

1 **Recombination facilitates adaptive evolution in rhizobial soil bacteria**

2

3 Maria Izabel A. Cavassim^{1,2,+,*}, Stig U. Andersen², Thomas Bataillon¹, and Mikkel Heide
4 Schierup^{1,*}

5

6 **Author affiliations:**

7 ¹ Bioinformatics Research Centre, Aarhus University, Aarhus, 8000, Denmark

8 ² Department of Molecular Biology and Genetics, Aarhus University, Aarhus, 8000, Denmark

9 ⁺Present address: Department of Ecology and Evolutionary Biology, University of California, Los
10 Angeles, 90095, United States

11

12 ***Corresponding authors:**

13 Maria Izabel A. Cavassim, izabelcavassim@gmail.com

14 Mikkel H. Schierup, mheide@birc.au.dk

15

16 **Author Contributions**

17 **Conceptualization:** M.I.A.C., T.B., and M.H.S.; **Methodology:** M.I.A.C., T.B., and M.H.S.;

18 **Formal Analysis:** M.I.A.C.; **Investigation:** M.I.A.C., T.B., and M.H.S.; **Resources:** S.U.A., and

19 M.H.S.; **Data curation:** M.I.A.C., T.B., and M.H.S.; **Writing - Original Draft:** M.I.A.C.; **Writing**

20 **- Review and Editing:** M.I.A.C., S.U.A., T.B., and M.H.S.; **Visualization:** M.I.A.C., T.B. and

21 M.H.S. **Supervision:** T.B. and M.H.S.; **Project administration:** S.U.A. and M.H.S.; **Funding**

22 **acquisition:** S.U.A. and M.H.S.

23

24 **Competing interest**

25 The authors declare that they have no competing interests.

26

27 **Classification:** Biological Sciences and Evolution

28 **Keywords:** Adaptive evolution, rhizobium, recombination, beneficial mutations

29

30 **This PDF file includes:**

31 **Main Text**

32 **Figures 1 to 3**

33 **Tables 1 to 1**

34

35

36 **Abstract**

37 Homologous recombination is expected to increase natural selection efficacy by decoupling the
38 fate of beneficial and deleterious mutations and by readily creating new combinations of beneficial
39 alleles. Here, we investigate how the proportion of amino acid substitutions fixed by adaptive
40 evolution (α) depends on the recombination rate in bacteria. We analyze 3086 core protein-coding
41 sequences from 196 genomes belonging to five closely-related *Rhizobium leguminosarum*
42 species. We find that α varies from 0.07 to 0.39 across species and is positively correlated with
43 the level of recombination. We then evaluate the impact of recombination within each species by
44 dividing genes into three equally sized recombination classes based on their average level of
45 intragenic linkage disequilibrium. Generally, we found a significant increase in α with an increased
46 recombination rate. This is both due to a higher estimated rate of adaptive evolution and a lower
47 estimated rate of non-adaptive evolution, suggesting that recombination both increases the
48 fixation probability of advantageous variants and decreases the probability of fixation of
49 deleterious variants. Our results demonstrate that recombination facilitates adaptive evolution
50 not only in eukaryotes, but also in prokaryotes. Adaptive evolution could thus be a selective force
51 that universally promotes recombination.

52

53 **Significance statement**

54 Whether homologous recombination has a net beneficial or detrimental effect on adaptive
55 evolution is largely unexplored in natural bacterial populations. We address this question by
56 evaluating polymorphism and divergence data across 196 bacterial genome sequences of five
57 closely-related *Rhizobium leguminosarum* species. We show that the proportion of amino acid
58 changes fixed due to adaptive evolution (α) increases with an increased recombination rate. This
59 correlation is observed both in the interspecies and intraspecific comparisons. These results
60 suggest that homologous recombination directly impacts the efficacy of natural selection in
61 prokaryotes, as it has been shown previously to be in eukaryotes.

62 Main text

63 Introduction

64 Genetic recombination is expected to facilitate adaptive evolution by increasing the fixation
65 probability of adaptive mutations and decreasing the probability of fixation of deleterious
66 mutations (1). This is because recombination decouples the fate of adaptive and deleterious
67 variants, decreasing the amount of selective interference throughout the genome (2, 3). Selective
68 interference—also termed as the Hill-Robertson (HR) effect—is expected to be strongest in
69 regions of the genome where recombination is low (4). The HR effect is expected to cause: (i) a
70 reduction in the number of neutral polymorphisms, (ii) the accumulation of slightly deleterious
71 polymorphisms, and (iii) a decrease in the probability of fixation of advantageous alleles (see (5)).
72 Homologous recombination is expected to mitigate the HR effect and increase the percentage of
73 amino acid substitutions that are due to adaptive evolution (α).

74 Empirical evidence based on population genomics data supports these expectations with
75 a positive correlation between recombination and α reported in diverse species of eukaryotes,
76 including flies (*Drosophila melanogaster* (6, 7)), fungi (*Zymoseptoria tritici* (8)), plants
77 (*Arabidopsis thaliana* (9)), and non-model animal species (10, 11).

