

# Natural selection and the advantage of recombination

Philip J Gerrish<sup>1,2,3,\*</sup>, Benjamin Galeota-Sprung<sup>4</sup>, Fernando Cordero<sup>5</sup>, Paul Sniegowski<sup>4</sup>, Alexandre Colato<sup>6</sup>, Nicholas Hengartner<sup>2</sup>, Varun Vejalla<sup>7</sup>, Julien Chevallier<sup>8</sup> & Bernard Ycart<sup>8</sup>

<sup>1</sup>*Department of Biology, University of New Mexico, Albuquerque, New Mexico, USA*

<sup>2</sup>*Theoretical Biology & Biophysics, Los Alamos National Lab, Los Alamos, New Mexico, USA*

<sup>3</sup>*Instituto de Ciencias Biomédicas, Universidad Autónoma de Ciudad Juárez, México,*

<sup>4</sup>*Department of Biology, University of Pennsylvania, Philadelphia, Pennsylvania, USA*

<sup>5</sup>*Biomathematics and Theoretical Bioinformatics, Technische Fakultät, Universität Bielefeld, Germany*

<sup>6</sup>*Ciências da Natureza, Matemática e Educação, Univ Fed de São Carlos, Araras SP, Brazil*

<sup>7</sup>*Thomas Jefferson High School for Science and Technology, Alexandria, Virginia, USA*

<sup>8</sup>*Mathématique Appliquée, Laboratoire Jean Kuntzmann, Université Grenoble Alpes, France*

1 **The ubiquity of recombination (and sex) in nature has defied explanation since the time**  
2 **of Darwin<sup>1-4</sup>. Conditions that promote the evolution of recombination, however, are well-**  
3 **understood and arise when genomes contain more selectively mismatched combinations of**  
4 **alleles across loci than can be explained by chance alone. Recombination remedies this**  
5 **across-loci imbalance by shuffling alleles across individuals. The great difficulty in explain-**  
6 **ing the ubiquity of recombination in nature lies in identifying a source of this imbalance that**  
7 **is comparably ubiquitous. Here, we look to natural selection itself as a possible source of per-**  
8 **vasive imbalance, with the rationale that the ubiquity of natural selection approximates the**  
9 **ubiquity of sex and recombination in nature. Natural selection is fed by heritable variation**  
10 **which may be produced by any number of factors, such as drift, founder effects, migration**  
11 **and mutation. We ask how natural selection, acting on this variation, affects the across-loci**  
12 **imbalance and hence the evolutionary potential of recombination. Remarkably, we find that**  
13 **the effect of natural selection is to *always* promote the evolution of recombination, on average,**  
14 **independently of the source of the variation that feeds it. We show this is true for both across-**  
15 **and within-population recombination. Our findings suggest that recombination evolved and**  
16 **is maintained more as an unavoidable byproduct of natural selection than as a catalyst.**

17       The ability to exchange genetic material through recombination (and sex) is a heritable trait<sup>5,6</sup>  
18 that is influenced by many different evolutionary and ecological factors, both direct and indirect,  
19 both positive and negative. Evidence from nature clearly indicates that the net effect of these factors  
20 must be positive: recombination across all levels of organismal size and complexity is undeniably  
21 the rule rather than the exception<sup>2-4,7</sup>. Theoretical studies, on the other hand, have revealed a va-  
22 riety of different mechanisms and circumstances that can promote the evolution of recombination,  
23 but each one by itself is of limited scope<sup>2,4,8</sup>. These studies would thus predict that the absence of  
24 recombination is the rule and its presence an exception. The sheer abundance of these exceptions,  
25 however, can be seen as amounting to a rule in its own right – a “pluralist” view that has been  
26 adopted by some authors to explain the ubiquity of recombination<sup>3,7,9</sup>. The necessity of this plu-  
27 ralist view, however, may be seen as pointing toward a fundamental shortcoming in existing theory:  
28 perhaps some very general factor that would favour recombination has been missing<sup>3,4,8,10</sup>.

29       Existing theories of the evolution and maintenance of sex and recombination can be divided  
30 into those that invoke *direct* vs *indirect* selection on recombination. Theories invoking direct se-  
31 lection propose that recombination evolved and is maintained by some physiological effect that  
32 mechanisms of recombination themselves have on survival or on replication efficiency. Such the-  
33 ories might speak to the origins of sex and recombination but they falter when applied to their  
34 maintenance<sup>1</sup>. Most theories invoke indirect selection: they assume that any direct effect of recom-  
35 bination mechanisms is small compared to the trans-generational consequences of recombination.

36       While differing on the causal factors involved, established theoretical approaches that invoke

37 indirect selection are unanimous in their identification of the fundamental selective environment  
38 required for sex and recombination to evolve: a population must harbour an excess of selectively  
39 mismatched combinations of alleles across loci and a deficit of selectively matched combinations.  
40 Recombination is favoured under these conditions because on average it breaks up the mismatched  
41 combinations and assembles matched combinations. Assembling selectively matched combina-  
42 tions increases the efficiency of natural selection: putting high-fitness alleles together can expedite  
43 their fixation<sup>11-15</sup>, and putting low-fitness alleles together can expedite their elimination<sup>16,17</sup>. Under  
44 these conditions therefore, populations with recombination have an evolutionary advantage over  
45 populations without.

46 Furthermore, competition among recombination-rate variants at a *modifier* locus under these  
47 conditions will tend to increase recombination rate<sup>6,18</sup>. A modifier locus can be affected by short-  
48 term and long-term indirect selection<sup>13</sup>. In the long term, modifiers that increase the recombination  
49 rate, or *up-modifiers*, can in theory be favoured because of the fitness variation they augment<sup>13</sup>.  
50 But this relies on an unlikely – and somewhat inconsistent – supposition that the modifier remains  
51 linked to that variation. Under the more plausible supposition that modifiers themselves can be-  
52 come dissociated from fitness-related loci through the very recombination they modify, it is the  
53 short-term advantage or disadvantage of recombinants that determines the fate of a modifier. In the  
54 short term, recombinants formed from randomly-chosen parents harbouring the aforementioned  
55 imbalance have an advantage that is relatively immediate, thereby creating indirect selection for  
56 *up-modifiers*. If the imbalance is “chronic” such that recombinants are always advantageous on av-  
57 erage, then up-modifiers with only fleeting linkage to the recombinants they produce will increase

58 in frequency a notch with each such linkage, on average. Overall recombination rate in such a  
59 population will increase as a byproduct of chronic selection for recombinants.

60 The great challenge in explaining the evolution of recombination has been to identify an  
61 evolutionary source of the aforementioned imbalance whose prevalence in nature is comparable to  
62 the prevalence of sex and recombination in nature. One feature of living things whose prevalence  
63 approximates that of sex and recombination is evolution by natural selection. In what follows,  
64 we assess the effects of natural selection on selective imbalance and hence on the evolution of  
65 recombination.

