



## Abstract

12

13 Organisms are often more likely to exchange genetic information with others that are similar to  
14 themselves. One of the most widely accepted mechanisms of RNA virus recombination requires  
15 substantial sequence similarity between the parental RNAs and is termed similarity-essential  
16 recombination. This mechanism may be considered analogous to assortative mating, an important  
17 form of non-random mating that can be found in animals and plants. Here we study the dy-  
18 namics of haplotype frequencies in populations evolving under similarity-essential recombination.  
19 Haplotypes are represented by a genome of  $B$  biallelic loci and the Hamming distance between  
20 individuals is used as a criterion for recombination. We derive the evolution equations for the  
21 haplotype frequencies assuming that recombination does not occur if the genetic distance is larger  
22 than a critical value  $G$  and that mutation occurs at a rate  $\mu$  per locus. Additionally, uniform  
23 crossover is considered. Although no fitness is directly associated to the haplotypes, we show  
24 that frequency-dependent selection emerges dynamically and governs the haplotype distribution.  
25 A critical mutation rate  $\mu_c$  can be identified as the error threshold transition, beyond which  
26 this selective information cannot be stored. For  $\mu < \mu_c$  the distribution consists of a dominant  
27 sequence surrounded by a cloud of closely related sequences, characterizing a quasispecies. For  
28  $\mu > \mu_c$  the distribution becomes uniform, with all haplotypes having the same frequency. In  
29 the case of extreme assortativeness, where individuals only recombine with others identical to  
30 themselves ( $G = 0$ ), the error threshold results  $\mu_c = 1/4$ , independently of the genome size. For  
31 weak assortativity ( $G = B - 1$ )  $\mu_c = 2^{-(B+1)}$  and for the case of no assortativity ( $G = B$ )  $\mu_c = 0$ .  
32 We compute the mutation threshold for  $0 < G < B$  and show that, for large  $B$ , it depends only  
33 on the ratio  $G/B$ . We discuss the consequences of these results for recombination in viruses and  
34 for speciation.

35

36 \*Corresponding author: [aguilar@ifi.unicamp.br](mailto:aguilar@ifi.unicamp.br)

37 Keywords: haplotype frequencies, frequency-dependent selection, quasispecies theory, neutral mod-  
38 els.

## 39 I. INTRODUCTION

40 Mate choice by phenotypic similarity, or assortative mating, is a form of non-random  
41 mating that plays important roles in evolution and speciation. The mechanism has been  
42 explored in several mathematical and computational models and is often implemented as  
43 occurring in a Mendelian trait determined by a single locus [1–3] or in a single quantitative  
44 trait determined by multiple loci with additive effects [4–8].

45 Mate choice, however, often relies on multiple cues that are determined genetically [9–12],  
46 such that the effect of the state of these traits on phenotypic similarity is additive [13–15]. In  
47 the past twenty years different models have described this type of assortative mating [16–22],  
48 which may also be interpreted as a multilocus generalization of the Bateson-Dobzhansky-  
49 Muller model of intrinsic genetic incompatibilities [23].

50 Interestingly, a form of assortativity also occurs in virus populations. Contrary to what  
51 was initially thought, recombination is now considered to be a general phenomenon in RNA  
52 viruses and might play a major role as a driving force in virus variability and evolution [24].  
53 Although the mechanisms of viral RNA recombination are only now beginning to be eluci-  
54 dated, in the most widely accepted mechanism of viral recombination the enzyme responsible  
55 for replication switches from one sequence to another during the synthesis generating a re-  
56 combinant genome [25, 26]. This sequence switch is known to be dependent of the extent of  
57 similarity between the recombining genomes [27–29] and referred to as similarity-essential  
58 recombination [30]. Although there is strong evidence that the genetic exchange promoted  
59 by recombination can offer advantages, random recombination destroys more good alleles  
60 than it creates, leading to a selective pressure towards close similarity in the process [31].

61 In this paper we develop a theory for the evolution of recombinant haplotypes subjected  
62 to point mutations and similarity-essential recombination, but no other selective pressures.  
63 Each sequence will be represented by a binary string of length  $B$  and we will assume that  
64 sequences differing in more than  $G$  loci do not recombine. This mode of recombination is  
65 analogous to the assortative mating that often appears in models of population genetics, but  
66 may also represent replication of RNA viruses as described above. We will write the evolution  
67 equations for the haplotype frequencies and find their equilibrium solutions. For the case  
68 of zero mutation probability we will show that the population evolves to an equilibrium  
69 quasispecies where  $G$  loci are polymorphic and all the remaining  $B - G$  loci get fixed. This

70 equilibrium configuration is equivalent to that of a population with only  $G$  loci evolving under  
71 unconstrained recombination. However, which loci will remain polymorphic depends on the  
72 initial conditions. For non-zero mutation probability the scenario is more subtle. Our results  
73 show that, at least qualitatively, the population under similarity-essential recombination  
74 behaves in a way similar to the original quasispecies model for replicating macromolecules  
75 [32].

76 The quasispecies theory was originally developed by Eigen and Schuster to study the  
77 evolution of prebiotic RNA molecules exploring the consequences of the mutation-selection  
78 dynamics in near-infinite populations. Mutation rates are thought to have been much higher  
79 in the early history of life. As a result, equilibrium populations can be described as a  
80 a distribution of related genotypes known as quasispecies [32, 33]. Later the theory has  
81 successfully been applied to the study of viral evolution, especially RNA viruses [34–36].

