogenetic relationships could be  
585 impaired when these methods are employed [61]. In our model, trees generated  
586 from genetic distance methods were more different from the true trees (SSEE)  
587 than MRCAT phylogenies possibly because of high divergence among simulated  
588 genomes. This also could explain the high similarity in tree summary statis-  
589 tics among distance methods (Fig. 7). Moreover, the worst performance of NJ  
590 method in recover tree balance might be due to the lack of an explicit optimiza-  
591 tion criterion in the selection of taxon pairs in the original method proposed  
592 by Saitou and Nei [24] and utilized here [30, 63]. In addition, the choice of a  
593 substitution model to compute the pairwise distance between sequences might  
594 be important to determine the efficacy of distance methods [61]. Here we used  
595 the Hamming distance to calculate differences between pairs of sequences, but  
596 other methods could yield different results [64–67].

597 Modifications of the model to include *loci* not linked to the computation  
598 of genetic threshold would be important to understand how phylogenetic trees  
599 computed from these *loci* would differ from the ones computed here. Changing  
600 parameters values such as genome size and mutation rate could also affect tree  
601 estimations from distance-based methods and are a possible direction to future  
602 research. Nevertheless, the incorporation of algorithms that record the evolu-  
603 tionary history of individuals and species in an IBM context is an important  
604 step to help understanding the patterns left by specific speciation mechanisms  
605 at the macroevolutionary level.

## 606 6. Conclusions

607 The recent interest in the role of evolutionary history to explain the spa-  
608 tial patterns of abundance and species diversity calls for the incorporation of  
609 phylogenetic trees in the speciation modeling approach. Phylogenetic trees are  
610 essential tools to understand macroevolutionary patterns of diversity. They re-  
611 veal how species are related to each other and the times between speciation  
612 events. Moreover, topological structure and branch length distribution also  
613 contain clues about processes originating a particular group of species. Previ-  
614 ous works have already considered this problem for simpler models where each  
615 mutation corresponds directly to a new species [16]. Our study provides the

616 first general attempt to extend individual-based models by incorporating the  
617 branching process using the ancestor-descendant relationships between individ-  
618 uals and species. We believe this methodology will help predict and classify the  
619 macroevolutionary branching process, as well as the corresponding macroeco-  
620 logical patterns (*e.g.*, species abundance distributions), resulting from different  
621 speciation models. The comparison of these results with empirical studies may  
622 clarify the role of different processes in generating the patterns observed in nature  
623 [4, 5]. Finally, the role of extinction in determining macroevolutionary  
624 patterns is an open field [19] which could be explored by using the full phyloge-  
625 netic trees generated from the SSEE algorithm introduced here.

#### 626 Acknowledgments

627 This research was supported by Sao Paulo Research Foundation (FAPESP),  
628 National Council for Scientific and Technological Development (CNPq), and  
629 Coordination for the Improvement of Higher Education Personnel (CAPES).  
630 We thank P. L. Costa for his constructive comments about the manuscript,  
631 and A. B. Martins and L. D. Fernandes, who provided expertise that greatly  
632 assisted the research, all contributing with their insights and comments about  
633 the research results.

634 Conflicts of interest: none.

- 635 [1] J. Coyne, H. Orr, *Speciation*, Sinauer Associates, Sunderland, MA, 2004.
- 636 [2] M. A. M. De Aguiar, M. Baranger, E. Baptestini, L. Kaufman, Y. Bar-  
637 Yam, Global patterns of speciation and diversity, *Nature* 460 (7253) (2009)  
638 384–387.
- 639 [3] S. Gavrillets, Models of speciation: Where are we now?, *Journal of heredity*  
640 105 (S1) (2014) 743–755.
- 641 [4] M. Turelli, N. H. Barton, J. A. Coyne, Theory and speciation, *Trends in*  
642 *Ecology & Evolution* 16 (7) (2001) 330–343.
- 643 [5] R. Field, B. A. Hawkins, H. V. Cornell, D. J. Currie, J. A. F. Diniz-Filho,  
644 J.-F. Guégan, D. M. Kaufman, J. T. Kerr, G. G. Mittelbach, T. Oberdorff,  
645 et al., Spatial species-richness gradients across scales: a meta-analysis,  
646 *Journal of Biogeography* 36 (1) (2009) 132–147.
- 647 [6] A. B. Martins, M. A. de Aguiar, Y. Bar-Yam, Evolution and stability of  
648 ring species, *Proceedings of the National Academy of Sciences* 110 (13)  
649 (2013) 5080–5084.
- 650 [7] F. May, A. Huth, T. Wiegand, Moving beyond abundance distributions:  
651 neutral theory and spatial patterns in a tropical forest, *Proc. R. Soc. B*  
652 282 (1802) (2015) 20141657.