66 We preface our developments with an essential technical point. In much of the relevant  
67 literature, the measure of selective mismatch across loci affecting the evolution of recombination  
68 is linkage disequilibrium (LD)<sup>8,12,13,19–22</sup>, which measures bias in allelic frequencies across loci  
69 but does not retain information about the selective value of those alleles. Here, the objectives  
70 of our study require a slight departure from tradition: our measure of selective mismatch will  
71 be covariance between genic fitnesses. This departure is necessary because covariance retains  
72 information about both the frequencies and selective value of alleles, and it is convenient because  
73 the mean selective advantage accrued by recombinants over the course of a single generation is  
74 equal to minus the covariance (Methods and Fig. S1). Our results will thus be given in terms of  
75 covariance and we recall: negative covariance means positive selection for recombinants.

## 76 **Setting**

77 For our two-part study, we reduce the problem to what we believe is its most essential form: we  
78 ask how the selective value of haploid recombinants is affected when natural selection simply acts  
79 on standing heritable variation. We ask this: 1) when recombination occurs between individuals  
80 from two different populations, and 2) when recombination occurs between two individuals within  
81 the same population.

82 To isolate the effects of natural selection, we consider large (effectively infinite) popula-  
83 tions, each of which consists of just two competing genotypes that differ in both of two genes (or  
84 two *loci*). This simple setting permits clean presentation and mathematical tractability and, more  
85 importantly, is biologically motivated by the observation that large clonal populations tend to be  
86 overwhelmingly dominated by one or two genotypes<sup>23</sup>. It further provides a connection to founda-  
87 tional evolution-of-sex studies: Fisher<sup>24</sup> considered the case of a single beneficial mutation arising  
88 on a variable background, thereby effectively giving rise to two competing genotypes – wildtype  
89 and beneficial mutant – that differ in both the gene with the beneficial mutation (call it the  $x$  gene)  
90 and its genetic background (call it the  $y$  gene); Muller<sup>25</sup> considered the case of two competing  
91 genotypes, one carrying a beneficial mutation in the  $x$  gene and the other in the  $y$  gene. Both of  
92 these approaches consider two competing genotypes that differ in both of two loci, and our en-  
93 compassing qualitative findings thus apply to these foundational models and others. Simulations  
94 further confirm the adequacy of this two-genotype setting: increasing the number of genotypes  
95 only accentuates the effects we describe (Fig. S2).

96 Figure 1 illustrates how the problem is posed analytically. We consider a clonal haploid  
97 organism whose genome consists of just two fitness-related loci labeled  $x$  and  $y$ . Genetically-  
98 encoded phenotypes at these two loci are quantified by random variables  $X$  and  $Y$ , both of which  
99 are positively correlated with fitness. In each large population of such organisms, two genotypes  
100 exist: one encodes the phenotype  $(X_1, Y_1)$ , has fitness  $Z_1 = \phi(X_1, Y_1)$  and exists at some arbitrary  
101 initial frequency  $p$ ; the other encodes phenotype  $(X_2, Y_2)$ , has fitness  $Z_2 = \phi(X_2, Y_2)$  and exists  
102 at initial frequency  $1 - p$ . We note that, in the absence of epistasis or dominance, the scenario we  
103 describe is formally equivalent to considering a diploid organism whose genome consists of one  
104 locus and two alleles available to each haploid copy. The question we ask is this: Does the action  
105 of natural selection, by itself, affect covariance between  $X$  and  $Y$ , denoted  $\sigma_{XY}$ , and if so, how?

## 106 **Natural selection promotes recombination *across* populations**

107 Here, we assess how natural selection affects the evolution of recombination *across* populations.  
108 Figure 2 illustrates the problem by analogy to a set of canoe races. On the surface, one might  
109 suspect that natural selection would promote well-matched combinations in which large values of  
110  $X$  are linked to large values of  $Y$ , thereby creating a positive association between  $X$  and  $Y$ . In  
111 fact, this notion is so intuitive that it is considered self-evident, explicitly or implicitly, in much  
112 of the literature<sup>1-3,7,9,14,26,27</sup>. If this notion were true, recombination would break up good allelic  
113 combinations, on average, and should thus be selectively suppressed. Such allele shuffling has  
114 been called “genome dilution”, a label that betrays its assumed costliness. We find, however, that  
115 the foregoing intuition is wrong. To the contrary, we find that natural selection will, on average,

116 promote an excess of mismatched combinations in which large values of  $X$  are linked to small  
117 values of  $Y$ , or vice versa, thereby creating a negative association between  $X$  and  $Y$ . Recombina-  
118 tion will on average break up the mismatched combinations created by natural selection, assemble  
119 well-matched combinations, and should thus be favoured.

120 Figure 3 illustrates why our initial intuition was wrong and why natural selection instead  
121 tends to create negative fitness associations among genes. For simplicity of presentation, we as-  
122 sume here that an individual's fitness is  $Z = \phi(X, Y) = X + Y$ , i.e., that  $X$  and  $Y$  are simply  
123 additive genic fitness contributions, and that  $X$  and  $Y$  are independent. In the absence of recom-  
124 bination, selection does not act independently on  $X$  and  $Y$  but on their sum,  $Z = X + Y$ . Perhaps  
125 counter-intuitively, this fact alone creates negative associations. To illustrate, we suppose that we  
126 know the fitness of successful genotypes to be some constant,  $z$ , such that  $X + Y = z$ ; here, we  
127 have the situation illustrated in Fig. 3a and we see that  $X$  and  $Y$  are negatively associated; indeed,  
128 covariance is immediate:  $\sigma_{XY} = -\sigma_X\sigma_Y \leq 0$ . Of course, in reality the fitnesses of successful  
129 genotypes will not be known *a priori* nor will they be equal to a constant; instead, they will follow  
130 a distribution of maxima of  $Z$  as illustrated in Fig. 3b. This is because, in large populations, the  
131 successful genotype will practically always be the genotype of maximum fitness. If populations  
132 consist of  $n$  contending genotypes, then  $X_{(n)} + Y_{(n)} = Z^{[n]}$ , the  $n^{th}$  order statistic of  $Z$  with genic  
133 components  $X_{(n)}$  and  $Y_{(n)}$  (called *concomitants* in the probability literature<sup>28,29</sup>). In general,  $Z^{[n]}$   
134 will have smaller variance than  $Z$ . Components  $X_{(n)}$  and  $Y_{(n)}$ , therefore, while not exactly follow-  
135 ing a line as in Fig. 3a, will instead be constrained to a comparatively narrow distribution about  
136 that straight line, illustrated by Fig. 3b, thereby creating a negative association. Figure 3c plots ten

137 thousand simulated populations evolving from their initial (green dots) to final (black dots) mean  
138 fitness components; this panel confirms the predicted negative association.