82 One of the main results of the quasispecies theory is the existence of a mutation rate  
83 above which selection cannot overcome the mutation load (*i.e.* the error threshold). The  
84 amount of information that can be encoded in such evolutionary systems is limited by the  
85 genome length, since longer sequences suffer from mutations more than shorter ones. This  
86 leads to a logical enigma called Eigen’s paradox [37, 38]: given the mutations rates of this  
87 prebiotic scenario, these early genomes would not be long enough to encode the enzymes  
88 required to increase replication accuracy [32, 33]. Different mechanisms have been proposed  
89 to overcome or alleviate the genome size constraint imposed by the error threshold and  
90 warrant stable integration of information contained in the self-replicative units, like the the-  
91 ory of hypercycles [32], group selection models [39] and models incorporating recombination  
92 [40–42], and/or more complex genotype-phenotype mapping [43–45].

93 Here we show that, as in the quasispecies theory, recombinant haplotypes evolving under  
94 similarity-essential recombination exhibit two equilibrium regimes separated by a critical  
95 mutation rate  $\mu_c$ . In the first regime, which takes place for  $\mu < \mu_c$ , a dominant haplotype  
96 coexists with a cloud of closely related haplotypes. In the second regime, there is an in-  
97 formation crisis and an uniform distribution of haplotypes is obtained. Depending on the  
98 degree of assortativity, described by the parameter  $G$ , the error threshold can be as high as  
99  $\mu_c = 1/4$  and independent of the genome size, or as low as  $\mu_c = 2^{-(B+1)}$ , in contrast with  
100 the  $1/B$  behavior obtained in the original quasispecies model. We compute the mutation  
101 threshold for all values of  $G$  and show that, for large  $B$ , it depends only on the ratio  $G/B$ .

102 **II. MATERIALS AND METHODS**

103 We consider a population of haploid individuals with  $B$  biallelic loci. The genome of each  
 104 individual is represented by a string of  $B$  binary digits

$$i = i_1 i_2 \dots i_B \quad (1)$$

105 where the alleles  $i_k$  are either 0 or 1. We introduce the genotypic distance between two  
 106 haplotypes as the number of different alleles between them:

$$d(i, j) = \sum_{k=1}^B |i_k - j_k|. \quad (2)$$

107 Similarity-essential recombination corresponds to forbidding mating if  $d(i, j) > G$ , where  
 108  $G \leq B$ .

109 **A. Unconstrained recombination ( $G = B$ )**

110 Using the compact notation  $p_i$  for the frequency of haplotype  $i = i_1 i_2 \dots i_B$ , the equation  
 111 determining the frequency  $p_i^{t+1}$  in terms of the frequencies at time  $t$  assumes the form

$$p_i^{t+1} = \sum_{j,k} c_\mu(j, k; i) p_j^t p_k^t \quad (3)$$

112 where  $c_\mu(j, k; i)$  is the probability that individuals with haplotypes  $j$  and  $k$  produce a recom-  
 113 binant haplotype  $i$  if the mutation rate is  $\mu$ , whereas  $p_j^t p_k^t$  is the probability of an encounter  
 114 of haplotypes  $j$  and  $k$  at time  $t$ . To determine these coefficients we assume independent  
 115 segregation (uniform crossover) and look at one locus at a time. For a given allele  $i_n$  there  
 116 are three possibilities for  $j_n$  and  $k_n$ :

117

118 (a)  $j_n = k_n \neq i_n$ .

119 In this case the allele transmitted to the recombinant sequence is  $(1 - i_n)$  and it contributes  
 120 to  $p_i$  only if it mutates to  $i_n$ . Therefore it contributes a factor  $\mu$  to the probability. We call  
 121  $\alpha$  the number of loci satisfying this condition:

$$\alpha = \sum_{n=1}^B [1 - |j_n - k_n|] |i_n - j_n|. \quad (4)$$

122 (b)  $j_n = k_n = i_n$ .

123 The allele transmitted is  $i_n$  if it does not mutate. It contributes a factor  $(1 - \mu)$  and the  
 124 number of loci in this case is  $\beta$ :

$$\beta = \sum_{n=1}^B [1 - |j_n - k_n|] [1 - |i_n - j_n|] \quad (5)$$

125 (c)  $j_n \neq k_n$ .

126 The allele transmitted is either  $i_n$  or  $1 - i_n$ . It contributes a factor  $\frac{1}{2}(1 - \mu) + \frac{1}{2}\mu = \frac{1}{2}$ . The  
 127 number of loci of this type is

128

$$\gamma = \sum_{n=1}^B |j_n - k_n| = d(j, k). \quad (6)$$

129 It can be checked that  $\alpha + \beta + \gamma = B$  and that

$$\alpha = \frac{d(i, j) + d(i, k) - d(j, k)}{2}. \quad (7)$$

130 With these considerations equation (3) becomes

$$p_i^{t+1} = \sum_{j,k} p_j^t p_k^t (1 - \mu)^{B-\alpha-\gamma} \mu^\alpha \left(\frac{1}{2}\right)^\gamma. \quad (8)$$

131 It is interesting to report the two limiting cases  $\mu = 0$  and  $\mu = 1/2$ . In the first case only  
 132 haplotypes with  $\alpha = 0$  contribute to offspring and  $d(j, k) = d(i, j) + d(i, k)$ , showing that the  
 133 sum of the genetic distances from the recombinant sequence to each original sequence is the  
 134 genetic distance between the parents. If  $\mu = 1/2$ , corresponding to maximum randomness,  
 135 each pair of parental genomes contribute equally with weight  $2^{-B}$ .

136 The normalization condition is

$$\sum_i c_\mu(j, k; i) = \sum_i (1 - \mu)^{B-\alpha-\gamma} \mu^\alpha \left(\frac{1}{2}\right)^\gamma = 1 \quad (9)$$

137 and can be easily verified explicitly (see Electronic Supplementary Material, section I).