- 653 [8] M. Kopp, Speciation and the neutral theory of biodiversity: Modes of  
654 speciation affect patterns of biodiversity in neutral communities., *BioEssays*  
655 : news and reviews in molecular, cellular and developmental biology 32 (7)  
656 (2010) 564–70.
- 657 [9] S. P. Hubbell, *The Unified Neutral Theory of Biodiversity and Biogeogra-*  
658 *phy*, Princeton University Press, Princeton, NJ, 2001.
- 659 [10] S. Gavrillets, H. Li, M. D. Vose, Patterns of parapatric speciation, *Evolution*  
660 54 (4) (2000) 1126–1134.
- 661 [11] U. Dieckmann, M. Doebeli, On the origin of species by sympatric specia-  
662 tion., *Nature* 400 (6742) (1999) 354–7.
- 663 [12] A. Rettelbach, M. Kopp, U. Dieckmann, J. Hermisson, Three modes of  
664 adaptive speciation in spatially structured populations, *The American Nat-*  
665 *uralist* 182 (6) (2013) E215–E234.
- 666 [13] S. Gavrillets, Perspective: models of speciation: what have we learned in  
667 40 years?, *Evolution* 57 (10) (2003) 2197–2215.
- 668 [14] P. R. A. Campos, E. D. C. Neto, V. M. d. Oliveira, M. A. F. Gomes, Neutral  
669 communities in fragmented landscapes, *Oikos* 121 (11) (2012) 1737–1748.
- 670 [15] H. G. Martín, N. Goldenfeld, On the origin and robustness of power-law  
671 species–area relationships in ecology, *Proceedings of the National Academy*  
672 *of Sciences* 103 (27) (2006) 10310–10315.
- 673 [16] M. Manceau, A. Lambert, H. Morlon, Phylogenies support out-of-  
674 equilibrium models of biodiversity, *Ecology letters* 18 (4) (2015) 347–356.
- 675 [17] A. L. Pigot, A. B. Phillimore, I. P. Owens, C. D. L. Orme, The shape and  
676 temporal dynamics of phylogenetic trees arising from geographic speciation,  
677 *Systematic biology* 59 (6) (2010) 660–673.
- 678 [18] O. Hagen, K. Hartmann, M. Steel, T. Stadler, Age-dependent speciation  
679 can explain the shape of empirical phylogenies, *Systematic Biology* 64 (3)  
680 (2015) 432–440.
- 681 [19] T. B. Quental, C. R. Marshall, The molecular phylogenetic signature of  
682 clades in decline, *PloS one* 6 (10) (2011) e25780.
- 683 [20] T. J. Davies, A. P. Allen, L. Borda-de Águas, J. Regetz, C. J. Melián,  
684 Neutral biodiversity theory can explain the imbalance of phylogenetic trees  
685 but not the tempo of their diversification, *Evolution* 65 (7) (2011) 1841–  
686 1850.
- 687 [21] J. Rosindell, L. J. Harmon, R. S. Etienne, Unifying ecology and macroevo-  
688 lution with individual-based theory, *Ecology letters* 18 (5) (2015) 472–482.