139 What we have shown so far is that, if recombination occurs across populations – or across  
140 *demes* in a structured metapopulation – the resulting offspring should be more fit than their par-  
141 ents, on average. This effect provides novel insight into established observations that population  
142 structure can favour recombination<sup>30–34</sup> and may even speak to notions that out-crossing can create  
143 hybrid vigour (heterosis).

144 Much of evolution indeed takes place in structured meta-populations providing ample op-  
145 portunity for cross-population (or cross-deme) recombination; it is thought, for example, that pri-  
146 mordial life forms evolved primarily on surfaces that provided spatial structure<sup>35,36</sup>. It is also true,  
147 however, that much of evolution takes place within unstructured (or “well-mixed”) populations;  
148 primitive life forms, for example, also existed in planktonic form<sup>37</sup>. We now turn to the question of  
149 how evolution by natural selection affects the selective value of recombinants in such unstructured  
150 populations.

### 151 **Natural selection promotes recombination *within* populations**

152 Here, we assess how natural selection affects the selective value of recombinants *within* unstruc-  
153 tured populations. Here again, Fig. 1 shows how the problem is posed analytically. Natural selec-  
154 tion will cause the two competing genotypes to change in frequency, causing  $\sigma_{XY}$  to change over  
155 time ( $\sigma_{XY} = \sigma_{XY}(t)$ ). Our measure of the net effect of natural selection on recombination is the



156 quantity  $\int_0^\infty \sigma_{XY}(t)dt$ ; if positive (negative), natural selection opposes (favours) recombination.

157 In Methods, we show that, in expectation, covariance over the long run is unconditionally  
158 non-positive,  $\mathbb{E}[\int_0^\infty \sigma_{XY}(t)dt] \leq 0$ , implying that the process of natural selection, on average, al-  
159 ways creates conditions that favour recombination. Remarkably, this finding requires no assump-  
160 tions about the distribution of allelic fitness contributions  $X$  and  $Y$ ; in fact, a smooth density is  
161 not required. Indeed, this distribution can have strongly positive covariance, and yet the net effect  
162 of natural selection is still to create negative time-integrated covariance. Put differently, this result  
163 is completely independent of the source of the heritable variation upon which natural selection  
164 acts – whether it be drift, migration, mutation, etc, or what the specific parameters, dynamics or  
165 interactions of these processes might be.

166 Our analyses further show that natural selection creates recombinant advantage even when  
167 recombinants are present in the initial variation upon which natural selection acts. Put differently,  
168 even in the presence of recombination, the effect of natural selection is to promote increased recom-  
169 bination. The implication is that natural selection not only promotes the evolution of recombination  
170 but also its maintenance (Fig. S3).

171 Finally, we show that it is primarily the additive component of fitness that causes time-  
172 integrated covariance to be negative. This fact stands in contrast to some previous indications that  
173 non-additive effects, specifically negative or fluctuating epistasis, are an essential ingredient in the  
174 evolution of recombination<sup>19,21,38–41</sup>.

## 175 **Discussion**

176 Some authors<sup>2,42</sup> have argued that negative associations build up within a population because pos-  
177 itive associations, in which alleles at different loci are selectively matched, are either removed  
178 efficiently (when they are both similarly deleterious), or fixed efficiently (when they are both sim-  
179 ilarly beneficial), thereby contributing little to overall within-population associations. Genotypes  
180 that are selectively mismatched, on the other hand, have longer sojourn times, as the less-fit loci  
181 effectively shield linked higher-fitness loci from selection. The net effect, it is argued, should be  
182 that alleles across loci will on average be selectively mismatched within a population. The find-  
183 ings from part one of our study differ from these arguments: we find that even genotypes that  
184 are ultimately fixed carry selectively mismatched alleles. The findings from part two of our study,  
185 however, are entirely consistent with these arguments; indeed, these arguments provide an intuitive  
186 way to understand our remarkable Proposition 7 (Methods).

187 We have identified a phenomenon that is an inherent consequence of natural selection and  
188 gives rise to selectively mismatched combinations of alleles across loci. Generally speaking, this  
189 pervasive phenomenon is an example of counter-intuitive effects caused by probabilistic condition-  
190 ing. For example, “Berkson’s paradox”<sup>43,44</sup> arises when a biased observational procedure produces  
191 spurious negative correlations. In the original context, among those admitted to hospital due to ill-  
192 ness, a negative correlation among potentially causative factors was observed because those with  
193 no illness (who tended to have no causative factors) were not admitted to the hospital and hence  
194 not observed. Similarly, negative correlations arise across genic fitnesses in part because genotypes

195 in which both loci have low genic fitness are purged by selection; here, however, the bias is not  
196 observational but actual, as these low-fitness genotypes no longer exist in the population.

197 Many previous studies, in one way or another, point to the increase in agility and efficiency  
198 of adaptation that recombination confers as the primary cause of its evolution. Here, we invert  
199 the perspective of those earlier studies, asking not whether recombination speeds adaptation, but  
200 whether adaptation via natural selection generally creates selective conditions that make recombi-  
201 nants directly and immediately advantageous. If so, as our findings indicate, then: 1) the ubiquity  
202 of recombination in nature might be less enigmatic than previously thought, and 2) perhaps recom-  
203 bination arose and is maintained more as an unavoidable byproduct than as a catalyst of natural  
204 selection.

## 205 **Methods**

206 **Notes.** In the main text, we employ the shorthand  $\sigma_{XY}$  to denote covariance. In what follows,  
207 however, we use  $\sigma_{XY}$  and  $\text{Cov}(X, Y)$  (for clarity) interchangeably. Some of the results presented  
208 here rely on some simplifying assumptions for compact presentation; generalized results that relax  
209 these assumptions are presented in the Supplementary Information (SI). Several of the proofs here  
210 are abridged; full proofs are in the SI, as well as alternative and supplemental proofs. Here, we  
211 restrict our analyses to the case of 2 loci and 2 alleles per locus; in the SI, we extend some of these  
212 analyses to  $m$  loci and  $n$  alleles per locus.