## 138 B. Similarity-essential recombination ( $G \neq B$ )

139 When genomes whose alleles differ in more than  $G$  loci are considered incompatible for  
 140 recombination, equations (8) have to be modified. In this case the sums over  $j$  and  $k$  on the  
 141 right hand side have to be restricted to parental sequences with  $\gamma = d(j, k) \leq G$  (see eq.(2))

142 and several terms are removed from the sum. Consequently, the equation has to be modified  
 143 in order to satisfy the normalization condition  $\sum_i p_i = 1$ . Normalization is ensured with the  
 144 introduction of an auxiliary function  $\Phi$  such that

$$p_i^{t+1} = \sum_{j,k,\gamma \leq G} p_j^t p_k^t c_\mu(j, k; i) - p_i^t (\Phi - 1). \quad (10)$$

145 Summing over  $i$  on both sides and using that  $\sum_i c_\mu(j, k; i) = \sum_i p_i = 1$  we find

$$\Phi = \sum_{j,k,\gamma \leq G} p_j^t p_k^t. \quad (11)$$

### 146 C. Analogy with the quasispecies theory and the error catastrophe

147 The quasispecies theory is originally a theory of molecular evolution [32]. In the Eigen  
 148 model molecules are represented by binary sequences of length  $L$  and the concentrations  $x_i$   
 149 of each type follow the equation

$$\frac{dx_i}{dt} = \sum_j x_j f_j q_{ji} - x_i \phi, \quad (12)$$

150 which assumes that the molecules replicate by cloning with mutations. In Eq. (12)  $f_j$  refers  
 151 to the replication rate (hereafter referred to as fitness) and the element of the mutation  
 152 matrix  $q_{ji} = (1 - \mu)^{B-d(i,j)} \mu^{d(i,j)}$  gives the probability that a molecule of type  $i$  produces a  
 153 molecule of type  $j$  if the mutation probability per digit is  $\mu$ . If the master string  $00\dots 0$  has  
 154 fitness  $f_0 > 1$  and all the remaining ones have fitness  $f_i = 1$ , it can be shown that a cloud of  
 155 mutant sequences surrounding and including the fittest master sequence (wild type) settles  
 156 in the population if

$$\mu < \frac{\log f_0}{B} \equiv \mu_c. \quad (13)$$

157 Above the mutation threshold  $\mu_c$  the population can no longer equilibrate in a mutation-  
 158 selection balance and the selection information is lost (error catastrophe).

159 Equation (10), which assumes similarity-essential recombination and discrete time, can  
 160 also be written in a similar form as the discrete time version of the Eigen's equation

$$p_i^{t+1} - p_i^t \equiv \sum_j p_j^t F_j Q_{ji} - p_i^t \Phi \quad (14)$$

161 with

$$Q_{ji} = \frac{\sum_{d(k,j) \leq G} p_k^t c_\mu(j, k; i)}{\sum_{d(k,j) \leq G} p_k^t} \quad (15)$$

162 and

$$F_i = \sum_{d(i,j) \leq G} p_j. \quad (16)$$

163  $Q_{ji}$  is the average probability that a haplotype  $j$  produces  $i$  by recombining with all com-  
164 patible haplotypes  $k$ . At this point the definition of  $F_i$  as the fitness of individuals of type  
165  $i$  comes about naturally, and is readily interpreted as the fraction of compatible individuals  
166 in the population. Note that  $\Phi = \sum_{j,k,\gamma \leq G} p_j^t p_k^t = \sum_j p_j F_j$  is, therefore, the average fitness of  
167 the population.

168 It is important to point out that the model assumes that all individuals are selectively  
169 equivalent regardless of their identities (neutral model). However, equation (14) demon-  
170 strates that frequency dependent selection arises from the similarity-essential recombination  
171 and can be quantified by equation (16). The fitness of an individual depends not on its  
172 specific haplotype but on the population composition. More importantly it is large if the  
173 individual is amongst compatible pairs (with which recombination is possible) and low if  
174 it is surrounded by incompatible mates. This idea concurs with the proper definition of  
175 quasispecies, at which natural selection is no longer directed toward a single variant but  
176 instead acts on the whole haplotype distribution [32].

### 177 III. RESULTS

#### 178 A. Unconstrained recombination

179 When  $G = B$  recombination is possible between every pair of individuals. In this case,  
180 since the alleles in each locus are segregated independently and there are no correlations  
181 between them, the result is that, for  $\mu = 0$  the allele frequencies remain constant from the  
182 first generation and the haplotype frequencies asymptotically reach the linkage equilibrium.  
183 For  $\mu \neq 0$  the haplotypes converge to the uniform distribution, and thereby all frequencies  
184 are equal to  $p_i = 2^{-B}$ . These results are well known [46, 47] and are demonstrated in the  
185 present context in the Electronic Supplementary Material, section II.



186 **B. Extreme assortativeness**

187 In the case of extreme assortativeness,  $G = 0$ , individuals only recombine with others  
 188 identical to themselves. In this case  $\gamma = d(j, k) = 0$ ,  $\alpha = d(i, j)$ ,  $F_j = p_j$  and  $Q_{ij} = q_{ij}$ .  
 189 Equation (14) becomes

$$p_i^{t+1} - p_i^t \equiv \sum_j (p_j^t)^2 q_{ji} - p_i^t \Phi \quad (17)$$

190 which is identical to (12) with  $F_i = p_i$ . For  $\mu = 0$  the equation simplifies to

$$p_i^{t+1} - p_i^t \equiv (p_i^t)^2 - p_i^t \Phi \quad (18)$$

191 and the only stationary solutions are:

- 192 (a) the single haplotype solution  $p_{i_0} = 1$  and  $p_i = 0$  for  $i \neq i_0$  and;  
 193 (b) the uniform solution  $p_i = 1/2^B$  for all  $i$ .