- 689 [22] D. L. DeAngelis, V. Grimm, Individual-based models in ecology after four  
690 decades, *F1000prime reports* 6.
- 691 [23] F. Murtagh, Complexities of hierarchic clustering algorithms: State of the  
692 art, *Computational Statistics Quarterly* 1 (2) (1984) 101–113.
- 693 [24] N. Saitou, M. Nei, The neighbor-joining method: a new method for recon-  
694 structing phylogenetic trees., *Molecular biology and evolution* 4 (4) (1987)  
695 406–425.
- 696 [25] A. Rzhetsky, M. Nei, Theoretical foundation of the minimum-evolution  
697 method of phylogenetic inference., *Molecular biology and evolution* 10 (5)  
698 (1993) 1073–1095.
- 699 [26] P. G. Higgs, B. Derrida, Genetic distance and species formation in evolving  
700 populations, *Journal of molecular evolution* 35 (5) (1992) 454–465.
- 701 [27] E. M. Baptestini, M. A. de Aguiar, Y. Bar-Yam, The role of sex separation  
702 in neutral speciation, *Theoretical ecology* 6 (2) (2013) 213–223.
- 703 [28] R. W. Hamming, Error detecting and error correcting codes, *Bell Labs*  
704 *Technical Journal* 29 (2) (1950) 147–160.
- 705 [29] R Core Team, *R: A Language and Environment for Statistical Computing*,  
706 *R Foundation for Statistical Computing*, Vienna, Austria (2017).  
707 URL <https://www.R-project.org/>
- 708 [30] E. Paradis, J. Claude, K. Strimmer, *Ape: analyses of phylogenetics and*  
709 *evolution in r language*, *Bioinformatics* 20 (2) (2004) 289–290.
- 710 [31] T. Britton, B. Oxelman, A. Vinnersten, K. Bremer, Phylogenetic dating  
711 with confidence intervals using mean path lengths, *Molecular phylogenetics*  
712 *and evolution* 24 (1) (2002) 58–65.
- 713 [32] L. J. Revell, S. A. Chamberlain, *Rphylic: an r interface for phylip*, *Methods*  
714 *in Ecology and Evolution* 5 (9) (2014) 976–981.
- 715 [33] J. Felsenstein, *PHYLIP: Phylogenetic inference program, version 3.6*  
716 (2005).
- 717 [34] D. F. Robinson, L. R. Foulds, Comparison of phylogenetic trees, *Mathe-*  
718 *matical biosciences* 53 (1-2) (1981) 131–147.
- 719 [35] O. G. Pybus, P. H. Harvey, Testing macro-evolutionary models using in-  
720 complete molecular phylogenies, *Proceedings of the Royal Society of Lon-*  
721 *don B: Biological Sciences* 267 (1459) (2000) 2267–2272.
- 722 [36] M. Sackin, "good" and "bad" phenograms, *Systematic Biology* 21 (2)  
723 (1972) 225–226.

- 724 [37] M. G. Blum, O. François, On statistical tests of phylogenetic tree im-  
725 balance: the sackin and other indices revisited, *Mathematical biosciences*  
726 195 (2) (2005) 141–153.
- 727 [38] K. P. Schliep, phangorn: phylogenetic analysis in r, *Bioinformatics* 27 (4)  
728 (2011) 592–593.
- 729 [39] B. L. Dearlove, S. D. Frost, Measuring asymmetry in time-stamped phylo-  
730 genies, *PLoS computational biology* 11 (7) (2015) e1004312.
- 731 [40] G. U. Yule, A mathematical theory of evolution, based on the conclusions of  
732 dr. jc willis, frs, *Philosophical transactions of the Royal Society of London.*  
733 Series B, containing papers of a biological character 213 (1925) 21–87.
- 734 [41] N. Bortolussi, E. Durand, M. Blum, O. Francois, apTreeshape: Analyses  
735 of Phylogenetic Treeshape, r package version 1.4-5 (2012).  
736 URL <https://CRAN.R-project.org/package=apTreeshape>
- 737 [42] M. Kirkpatrick, V. Ravigné, Speciation by natural and sexual selection:  
738 models and experiments, *The American Naturalist* 159 (S3) (2002) S22–  
739 S35.
- 740 [43] G. A. Hoelzer, R. Drewes, J. Meier, R. Doursat, Isolation-by-distance  
741 and outbreeding depression are sufficient to drive parapatric speciation in  
742 the absence of environmental influences, *PLoS Comput Biol* 4 (7) (2008)  
743 e1000126.
- 744 [44] P. Desjardins-Proulx, D. Gravel, How likely is speciation in neutral ecol-  
745 ogy?, *The American Naturalist* 179 (1) (2011) 137–144.
- 746 [45] C. J. Melián, D. Alonso, S. Allesina, R. S. Condit, R. S. Etienne, Does sex  
747 speed up evolutionary rate and increase biodiversity?, *PLoS Comput Biol*  
748 8 (3) (2012) e1002414.
- 749 [46] E. M. Baptestini, M. A. de Aguiar, Y. Bar-Yam, Conditions for neutral spe-  
750 ciation via isolation by distance, *Journal of theoretical biology* 335 (2013)  
751 51–56.
- 752 [47] G. S. van Doorn, P. Edelaar, F. J. Weissing, On the origin of species by  
753 natural and sexual selection, *Science* 326 (5960) (2009) 1704–1707.
- 754 [48] J. C. Uyeda, S. J. Arnold, P. A. Hohenlohe, L. S. Mead, Drift promotes  
755 speciation by sexual selection, *Evolution* 63 (3) (2009) 583–594.
- 756 [49] L. K. MGonigle, R. Mazzucco, S. P. Otto, U. Dieckmann, Sexual selec-  
757 tion enables long-term coexistence despite ecological equivalence, *Nature*  
758 484 (7395) (2012) 506–509.
- 759 [50] P. Nosil, *Ecological speciation*, Oxford University Press, 2012.