213 **Covariance and recombinant advantage.** Much work on the evolution of recombination employs

214 linkage disequilibrium (LD) as the measure of across-loci associations. It is straight-forward to  
215 estimate LD from genomic sequence data, which likely explains the popularity of this measure.  
216 LD, however, contains no information about the selective cost of such associations. Covariance,  
217 on the other hand, retains all of the information regarding both the prevalence of linkage and its  
218 selective cost (i.e., recombinant advantage), and is thus the measure we employ. We note that  
219 when the fitness function is a bivariate Bernoulli distribution ( $\phi(X, Y) = \mathbb{P}\{X = i, Y = j\} =$   
220  $p_{i,j}$ ,  $i, j \in \{0, 1\}$ ) then covariance and disequilibrium are equivalent ( $\sigma_{XY} = D = p_{1,1} - p_{1,\bullet}p_{\bullet,1}$ ).  
221 Recombinants are formed from two randomly-chosen contemporaneous parents such that their  
222 genetic makeup is simply an unbiased random sampling of the pool of available alleles at the  $x$  and  
223  $y$  loci. As such, their instantaneous advantage is zero on average:  $\mathbb{E}_R[X + Y] - \mathbb{E}[X + Y] = 0$ ,  
224 where subscript  $R$  denotes recombinant and no subscript denotes wildtype. Recombinants and  
225 wildtype, however, gain fitness at different rates:  $\partial_t \mathbb{E}_R[X + Y] = \sigma_X^2 + \sigma_Y^2$  and  $\partial_t \mathbb{E}[X + Y] =$   
226  $\sigma_X^2 + \sigma_Y^2 + 2\sigma_{XY}$ . A first order expansion thus reveals that the selective advantage of recombinants  
227 after a single generation of growth is  $\partial_t \mathbb{E}_R[X + Y] - \partial_t \mathbb{E}[X + Y] = -2\sigma_{XY}$ . A single-generation  
228 Moran model (Fig. S1) shows this prediction to be accurate and that covariance increases linearly  
229 in the first generation, implying that the mean selective advantage of recombinants over that first  
230 generation is  $-\sigma_{XY}$ . We note that a generalized linear relationship between fitness and phenotypes  
231  $X$  and  $Y$ , i.e.,  $Z = k_0 + k_X X + k_Y Y$ , yields a recombinant advantage of  $-k_X k_Y \sigma_{XY}$ . A full  
232 treatment of the relation between covariance and recombinant advantage is found in the SI, as well  
233 as the relation between our approach and classical population genetics.

234 **Natural selection promotes recombination across populations.** The setting for this problem is

235 shown in Fig 1. No hypothesis on the fitness function  $\phi$  is made at this point, apart from being  
 236 measurable. For the sake compact presentation we assume here (relaxed in SI) that  $(X_1, Y_1, X_2, Y_2)$   
 237 are i.i.d.; departures from this and other simplifying assumptions are dealt with in the SI. As defined  
 238 in Fig 1,  $Z_i = \phi(X_i, Y_i)$ ,  $Z^{[i]} = \phi(X_{(i)}, Y_{(i)})$ , and  $Z^{[2]} > Z^{[1]}$ .

239 **PROPOSITION 1.** *Let  $\psi$  be any measurable function from  $\mathbb{R}^2$  into  $\mathbb{R}$ . Then:  $\frac{1}{2}\mathbb{E}(\psi(X_{(1)}, Y_{(1)})) +$   
 240  $\frac{1}{2}\mathbb{E}(\psi(X_{(2)}, Y_{(2)})) = \mathbb{E}(\psi(X_1, Y_1))$ . In particular, the arithmetic mean of  $\mathbb{E}(X_{(1)})$  and  $\mathbb{E}(X_{(2)})$  is  
 241  $\mathbb{E}(X_1)$ .*

242 **PROOF:** Consider a random index  $I \in \{1, 2\}$ , and for now  $\mathbb{P}(I = 1) = \mathbb{P}(I = 2) =$   
 243  $1/2$ , and  $I$  is independent of  $(X_1, Y_1, X_2, Y_2)$ . The couple  $(X_I, Y_I)$  is distributed as  $(X_1, Y_1)$ .  
 244 Hence,  $\mathbb{E}(\psi(X_I, Y_I)) = \mathbb{E}(\psi(X_1, Y_1))$ , however,  $\mathbb{E}(\psi(X_I, Y_I)) = \mathbb{E}(\mathbb{E}(\psi(X_I, Y_I) | I)) =$   
 245  $\frac{1}{2}\mathbb{E}(\psi(X_{(1)}, Y_{(1)})) + \frac{1}{2}\mathbb{E}(\psi(X_{(2)}, Y_{(2)}))$ . □

247 **PROPOSITION 2.** *We have:  $\text{Cov}(X_{(1)}, Y_{(1)}) + \text{Cov}(X_{(2)}, Y_{(2)}) = -(\text{Cov}(X_{(1)}, Y_{(2)}) +$   
 246  $\text{Cov}(X_{(2)}, Y_{(1)})) = -\frac{1}{2}\mathbb{E}(X_{(2)} - X_{(1)})\mathbb{E}(Y_{(2)} - Y_{(1)})$ .*

249 **PROOF:** The couples  $(X_{(I)}, Y_{(I)})$  and  $(X_{(I)}, Y_{(3-I)})$  are both distributed as  $(X_1, Y_1)$ . There-  
 250 fore their covariances are null. These covariances can also be computed by condition-  
 251 ing on  $I$  (see e.g. formula (1.1) in <sup>45</sup>). For  $(X_{(I)}, Y_{(I)})$  we have:  $\text{Cov}(X_{(I)}, Y_{(I)}) =$   
 252  $\mathbb{E}(\text{Cov}(X_{(I)}, Y_{(I)} | I)) + \text{Cov}(\mathbb{E}(X_{(I)} | I), \mathbb{E}(Y_{(I)} | I))$ . On the right-hand side, the first term  
 253 is:  $\mathbb{E}(\text{Cov}(X_I, Y_I | I)) = \frac{1}{2}\text{Cov}(X_{(1)}, Y_{(1)}) + \frac{1}{2}\text{Cov}(X_{(2)}, Y_{(2)})$ . The second term is:  
 254  $\text{Cov}(\mathbb{E}(X_I | I), \mathbb{E}(Y_I | I)) = \frac{1}{4}\mathbb{E}(X_{(2)} - X_{(1)})\mathbb{E}(Y_{(2)} - Y_{(1)})$ . Similarly, we have:  $\text{Cov}(X_{(I)}, Y_{(3-I)}) =$   
 255  $\mathbb{E}(\text{Cov}(X_{(I)}, Y_{(3-I)} | I)) + \text{Cov}(\mathbb{E}(X_{(I)} | I), \mathbb{E}(Y_{(3-I)} | I))$ . The first term in the right-hand side is:  
 256  $\mathbb{E}(\text{Cov}(X_{(I)}, Y_{(3-I)} | I)) = \frac{1}{2}\text{Cov}(X_{(1)}, Y_{(2)}) + \frac{1}{2}\text{Cov}(X_{(2)}, Y_{(1)})$ . The second term in the right-

257 hand side is:  $\text{Cov}(\mathbb{E}(X_{(I)}|I), \mathbb{E}(Y_{(3-I)}|I)) = -\frac{1}{4}\mathbb{E}(X_{(2)} - X_{(1)})\mathbb{E}(Y_{(2)} - Y_{(1)})$ . Hence the result.