194 For small mutations a cloud of haplotypes similar to  $i_0$  is generated. Which haplotype  
 195 survives, along with its mutant cloud, is determined by the initial population [47–49]. As  
 196 the mutation rate increases the cloud spreads and, at  $\mu = \mu_c$ , the uniform solution becomes  
 197 stable. The mutation threshold can be calculated and results  $\mu_c = 1/4$ , independent of  
 198 the size of the genome  $B$  (see Electronic Supplementary Material, section VI). This should  
 199 be compared with equation (13), where the threshold becomes small as the genome size  $B$   
 200 increases. If we define an effective fitness  $F_0$  as the fitness of the corresponding quasispecies,  
 201 whose value is constant, that will result in the same error threshold, we find that  $1/4 =$   
 202  $\log(F_0)/B$  and so

$$F_0 = e^{B/4}. \quad (19)$$

203 The upper panel in figure 1 shows the frequencies of the haplotypes as a function of the  
 204 mutation probability  $\mu$ . The scenario displayed in the plot is essentially the pattern exhibited  
 205 by the quasispecies model: a dominant haplotype surrounded by a cloud of closely related  
 206 haplotypes.

207 Interestingly, for  $G = 1$  both expressions (12) and (15) work, even though  $Q_{ij} \neq q_{ij}$ . The  
 208 reason for this coincidence is that for  $G = 1$  reproduction occurs effectively with a single  
 209 locus and is equivalent to cloning one of the original haplotypes with equal probability.  
 210 Recombination affects only genetic exchanges between individuals with  $d(j, k) \geq 2$  [47].  
 211 Indeed, for  $G \geq 2$  only (14) is true.

212 A proof that the uniform distribution is always a solution for any value of  $G$  is presented  
 213 in the Electronic Supplementary Material, section V.

### 214 C. Error threshold for arbitrary $G$ and $B$

215 Explicit stationary solutions of equations (10) or (14) are not known, except for  $B = 2$   
 216 [47]. Because of the mating restriction imposed by the condition  $d(i, j) \leq G$ , the loci are  
 217 not independent and the dynamics of allele and haplotype frequencies are more complex and  
 218 richer. For zero mutation probability the haplotype frequencies converge to a distribution  
 219 where  $B - G$  loci are monomorphic (fixed in either 0 or 1) and the remaining  $G$  loci are  
 220 polymorphic in linkage equilibrium. Which loci become polymorphic depends on the initial  
 221 conditions and there are many possibilities. Indeed, for a population with this type of hap-  
 222 lotype distribution all individuals have maximum fitness  $F_i = 1$ , since the genetic distance  
 223 between any pair satisfies  $d(i, j) \leq G$  (see equation (16)). Moreover,  $\Phi = 1$  and equations  
 224 (10) become identical to (8) with  $B \rightarrow G$ .

225 However, the introduction of a small mutation rate  $\mu > 0$  generates mutants that decrease  
 226 the fitness of all resident types, resulting in further dynamics that converges to a single  
 227 dominant type plus a set of low frequency mutants. As  $\mu$  increases the distribution widens  
 228 and the uniform distribution, where  $p_i = 2^{-B}$  for all haplotypes, eventually takes over. The  
 229 mutation threshold for small values of  $G$  and for the limit case  $G = B - 1$  are:

$$\begin{aligned} \mu_c(B, G = 0) &= \frac{1}{4} \\ \mu_c(B, G = 1) &= \frac{(B - 1)}{4B} \\ \mu_c(B, G = 2) &= \frac{(B - 1)(B - 2)}{4(B^2 - B + 2)} \\ \mu_c(B, G = 3) &= \frac{(B - 1)(B - 2)(B - 3)}{4(B^3 - 3B^2 + 8B)} \\ \mu_c(B, G = 4) &= \frac{(B - 1)(B - 2)(B - 3)(B - 4)}{4(B^4 - 6B^3 + 23B^2 - 18B + 24)} \end{aligned} \tag{20}$$

230 and

$$\mu_c(B, G = B - 1) = 2^{-(B+1)}. \tag{21}$$

231 A detailed discussion is presented the Electronic Supplementary Material, section VI. These  
232 results are in qualitative agreement with previous numerical simulations in similar systems  
233 [50, 51].

234 Figure 2 shows  $\mu_c$  as a function of  $G/B$  for several values of  $B$  as calculated with the  
235 procedure indicated in the SI. For a fixed number of loci  $B$  the region under the corresponding  
236 curve indicates where a single haplotype or a set of closely related haplotypes dominates the  
237 population. It is interesting to note that this region is roughly independent of  $B$  and that it  
238 shrinks fast for  $G/B > 1/2$ . The more restrictive is the criterion for recombination (smaller  
239  $G/B$ ) the larger the interval of mutations leading to non-uniform distribution of alleles.  
240 This means that similarity-essential recombination turns the population less susceptible to  
241 the error-prone replication, i.e., the selective information can be kept for a broader range  
242 of values of mutation probabilities. This result is in contrast with the effects of standard  
243 recombination, which tends to lower the value of the critical mutation probability [49].  
244 Figure 1 also shows how the distribution of haplotypes changes with  $\mu$ . It is quite noticeable  
245 the shift of the error threshold to lower values as the assortativity  $G$  is reduced.