- 760 [51] J. L. Fierst, T. F. Hansen, Genetic architecture and postzygotic reproduc-  
761 tive isolation: evolution of bateson–dobzhansky–muller incompatibilities in  
762 a polygenic model, *Evolution* 64 (3) (2010) 675–693.
- 763 [52] S. Gourbiere, J. Mallet, Are species real? the shape of the species boundary  
764 with exponential failure, reinforcement, and the missing snowball, *Evolution*  
765 64 (1) (2010) 1–24.
- 766 [53] C. Fraïsse, J. Elderfield, J. Welch, The genetics of speciation: are complex  
767 incompatibilities easier to evolve?, *Journal of evolutionary biology* 27 (4)  
768 (2014) 688–699.
- 769 [54] R. Yamaguchi, Y. Iwasa, First passage time to allopatric speciation, *Inter-  
770 face Focus* 3 (6) (2013) 20130026.
- 771 [55] C. Bank, R. Bürger, J. Hermisson, The limits to parapatric speciation:  
772 Dobzhansky–muller incompatibilities in a continent–island model, *Genetics*  
773 191 (3) (2012) 845–863.
- 774 [56] R. Bürger, K. A. Schneider, M. Willensdorfer, S. Otto, The conditions  
775 for speciation through intraspecific competition, *Evolution* 60 (11) (2006)  
776 2185–2206.
- 777 [57] P. S. Pennings, M. Kopp, G. Meszéna, U. Dieckmann, J. Hermisson, An  
778 analytically tractable model for competitive speciation, *The American Nat-  
779 uralist* 171 (1) (2007) E44–E71.
- 780 [58] D. I. Bolnick, B. M. Fitzpatrick, Sympatric speciation: models and empir-  
781 ical evidence, *Annu. Rev. Ecol. Evol. Syst.* 38 (2007) 459–487.
- 782 [59] S. Gavrilets, J. B. Losos, Adaptive radiation: contrasting theory with data,  
783 *Science* 323 (5915) (2009) 732–737.
- 784 [60] J. Hein, M. Schierup, C. Wiuf, *Gene genealogies, variation and evolution:  
785 a primer in coalescent theory*, Oxford University Press, USA, 2004.
- 786 [61] Z. Yang, B. Rannala, *Molecular phylogenetics: principles and practice*,  
787 *Nature Reviews Genetics* 13 (5) (2012) 303.
- 788 [62] M. A. de Aguiar, Speciation in the derrida–higgs model with finite genomes  
789 and spatial populations, *Journal of Physics A: Mathematical and Theoret-  
790 ical* 50 (8) (2017) 085602.
- 791 [63] O. Gascuel, M. Steel, Neighbor-joining revealed, *Molecular biology and  
792 evolution* 23 (11) (2006) 1997–2000.
- 793 [64] T. H. Jukes, C. R. Cantor, et al., Evolution of protein molecules, *Mam-  
794 malian protein metabolism* 3 (21) (1969) 132.

- 795 [65] M. Kimura, A simple method for estimating evolutionary rates of base  
796 substitutions through comparative studies of nucleotide sequences, *Journal*  
797 *of molecular evolution* 16 (2) (1980) 111–120.
- 798 [66] M. Hasegawa, H. Kishino, T.-a. Yano, Dating of the human-ape splitting  
799 by a molecular clock of mitochondrial dna, *Journal of molecular evolution*  
800 22 (2) (1985) 160–174.
- 801 [67] Z. Yang, Estimating the pattern of nucleotide substitution, *Journal of*  
802 *molecular evolution* 39 (1) (1994) 105–111.

ACCEPTED MANUSCRIPT