258

□

260 PROPOSITION 3. Assume that the fitness function  $\phi$  is symmetric:  $\phi(x, y) = \phi(y, x)$ . Then the  
259  
261 couple  $(X_{(1)}, Y_{(2)})$  has the same distribution as the couple  $(Y_{(1)}, X_{(2)})$ .

262 As a consequence,  $X_{(1)}$  and  $Y_{(1)}$  have the same distribution, so do  $X_{(2)}$  and  $Y_{(2)}$ . Thus:  $\mathbb{E}(X_{(2)} -$   
263  $X_{(1)}) = \mathbb{E}(Y_{(2)} - Y_{(1)}) = \frac{1}{2}\mathbb{E}(Z^{[2]} - Z^{[1]})$ . Another consequence is that:  $\text{Cov}(X_{(1)}, Y_{(2)}) =$   
264  $\text{Cov}(X_{(2)}, Y_{(1)})$ . Thus by Proposition 2:  $\text{Cov}(X_{(1)}, Y_{(2)}) = \text{Cov}(X_{(2)}, Y_{(1)}) = \frac{1}{16}\mathbb{E}^2(Z^{[2]} - Z^{[1]})$ .

265 PROOF: Since  $\phi$  is symmetric, the change of variable  $(X_1, Y_1, X_2, Y_2) \mapsto (Y_1, X_1, Y_2, X_2)$  leaves  
266 unchanged the couple  $(Z_1, Z_2)$ . □

268 PROPOSITION 4. Assume that the ranking function  $\phi$  is the sum:  $\phi(x, y) = x + y$ . Then:  $\mathbb{E}(X_{(1)}) =$   
269  
269  $\mathbb{E}(Y_{(1)})$ ,  $\mathbb{E}(X_{(2)}) = \mathbb{E}(Y_{(2)})$ , and  $\mathbb{E}(X_{(1)}) < \mathbb{E}(X_{(2)})$ .

270 PROOF: The first two equalities come from Proposition 3. By definition,  $\mathbb{E}(X_{(1)} + Y_{(1)}) < \mathbb{E}(X_{(2)} +$   
271  $Y_{(2)})$ . Hence the inequality. □

272 PROPOSITION 5. Assume that the ranking function  $\phi$  is the sum, and that the common distribution  
273  
274 of  $X_1, Y_1, X_2, Y_2$  is symmetric: there exists  $a$  such that  $f(x - a) = f(a - x)$ . Then  $(a - X_{(1)}, a -$   
275  $Y_{(1)})$  has the same distribution as  $(X_{(2)} - a, Y_{(2)} - a)$ . As a consequence,  $\text{Cov}(X_{(1)}, Y_{(1)}) =$   
276  $\text{Cov}(X_{(2)}, Y_{(2)})$ .

277 PROOF: The change of variable  $(X_1, Y_1, X_2, Y_2) \mapsto (2a - X_1, 2a - Y_1, 2a - X_2, 2a - Y_2)$  leaves  
278 the distribution unchanged. It only swaps the indices  $i$  and  $s$  of minimal and maximal sum. □

280 If we summarize Propositions 1, 2, 3, 4, 5 for the case where the ranking function is the sum, and  
 279  
 281 the distribution is symmetric, one gets:

$$\begin{aligned} \text{Cov}(X_{(1)}, Y_{(1)}) &= \text{Cov}(X_{(2)}, Y_{(2)}) < 0 \\ \text{Cov}(X_{(1)}, Y_{(2)}) &= \text{Cov}(X_{(2)}, Y_{(1)}) > 0 \\ |\text{Cov}(X_{(1)}, Y_{(1)})| &= \text{Cov}(X_{(1)}, Y_{(2)}) = \frac{1}{16} \mathbb{E}^2(Z^{[2]} - Z^{[1]}). \end{aligned}$$

282 **Natural selection promotes recombination *within* populations.** We recall that recombinant ad-  
 283 vantage is  $-\sigma_{XY}$ . Here, we study how the selection-driven changes in types  $(X_1, Y_1)$  and  $(X_2, Y_2)$   
 284 *within a single unstructured population* change  $\sigma_{XY} = \sigma_{XY}(t)$  over time. We are interested in  
 285 the net effect of these changes, given by  $\int_0^\infty \sigma_{XY}(t) dt$ ; in particular, we are interested in know-  
 286 ing whether this quantity is positive (net recombinant disadvantage) or negative (net recombinant  
 287 advantage).

288 **PROPOSITION 6.** *Within-population covariance integrated over time is:*

$$\int_0^\infty \sigma_{XY}(t) dt = q \mathbb{E} \left[ \frac{(X_2 - X_1)(Y_2 - Y_1)}{|Z_2 - Z_1|} \right] \quad (1)$$

289 *where  $q$  is the initial frequency of the inferior genotype. No assumption about the distribution of*  
 290  *$(X, Y)$  is required. And  $Z_i = \phi(X_i, Y_i)$  where fitness function  $\phi$  can be any function.*

291 **PROOF:** We let  $p$  denote initial frequency of the superior of the two genotypes, and we let  $q = 1 - p$   
 292 denote initial frequency of the inferior genotype. Time-integrated covariance is:

$$\int_0^\infty \sigma_{X,Y}(t) dt = \mathbb{E} \left[ (X_{(2)} - X_{(1)})(Y_{(2)} - Y_{(1)}) \int_0^\infty \frac{p q e^{(Z^{[1]} + Z^{[2]})t}}{(p e^{Z^{[2]}t} + q e^{Z^{[1]}t})^2} dt \right]$$

293 **Integration by parts yields:**

$$\int_0^\infty \sigma_{XY}(t) dt = q \mathbb{E} \left[ \frac{(X_{(2)} - X_{(1)})(Y_{(2)} - Y_{(1)})}{Z^{[2]} - Z^{[1]}} \right]$$

294 where  $q$  in Prop 6 is written as  $1 - p_0$ . We observe that:

$$(X_{(2)} - X_{(1)})(Y_{(2)} - Y_{(1)}) = (X_{(1)} - X_{(2)})(Y_{(1)} - Y_{(2)}) = (X_2 - X_1)(Y_2 - Y_1)$$

295 and that

$$Z^{[2]} - Z^{[1]} = |Z_2 - Z_1|$$

296 from which we have:

$$\mathbb{E}\left[\frac{(X_{(2)} - X_{(1)})(Y_{(2)} - Y_{(1)})}{Z^{[2]} - Z^{[1]}}\right] = \mathbb{E}\left[\frac{(X_2 - X_1)(Y_2 - Y_1)}{|Z_2 - Z_1|}\right]$$