#### 246 IV. DISCUSSION

247 Here we studied the evolution of haplotype frequencies in an infinite population with  
248 similarity-essential recombination. We assumed that recombination is not possible if the  
249 genetic distance between two sequences is greater than a certain threshold  $G$ . In viral  
250 populations recombination often occurs when the replication enzyme switches from one  
251 molecule to another, and reducing  $G$  is equivalent to increasing the extent of similarity  
252 required for template switching. Depending on the similarity threshold the recombination  
253 is classified as precise or imprecise [24] and its value may be interpreted as a by-product of  
254 physical-chemical properties of the molecules involved in nucleic acid replication.

255 We derived the evolution equations for the haplotype frequencies assuming that each  
256 locus segregates independently (uniform crossover, as in [40, 41, 50]). This assumption is  
257 known not to be realistic. However, it is the simplest case to consider in order to highlight  
258 the effects of assortativity in the process of recombination. The correspondence between our  
259 equations and the quasispecies model allowed us to quantitatively determine the contribution  
260 of this type of recombination constraint to fitness, which we have shown to be equal to the

261 fraction of all compatible sequences in the population. It is important to highlight that, aside  
262 from the differences in fitness resulting from this constraint, our model is essentially neutral,  
263 since all individuals are assumed to be selectively identical. In spite of this neutrality at  
264 the individual level, natural selection arises as an outcome of the internal dynamics, which  
265 favors the selection of common haplotypes.

266 In the case of RNA viruses, most experimental studies are performed under strong selec-  
267 tive pressures, so that only the higher fitness types are detected [24]. Adding intrinsic fitness  
268 to the haplotypes would be a natural, though non-trivial, extension of our work. Depending  
269 on the initial conditions and on the strength of selection, the dynamics could give rise to  
270 a competition between the higher fitness types and those starting with large fractions of  
271 compatible individuals, leading to interesting properties of the haplotype distributions.

272 One of our main results is the observation and calculation of the error catastrophe in  
273 a population under similarity-based recombination. We demonstrated that for small  $G/B$   
274 the threshold tends to  $1/4$ , which is very large and independent of the genome size. At the  
275 other extreme, where  $G = B - 1$  the error threshold decreases exponentially as  $2^{-(B+1)}$ .  
276 Both these behaviors are in contrast with the  $1/B$  formula of the original Eigen model.  
277 The main consequence of the information crisis is the prediction of a maximum length for  
278 the sequence size beyond which the selective information is lost. In the context of prebi-  
279 otic evolution, this limits our understanding about how complex molecular structures can  
280 emerge from the prebiotic scenario. Here we have shown that replicating units subjected  
281 to similarity-essential recombination are able to safely transmit information at higher mu-  
282 tations through the emergence of a stable distribution of closely related haplotypes. These  
283 haplotypes naturally arise from the dynamics without defining a priori the set with large  
284 fitness. Because the error threshold is large and independent of the genome size for strong  
285 assortativity there is no restriction on the amount of information that can be stored in the  
286 system. The framework developed here provides exact results on the mutation thresholds  
287 for any values of genome size and assortativeness that could be applied to these populations.

288 Taken in a broader sense, the quasispecies framework may describe the evolution of any  
289 population of reproducing organisms [40]. In this case, similarity-based recombination is  
290 equivalent to assortative mating. This form of non-random mating plays an important  
291 role in the reproductive behavior of many populations. In mathematical models it can be  
292 introduced by prohibiting mating between individuals whose phenotypes are too dissimilar.

293 Depending on how genotypes are related to phenotypes, mate choice may be translated into  
294 a rule to be applied directly to genotypes. For haploid individuals with  $B$  biallelic loci where  
295 each locus represents a trait, assortative mating can be implemented by preventing mating  
296 between individuals whose haplotypes differ in more than  $G$  loci [4, 17, 18, 20–22]. These  
297 models assume that there are many characteristics controlling mating preference. In birds,  
298 for example, important traits are the color of plumage, patterns on plumage, song length,  
299 song complexity, beak size, body size, etc. Modeling each trait by a binary label results  
300 in  $2^B$  different phenotypes, where  $B$  is the genome size. For example, color is blue or red,  
301 beak is small or large, song is short or long, etc. Similar associations between haplotype and  
302 phenotype have also been used to study the branching of languages [52].

303 Our results also have implications for neutral models of speciation where reproductive  
304 isolation results from incompatibilities between individuals at the boundary between species  
305 [21, 22]. Individuals from each species have high fitness when amongst their own kind, but  
306 lower fitness at the boundary with other species, where the fraction of compatible mates  
307 drops. This feature keeps the populations isolated and prevents mixing in the absence of  
308 environmental selection.