297

□

298

299 **PROPOSITION 7.** *We define spacings  $\Delta X = X_2 - X_1$ ,  $\Delta Y = Y_2 - Y_1$ , and  $\Delta Z = Z_2 - Z_1 =$*   
300  *$\Delta X + \Delta Y$ . If the pairs  $(X_i, Y_i)$  are independently drawn from any distribution, then  $\Delta X$  and  $\Delta Y$*   
301 *are symmetric about zero, and time-integrated covariance is unconditionally non-positive:*

$$\int_0^\infty \sigma_{X,Y}(t) dt = \mathbb{E}\left[\frac{\Delta X \Delta Y}{|\Delta Z|}\right] \leq 0$$

**PROOF:** There is no need to assume that  $(\Delta X, \Delta Y)$  has a density. This proof also reveals that the result also holds for discrete random variables. Let  $\Delta X, \Delta Y$  be two real-valued random variables



such that:  $(-\Delta X, \Delta Y)$  has the same distribution as  $(\Delta X, \Delta Y)$ . We have:

$$\begin{aligned}
 \mathbb{E}[\Delta X \Delta Y / |\Delta X + \Delta Y|] &= \mathbb{E}[\mathbb{1}_{\Delta X \Delta Y > 0} \Delta X \Delta Y / |\Delta X + \Delta Y|] + \mathbb{E}[\mathbb{1}_{\Delta X \Delta Y < 0} \Delta X \Delta Y / |\Delta X + \Delta Y|] \\
 &= \mathbb{E}[\mathbb{1}_{\Delta X \Delta Y > 0} \Delta X \Delta Y / |\Delta X + \Delta Y|] + \mathbb{E}[\mathbb{1}_{-\Delta X \Delta Y < 0} (-\Delta X) \Delta Y / |\Delta Y - \Delta X|] \\
 &= \mathbb{E}[\mathbb{1}_{\Delta X \Delta Y > 0} \Delta X \Delta Y / |\Delta X + \Delta Y|] - \mathbb{E}[\mathbb{1}_{\Delta X \Delta Y > 0} \Delta X \Delta Y / |\Delta Y - \Delta X|] \\
 &= \mathbb{E}[\mathbb{1}_{\Delta X \Delta Y > 0} \Delta X \Delta Y (1/|\Delta X + \Delta Y| - 1/|\Delta Y - \Delta X|)] \\
 &\leq 0
 \end{aligned}$$

302 When  $\Delta X$  and  $\Delta Y$  have the same sign as imposed by the indicator function in the last expectation,  
 303 we have  $|\Delta X + \Delta Y| > |\Delta Y - \Delta X|$ , from which the inequality derives.  $\square$

304

305 **COROLLARY 1.** *Proposition 7 holds for divergent expectations.*

**PROOF:** Set  $U = |\Delta X|$  and  $V = |\Delta Y|$ ;  $M = \text{Max}(U, V)$ ,  $m = \text{Min}(U, V)$ . Then you can rewrite the expectation as:

$$\begin{aligned}
 \mathbb{E}[UV\{1/(U + V) - 1/(|U - V|)\}] &= \mathbb{E}[mM\{-2m/(M^2 - m^2)\}] \\
 &= -2\mathbb{E}[Mm^2/(M^2 - m^2)] \leq 0
 \end{aligned}$$

306 Indeed, if the expectation is divergent, then it is always  $-\infty$ . This approach removes the need to  
 307 make the argument that  $U + V > |U - V|$  and avoids the need to take a difference of expectations.  
 308 An alternative approach is given in an expanded statement and proof of Proposition 7 in the SI.  $\square$

309

## References

1. de Visser, J. A. G. M. & Elena, S. F. The evolution of sex: empirical insights into the roles of epistasis and drift. *Nat. Rev. Genet.* **8**, 139–149 (2007).
2. Otto, S. P. The evolutionary enigma of sex. *Am. Nat.* **174 Suppl 1**, S1–S14 (2009).
3. Otto, S. P. & Lenormand, T. Resolving the paradox of sex and recombination. *Nat. Rev. Genet.* **3**, 252–261 (2002).
4. Barton, N. H. & Charlesworth, B. Why sex and recombination? *Science* **281**, 1986–1990 (1998).
5. Bodmer, W. F. & Parsons, P. A. Linkage and recombination in evolution. In Caspari, E. W. & Thoday, J. M. (eds.) *Advances in Genetics*, vol. 11, 1–100 (Academic Press, 1963).
6. Nei, M. Modification of linkage intensity by natural selection. *Genetics* **57**, 625–641 (1967).
7. Hartfield, M. & Keightley, P. D. Current hypotheses for the evolution of sex and recombination. *Integr. Zool.* **7**, 192–209 (2012).
8. Felsenstein, J. The evolutionary advantage of recombination. *Genetics* **78**, 737–756 (1974).
9. West, S. A., Lively, C. M. & Read, A. F. A pluralist approach to sex and recombination. *J. Evol. Biol.* **12**, 1003–1012 (1999).
10. Otto, S. P. & Barton, N. H. Selection for recombination in small populations. *Evolution* **55**, 1921–1931 (2001).

11. Gerrish, P. J. & Lenski, R. E. The fate of competing beneficial mutations in an asexual population. *Genetica* **102-103**, 127–144 (1998).
12. Otto, S. P. & Barton, N. H. The evolution of recombination: removing the limits to natural selection. *Genetics* **147**, 879–906 (1997).
13. Barton, N. H. Linkage and the limits to natural selection. *Genetics* **140**, 821–841 (1995).
14. Agrawal, A. F. Evolution of sex: Why do organisms shuffle their genotypes? *Curr. Biol.* **16**, R696–R704 (2006).
15. Arjan, J. A. *et al.* Diminishing returns from mutation supply rate in asexual populations. *Science* **283**, 404–406 (1999).
16. Kondrashov, A. S. Deleterious mutations and the evolution of sexual reproduction. *Nature* **336**, 435–440 (1988).
17. Keightley, P. D. & Otto, S. P. Interference among deleterious mutations favours sex and recombination in finite populations. *Nature* **443**, 89–92 (2006).
18. Felsenstein, J. & Yokoyama, S. The evolutionary advantage of recombination. II. individual selection for recombination. *Genetics* **83**, 845–859 (1976).
19. Barton, N. H. A general model for the evolution of recombination. *Genet. Res.* **65**, 123–145 (1995).
20. Barton, N. H. Genetic linkage and natural selection. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **365**, 2559–2569 (2010).

21. Otto, S. P. & Feldman, M. W. Deleterious mutations, variable epistatic interactions, and the evolution of recombination. *Theor. Popul. Biol.* **51**, 134–147 (1997).
22. Slatkin, M. Linkage disequilibrium—understanding the evolutionary past and mapping the medical future. *Nat. Rev. Genet.* **9**, 477–485 (2008).
23. Baake, E., González Casanova, A., Probst, S. & Wakolbinger, A. Modelling and simulating lenski’s long-term evolution experiment. *Theor. Popul. Biol.* **127**, 58–74 (2019).
24. Fisher, R. A. *The genetical theory of natural selection* (Oxford Clarendon Press, 1930).
25. Muller, H. J. Some genetic aspects of sex. *Am. Nat.* **66**, 118–138 (1932).
26. Jaffe, K. Emergence and maintenance of sex among diploid organisms aided by assortative mating. *Acta Biotheor.* **48**, 137–147 (2000).
27. Redfield. A truly pluralistic view of sex and recombination. *J. Evol. Biol.* **12**, 1043–1046 (1999).
28. Yang, S. S. General distribution theory of the concomitants of order statistics. *Ann. Stat.* **5**, 996–1002 (1977).
29. David, H. A. Concomitants of extreme order statistics. In Galambos, J., Lechner, J. & Simiu, E. (eds.) *Extreme Value Theory and Applications: Proceedings of the Conference on Extreme Value Theory and Applications, Volume 1 Gaithersburg Maryland 1993*, 211–224 (Springer US, Boston, MA, 1994).

30. Martin, G., Otto, S. P. & Lenormand, T. Selection for recombination in structured populations. *Genetics* **172**, 593–609 (2006).
31. Becks, L. & Agrawal, A. F. Higher rates of sex evolve in spatially heterogeneous environments. *Nature* **468**, 89–92 (2010).
32. Hartfield, M., Otto, S. P. & Keightley, P. D. The maintenance of obligate sex in finite, structured populations subject to recurrent beneficial and deleterious mutation. *Evolution* **66**, 3658–3669 (2012).
33. Whitlock, A. O. B., Azevedo, R. B. R. & Burch, C. L. Population structure promotes the evolution of costly sex in artificial gene networks. *Evolution* **73**, 1089–1100 (2019).
34. Lenormand, T. & Otto, S. P. The evolution of recombination in a heterogeneous environment. *Genetics* **156**, 423–438 (2000).
35. Trevors, J. T. Hypothesized origin of microbial life in a prebiotic gel and the transition to a living biofilm and microbial mats. *C. R. Biol.* **334**, 269–272 (2011).
36. Flemming, H.-C. & Wuertz, S. Bacteria and archaea on earth and their abundance in biofilms. *Nat. Rev. Microbiol.* **17**, 247–260 (2019).
37. Miller, S. L., Schopf, J. W. & Lazcano, A. Oparin’s “origin of life”: Sixty years later. *J. Mol. Evol.* **44**, 351–353 (1997).
38. Kouyos, R. D., Otto, S. P. & Bonhoeffer, S. Effect of varying epistasis on the evolution of recombination. *Genetics* **173**, 589–597 (2006).

39. Otto, S. P. & Michalakis, Y. The evolution of recombination in changing environments. *Trends Ecol. Evol.* **13**, 145–151 (1998).
40. Gandon, S. & Otto, S. P. The evolution of sex and recombination in response to abiotic or coevolutionary fluctuations in epistasis. *Genetics* **175**, 1835–1853 (2007).
41. Peters, A. D. & Lively, C. M. The red queen and fluctuating epistasis: A population genetic analysis of antagonistic coevolution. *Am. Nat.* **154**, 393–405 (1999).
42. Barton, N. H. & Otto, S. P. Evolution of recombination due to random drift. *Genetics* **169**, 2353–2370 (2005).
43. Miller, J. B. & Sanjurjo, A. A bridge from monty hall to the hot hand: The principle of restricted choice. *J. Econ. Perspect.* **33**, 144–162 (2019).
44. Berkson, J. Limitations of the application of fourfold table analysis to hospital data. *Biometrics* **2**, 47–53 (1946).
45. Joag-Dev, K. & Proschan, F. Negative association of random variables with applications. *Ann. Statist.* **11**, 286–295 (1983).
46. Roze, D. & Barton, N. H. The Hill-Robertson effect and the evolution of recombination. *Genetics* **173**, 1793–1811 (2006).
47. Whitlock, A. O. B., Peck, K. M., Azevedo, R. B. R. & Burch, C. L. An evolving genetic architecture interacts with Hill–Robertson interference to determine the benefit of sex .

48. Hill, W. G. & Robertson, A. The effect of linkage on limits to artificial selection. *Genet. Res.* **8**, 269–294 (1966).

**Acknowledgements** The authors thank S. Otto and N. Barton for their thoughts on early stages of this work. Special thanks go to E. Baake for her thoughts on later stages of this work and help with key mathematical aspects. Much of this work was performed during a CNRS-funded visit (P.G.) to the Laboratoire Jean Kuntzmann, University of Grenoble Alpes, France, and during a visit to Bielefeld University (P.G.) funded by Deutsche Forschungsgemeinschaft (German Research Foundation, DFG) via Priority Programme SPP 1590 Probabilistic Structures in Evolution, grants BA 2469/5-2 and WA 967/4-2. The authors thank D. Chench, J. Streelman, R. Rosenzweig and the Biology Department at Georgia Institute of Technology for critical infrastructure and computational support. P.G. and A.C. received financial support from the USA/Brazil Fulbright scholar program. P.G. and P.S. received financial support from National Aeronautics and Space Administration grant NNA15BB04A.

**Author contributions** P.G. conceived the theory conceptually; P.G., P.S., B.S. and A.C. developed the theory verbally and with simulation; P.G, B.Y. and J.C. developed the theory mathematically; B.Y. and J.C. provided mathematical proofs for the across-population part; P.G., V.V., F.C. and N.H. provided mathematical proofs for the within-population part. P.G. wrote the paper with critical help and guidance from B.S., P.S. and B.Y.

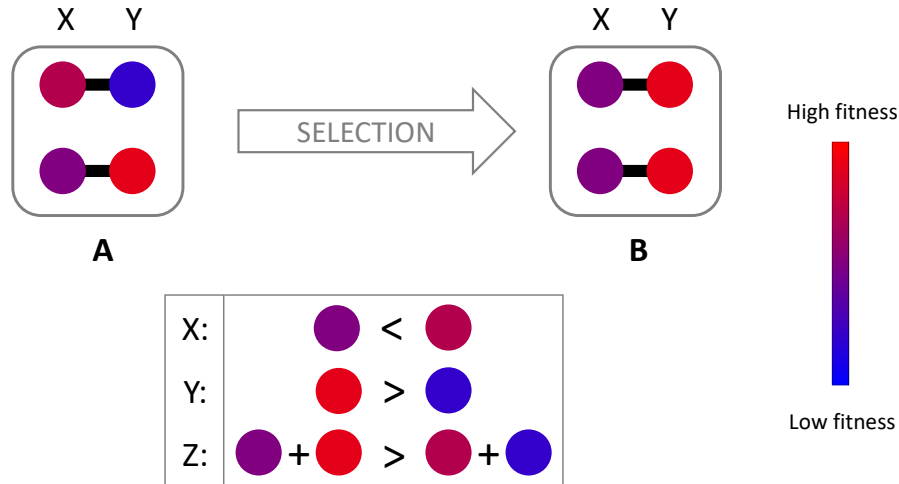
**Competing interests** The authors declare no competing interests.