309

#### 310 ACKNOWLEDGMENTS

311 This work was partly supported by FAPESP (MAMA, DMS and ABM) and CNPq (MAMA,  
312 PRAC).

- 
- 313 [1] H. S. Jennings, “The numerical results of diverse systems of breeding,” *Genetics*, vol. 1,  
314 pp. 53–89, Jan. 1916.
- 315 [2] S. Wright, “Systems of mating. III. assortative mating based on somatic resemblance,” *Ge-*  
316 *netics*, vol. 6, pp. 144–161, Mar. 1921.
- 317 [3] S. P. Otto, M. R. Servedio, and S. L. Nuismer, “Frequency-dependent selection and the evo-  
318 lution of assortative mating,” *Genetics*, vol. 179, pp. 2091–2112, Aug. 2008.
- 319 [4] A. S. Kondrashov and M. Shpak, “On the origin of species by means of assortative mating.,”  
320 *Proceedings of the Royal Society B: Biological Sciences*, vol. 265, pp. 2273–2278, Dec. 1998.
- 321 [5] M. Higashi, G. Takimoto, and N. Yamamura, “Sympatric speciation by sexual selection,”  
322 *Nature*, vol. 402, pp. 523–526, Dec. 1999.
- 323 [6] U. Dieckmann and M. Doebeli, “On the origin of species by sympatric speciation,” *Nature*,  
324 vol. 400, pp. 354–357, July 1999.
- 325 [7] D. I. Bolnick, “Waiting for sympatric speciation,” *Evolution*, vol. 58, pp. 895–899, Apr. 2004.  
326 ArticleType: research-article / Full publication date: Apr., 2004 / Copyright 2004 Society  
327 for the Study of Evolution.
- 328 [8] V. Schwmmle, K. Luz-Burgoa, J. S. S. Martins, and S. M. de Oliveira, “Phase transition  
329 in a mean-field model for sympatric speciation,” *Physica A: Statistical Mechanics and its*  
330 *Applications*, vol. 369, pp. 612–618, Sept. 2006.
- 331 [9] N. Burley, “Mate choice by multiple criteria in a monogamous species,” *The American Natu-*  
332 *ralist*, vol. 117, pp. 515–528, Apr. 1981. ArticleType: research-article / Full publication date:  
333 Apr., 1981 / Copyright 1981 The University of Chicago.
- 334 [10] U. Candolin, “The use of multiple cues in mate choice,” *Biological Reviews*, vol. 78, no. 4,  
335 pp. 575–595, 2003.
- 336 [11] S. F. Chenoweth and M. W. Blows, “Dissecting the complex genetic basis of mate choice,”  
337 *Nature Reviews Genetics*, vol. 7, pp. 681–692, Sept. 2006.
- 338 [12] P. A. Hohenlohe and S. J. Arnold, “Dimensionality of mate choice, sexual isolation, and  
339 speciation,” *Proceedings of the National Academy of Sciences of the United States of America*,  
340 vol. 107, pp. 16583–16588, Sept. 2010.
- 341 [13] J. M. Jawor, S. U. Linville, S. M. Beall, and R. Breitwisch, “Assortative mating by multiple

- 342 ornaments in northern cardinals (*cardinalis cardinalis*),” *Behavioral Ecology*, vol. 14, pp. 515–  
343 520, July 2003.
- 344 [14] J. Blais, M. Plenderleith, C. Rico, M. I. Taylor, O. Seehausen, C. v. Oosterhout, and G. F.  
345 Turner, “Assortative mating among lake malawi cichlid fish populations is not simply pre-  
346 dictable from male nuptial colour,” *BMC Evolutionary Biology*, vol. 9, p. 53, Mar. 2009.
- 347 [15] G. L. Conte and D. Schluter, “Experimental confirmation that body size determines mate  
348 preference via phenotype matching in a stickleback species pair,” *Evolution*, vol. 67, no. 5,  
349 pp. 1477–1484, 2013.
- 350 [16] M. Serva and L. Peliti, “A statistical model of an evolving population with sexual reproduc-  
351 tion,” *Journal of Physics A: Mathematical and General*, vol. 24, p. L705, July 1991.
- 352 [17] P. G. Higgs and B. Derrida, “Stochastic models for species formation in evolving populations,”  
353 *Journal of Physics A: Mathematical and General*, vol. 24, p. L985, Sept. 1991.
- 354 [18] P. G. Higgs and B. Derrida, “Genetic distance and species formation in evolving populations,”  
355 *Journal of Molecular Evolution*, vol. 35, pp. 454–465, Nov. 1992.
- 356 [19] F. Manzo and L. Peliti, “Geographic speciation in the derrida-higgs model of species forma-  
357 tion,” *Journal of Physics A: Mathematical and General*, vol. 27, p. 7079, Nov. 1994.
- 358 [20] S. Gavrillets, H. Li, and M. D. Vose, “Rapid parapatric speciation on holey adaptive land-  
359 scapes.,” *Proceedings of the Royal Society B: Biological Sciences*, vol. 265, pp. 1483–1489,  
360 Aug. 1998.
- 361 [21] M. A. M. de Aguiar, M. Baranger, E. M. Baptestini, L. Kaufman, and Y. Bar-Yam, “Global  
362 patterns of speciation and diversity,” *Nature*, vol. 460, pp. 384–387, July 2009.
- 363 [22] A. B. Martins, M. A. M. d. Aguiar, and Y. Bar-Yam, “Evolution and stability of ring species,”  
364 *Proceedings of the National Academy of Sciences*, vol. 110, pp. 5080–5084, Mar. 2013.
- 365 [23] S. Gavrillets, *Fitness Landscapes and the Origin of Species (MPB-41)*. Princeton University  
366 Press, July 2004.
- 367 [24] R. Aaziz and M. Tepfer, “Recombination in RNA viruses and in virus-resistant transgenic  
368 plants,” *Journal of General Virology*, vol. 80, pp.1339–1346, 1999.
- 369 [25] K. Delviks-Frankenberry, A. Galli, O. Nikolaitchik, H. Mens, V. K. Pathak, and W.-S. Hu,  
370 “Mechanisms and factors that influence high frequency retroviral recombination,” *Viruses*,  
371 vol. 3, pp. 1650–1680, Sept. 2011.
- 372 [26] E. Simon-Loriere and E. C. Holmes, “Why do RNA viruses recombine?,” *Nature Reviews*