## **Additional information**

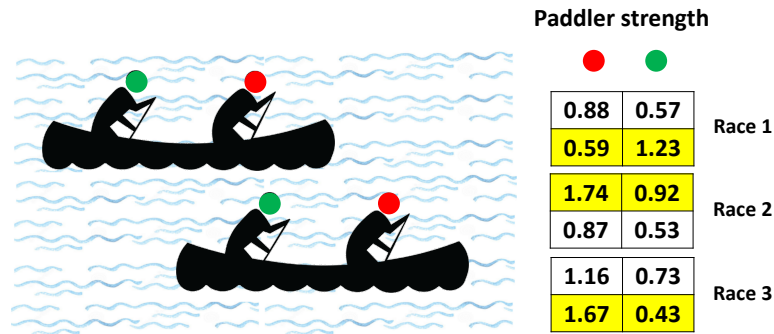
**Supplementary information** is available for this paper at <https://doi.org/10.1038/s....>

**Correspondence and requests for materials** should be addressed to P.G.

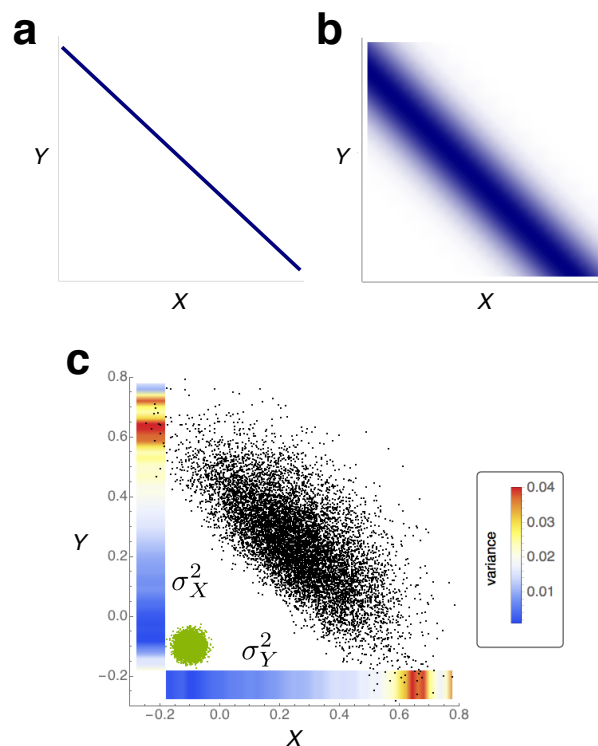




**Fig 1 | Two loci, two alleles.** Here, a large (infinite) population consists of individuals whose genome has only two loci  $x$  and  $y$ , each of which carries one of two alleles: genotype 1 encodes quantified phenotype  $X_1$  at the  $x$  locus and  $Y_1$  at the  $y$  locus, and genotype 2 carries quantified phenotype  $X_2$  at the  $x$  locus and  $Y_2$  at the  $y$  locus. Fitness is indicated by color. An individual's fitness is a function of the two phenotypes:  $Z = \phi(X, Y)$ ; here we make the simplifying assumption that  $\phi(X, Y) = X + Y$ , so that the fitnesses of genotypes 1 and 2 are  $Z_1 = X_1 + Y_1$  and  $Z_2 = X_2 + Y_2$ , respectively. The fitter of these two genotypes has total fitness denoted  $Z^{[2]}$  (i.e.,  $Z^{[2]} = \text{Max}\{Z_1, Z_2\}$ ) and genic fitnesses  $X_{(2)}$  and  $Y_{(2)}$  (i.e.,  $Z^{[2]} = X_{(2)} + Y_{(2)}$ ). Similarly, the less-fit of these two genotypes has total fitness  $Z^{[1]} = X_{(1)} + Y_{(1)}$ . We note:  $Z^{[2]} > Z^{[1]}$  by definition, but this does *not* guarantee that  $X_{(2)} > X_{(1)}$  or that  $Y_{(2)} > Y_{(1)}$ , as illustrated in the lower box. The population labeled  $A$  consists of two distinct genotypes but selection acts to remove the inferior genotype leaving a homogeneous population in which individuals are all genetically identical (with fitness  $Z^{[2]}$ ) as illustrated in the population labeled  $B$ . We derive selective mismatch measured by covariance  $\sigma_{XY}$ : 1) across populations (among different  $B$ ), given by  $\sigma_{X_{(2)}Y_{(2)}}$ , and 2) within populations (going from  $A$  to  $B$ ), given by  $\int_0^\infty \sigma_{XY}(t)dt$ .



**Fig 2 | Canoe race analogy.** Each canoe contains two paddlers. The strength of each paddler is measured and reported in the table. In any given canoe race, there is no correlation between paddler strengths  $A$  (green) and  $B$  (red). In each race, paddler strengths are recorded (tables on right), and the winning canoe is that in which the sum of the strengths of the two paddlers is the greatest (highlighted). Three such canoe races are conducted. We ask: what is the covariance between the strengths of paddlers  $A$  and  $B$  among winning canoes only? While it seems reasonable to suppose that winning canoes would carry two strong paddlers thereby resulting in positive covariance, the counter-intuitive answer we find is that the covariance is, for all practical purposes, unconditionally negative in expectation. By analogy, paddlers are genes, paddler strength is genic fitness, and canoes are genotypes. Natural selection picks the winner.



**Fig 3 | Natural selection promotes negative associations.** In the absence of recombination, selection does not act independently on  $X$  and  $Y$  but organismal fitness which, for simplicity, we here assume to be their sum,  $Z = \phi(X, Y) = X + Y$ . Perhaps counterintuitively, this fact alone creates negative associations. As discussed in the main text, this fact gives rise to a correlation of exactly negative one when the sum is a constant (**a**) and something intuitively negative when the sum is distributed as expected (**b**), i.e., as an order statistic. (**c**), Ten thousand simulated populations move from their initial (green dots) to final (black dots) mean fitnesses. Here, the predicted negative covariance in the final state is apparent. The heatmap bars indicate variance in  $Y$  along the  $x$ -axis and variance in  $X$  along the  $y$ -axis, a manifestation of Hill-Robertson interference<sup>8, 17, 30, 46–48</sup>: larger genic fitness at one locus relaxes selection on the other locus allowing for larger fitness variance at the that locus.