- 373 *Microbiology*, vol. 9, pp. 617–626, Aug. 2011.
- 374 [27] J. Zhang and H. M. Temin, “Retrovirus recombination depends on the length of sequence  
375 identity and is not error prone.,” *Journal of Virology*, vol. 68, pp. 2409–2414, Apr. 1994.
- 376 [28] H. A. Baird, R. Galetto, Y. Gao, E. Simon-Loriere, M. Abreha, J. Archer, J. Fan, D. L.  
377 Robertson, E. J. Arts, and M. Negroni, “Sequence determinants of breakpoint location during  
378 HIV-1 intersubtype recombination,” *Nucleic Acids Research*, vol. 34, pp. 5203–5216, Oct.  
379 2006.
- 380 [29] E. Petterson, M. Stormoen, O. Evensen, A. B. Mikalsen, and O. Haugland, “Natural infection  
381 of atlantic salmon (*salmo salar* l.) with salmonid alphavirus 3 generates numerous viral deletion  
382 mutants,” *Journal of General Virology*, vol. 94, pp. 1945–1954, Sept. 2013.
- 383 [30] P. D. Nagy and A. E. Simon, “New insights into the mechanisms of RNA recombination,”  
384 *Virology*, vol. 235, pp. 1–9, Aug. 1997.
- 385 [31] M. Worobey and E. C. Holmes, “Evolutionary aspects of recombination in RNA viruses,”  
386 *Journal of General Virology*, vol. 80, pp. 2535–2543, 1999.
- 387 [32] M. Eigen and P. Schuster, “A principle of natural self-organization,” *Naturwissenschaften*,  
388 vol. 64, pp. 541–565, Nov. 1977.
- 389 [33] M. Eigen, “Selforganization of matter and the evolution of biological macromolecules,” *Natur-*  
390 *wissenschaften*, vol. 58, pp. 465–523, Oct. 1971.
- 391 [34] E. Domingo, D. Sabo, T. Taniguchi, and C. Weissmann, “Nucleotide sequence heterogeneity  
392 of an RNA phage population,” *Cell*, vol. 13, pp. 735–744, Apr. 1978.
- 393 [35] D. A. Steinhauer, J. C. d. l. Torre, E. Meier, and J. J. Holland, “Extreme heterogeneity in  
394 populations of vesicular stomatitis virus.,” *Journal of Virology*, vol. 63, pp. 2072–2080, May  
395 1989.
- 396 [36] E. Domingo, J. Sheldon, and C. Perales, “Viral quasispecies evolution,” *Microbiology and*  
397 *Molecular Biology Reviews*, vol. 76, pp. 159–216, June 2012.
- 398 [37] J. M. Smith and J. F. Y. Brookfield, “Models of evolution [and discussion],” *Proceedings of*  
399 *the Royal Society of London. Series B. Biological Sciences*, vol. 219, pp. 315–325, Oct. 1983.
- 400 [38] E. Szathmary, “The integration of the earliest genetic information,” *Trends in Ecology &*  
401 *Evolution*, vol. 4, pp. 200–204, July 1989.
- 402 [39] E. Szathmary and L. Demeter, “Group selection of early replicators and the origin of life,”  
403 *Journal of Theoretical Biology*, vol. 128, pp. 463–486, Oct. 1987.



- 404 [40] M. C. Boerlijst, S. Bonhoeffer, and M. A. Nowak, “Viral quasi-species and recombination,”  
 405 *Proceedings of the Royal Society of London. Series B: Biological Sciences*, vol. 263, pp. 1577–  
 406 1584, Nov. 1996.
- 407 [41] M. N. Jacobi and M. Nordahl, “Quasispecies and recombination,” *Theoretical Population*  
 408 *Biology*, vol. 70, pp. 479–485, Dec. 2006.
- 409 [42] J.-M. Park and M. W. Deem, “Phase diagrams of quasispecies theory with recombination and  
 410 horizontal gene transfer,” *Physical Review Letters*, vol. 98, p. 058101, Jan. 2007.
- 411 [43] F. K. d. Boer and P. Hogeweg, “Eco-evolutionary dynamics, coding structure and the infor-  
 412 mation threshold,” *BMC Evolutionary Biology*, vol. 10, p. 361, Nov. 2010.
- 413 [44] J. R. Peck and D. Waxman, “Is life impossible? information, sex, and the origin of complex  
 414 organisms,” *Evolution*, vol. 64, pp. 3300–3309, Nov. 2010.
- 415 [45] E. S. Colizzi and P. Hogeweg, “Evolution of functional diversification within quasispecies,”  
 416 *Genome Biology and Evolution*, vol. 6, pp. 1990–2007, Aug. 2014.
- 417 [46] W.J. Ewens. *Mathematical Population Genetics. I. Theoretical Introduction*. Series: Biomath-  
 418 ematics, Vol. 9 (New York: Springer Verlag, 1979).
- 419 [47] D. M. Schneider, E. do Carmo, A. B. Martins, and M. A. M. de Aguiar, “Toward a the-  
 420 ory of topopatric speciation: The role of genetic assortative mating,” *Physica A: Statistical*  
 421 *Mechanics and its Applications*, vol. 409, pp. 35–47, Sept. 2014.
- 422 [48] D. M. Schneider, E. do Carmo, and M. A. M. de Aguiar, “A dynamical analysis of allele  
 423 frequencies in populations evolving under assortative mating and mutations,” *Physica A:*  
 424 *Statistical Mechanics and its Applications*, vol. 421, pp. 54–68, 2015.
- 425 [49] M. N. Jacobi, and M. Nordahl, “Quasispecies and recombination,” *Theor. Popul. Biol.*,  
 426 vol. 70(4), pp. 479–485, 2006.
- 427 [50] G. Ochoa and K. Jaffe, “Assortative mating drastically alters the magnitude of error thresh-  
 428 olds,” in *Parallel Problem Solving from Nature - PPSN IX* (T. P. Runarsson, H.-G. Beyer,  
 429 E. Burke, J. J. Merelo-Guervs, L. D. Whitley, and X. Yao, eds.), no. 4193 in Lecture Notes  
 430 in Computer Science, pp. 890–899, Springer Berlin Heidelberg, Jan. 2006.
- 431 [51] E. Tannenbaum and J. F. Fontanari, “A quasispecies approach to the evolution of sexual  
 432 replication in unicellular organisms,” *Theory in Biosciences*, vol. 127, pp. 53–65, Mar. 2008.
- 433 [52] V. Schwmmle and P. M. C. de Oliveira, “A simple branching model that reproduces lan-  
 434 guage family and language population distributions,” *Physica A: Statistical Mechanics and*



436 FIGURE CAPTIONS

437

438 Figure 1: (Online version in color.) Equilibrium haplotype frequencies for  $B = 3$  and  $G = 0$ ,  
439  $G = 1$  and  $G = 2$  as a function of the mutation rate. Lines correspond, from top to bottom,  
440 to: 000 (black); 100, 010 and 001 (blue); 110, 101 and 011 (red); 111 (green). The initial  
441 condition is  $p_{000} = 1$  and the other frequencies zero.

442

443 Figure 2: (Online version in color.) Critical mutation as a function of  $G/B$ . The uniform  
444 distribution of haplotype frequencies becomes stable for  $\mu > \mu_c$ . The area below the curves  
445 correspond to a stable distribution of closely related haplotypes.

446

447

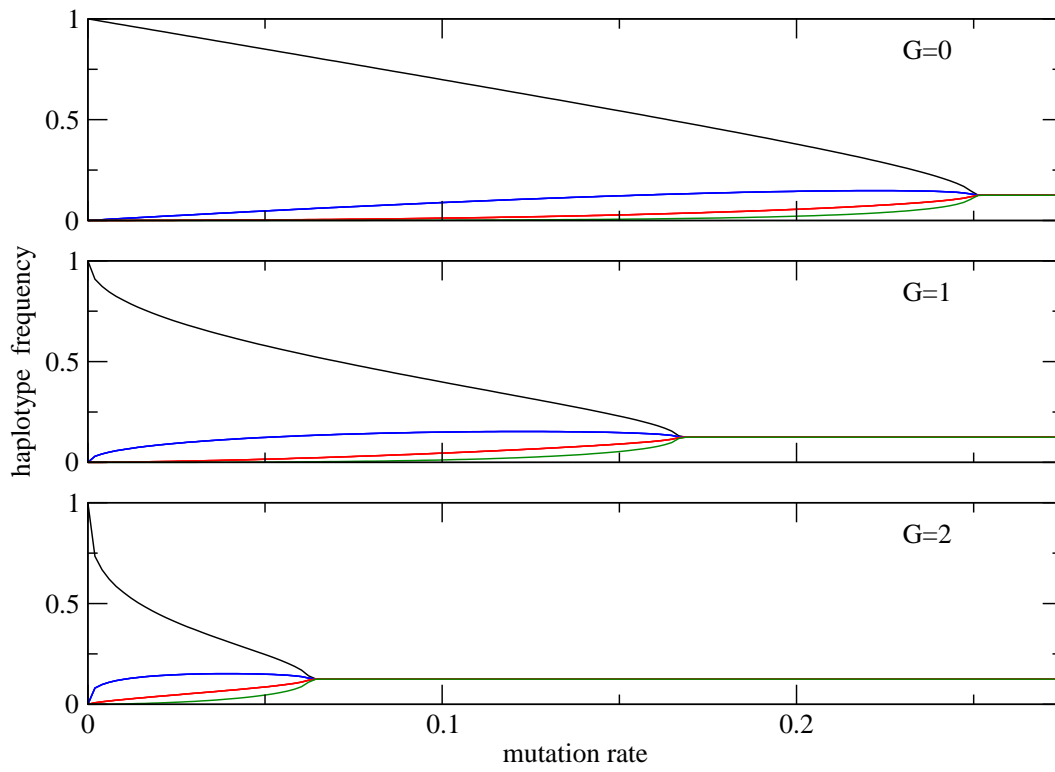


Figure 1.

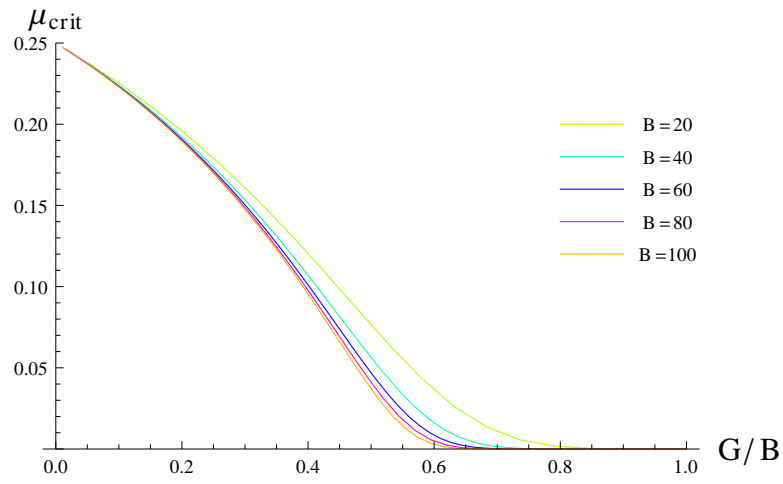


Figure 2.