78 Whereas recombination is ubiquitous in eukaryotes, this is not the case for prokaryotes.
79 Nevertheless, most studied prokaryotes show high rates of genetic exchange (12) and it is
80 therefore of interest to explore whether such recombination also facilitates adaptive evolution in
81 prokaryotes. Here we investigate adaptive evolution and recombination in a species complex of
82 *Rhizobium leguminosarum* responsible for nitrogen fixation in white clover (*Trifolium repens*)
83 nodules. We have previously reported the full genomic sequence of 196 isolates (13). Our
84 analyses showed that they cluster into five closely related species (2-5% divergent) with horizontal
85 gene transfer only affecting the nitrogen fixation genes and a few well-defined genomic regions.
86 This species complex thus offers a unique opportunity among prokaryotes to estimate the rates
87 of fixation of amino acid changes by adaptive evolution from isolates sampled from natural
88 populations—enabling multiple comparisons of polymorphism and divergence patterns among
89 species. Our analyses provide evidence that the rate of adaptive protein evolution increases with
90 the recombination rate in this species complex.

91

92 Results & Discussion

93 To estimate the proportion of adaptive evolution (α) across this *Rhizobium* species
94 complex, and study how α covaries with recombination rate estimates, we restricted analyses to
95 polymorphism data from regions of the core genome without evidence of horizontal gene transfer
96 (HGT).

97 Across all five species (196 strains, gsA:32, gsB:32, gsC:112, gsD:5, gsE:11), a total of
98 22115 orthologous genes were previously identified (13); of those, 4204 genes are present in all
99 strains (core genes). Most core genes are found in the large chromosome (3304 genes), but some
100 were located in the chromids (Rh01, Rh02) and in one of the plasmids (Rh03) (see (14), and (13)).
101 The chromosome, chromids, and the plasmid are hereafter referred to as genomic compartments.
102 We filtered out genes that showed evidence of interspecies HGT or unexpectedly high rates of
103 nucleotide diversity (see Methods) (**Fig. S1**), leaving a total of 3086 genes (total alignment length:
104 3091179) and 334040 variable sites for analysis (**Fig. S2**).

105 First, we estimated nucleotide diversity, intragenic linkage disequilibrium (LD), and the site
106 frequency spectrum (SFS) (see Methods) within each species (**Fig. 1a-c**). The average nucleotide
107 diversity, π , an estimator of $2N_e\mu$ in haploids, is significantly different among genomic
108 compartments (**Fig. 1a, Table S1**). Across the species, π differs by up to a factor of 4.5 (gsA:
109 0.018, gsB:0.0045, gsC:0.0140, gsD:0.00512, gsE:0.008), with the most polymorphic species
110 being gsA and the least gsB. If we assume similar mutation rates among these closely related
111 species, nucleotide diversity differences reflect interspecies differences in long-term effective
112 population size, N_e .

113 When recombination occurs, we expect that levels of non-random association between
114 pairs of alleles, quantified by r^2 (see Methods), decay with genomic distance (LD decay). To
115 evaluate the recombination rate differences among the five species, we used within-species
116 polymorphism data and computed the average intragenic LD decay for each gene in each
117 species. We observed a rapid decay of LD within the first 1000 base pairs for all species,
118 suggesting substantial amounts of within-species homologous recombination (**Fig. 1b**). The
119 slower decay observed in species gsB either reflects a lower per generation recombination rate
120 or a smaller effective population size (N_e). The latter is consistent with the low level of nucleotide
121 diversity measured in gsB. To reliably estimate interspecies differences in r^2 , we used genes with
122 at least ten informative sites within each species and evaluated their r^2 distributions separately
123 (**Fig. 1c**). As expected, the species with the most striking LD decay (gsC) has the lowest r^2
124 (median r^2 : 0.248), and the opposite is also true (gsD, median r^2 : 0.7131).

125 Species gsB and gsE have similar LD decay and r^2 distributions (**Fig. 1b-c**). We tested
126 whether their average r^2 was statistically different. Based on the total number of shared genes
127 (3086), the average r^2 of gsB was not significantly different from the average r^2 of gsE (Wilcox's
128 test, p-value=0.10). All the other pairwise comparisons were statistically significant (**Table S2**).
129 Thus, we considered gsB and gsE to have similar recombination rates. In summary, these species
130 can be ranked by their recombination levels, from the most recombining to the least, as follows:
131 gsC (median r^2 : 0.248) > gsA (median r^2 : 0.293) > gsB (median r^2 : 0.48) and gsE (median r^2 :
132 0.43) > gsD (median r^2 : 0.71).

133 Next, we computed the folded site frequency spectrum (SFS) of synonymous and
134 nonsynonymous mutations within each species. Overall, both synonymous and nonsynonymous
135 SFSs differ from the "L" shaped patterns (many rare alleles and fewer frequent alleles) expected
136 in a stationary population at mutation-selection-drift equilibrium (**Fig. 1d**). The observed excess
137 of intermediate frequency SNPs indicates the presence of population structure in some of the
138 species. The effect of population structure is particularly evident in gsC, and this excess is likely
139 driven by strains isolated from French soils (**Fig. S3**). Differences among species suggest distinct
140 demographic histories, with gsC showing an SFS compatible with population expansion and gsA
141 with population decline.

142 Using the counts of polymorphism in synonymous and nonsynonymous SFS within each
143 species, we can estimate the overall strength of purifying selection via p_iN/p_iS . The strength of
144 purifying selection ranks species similarly to their average recombination rate, with more
145 recombining species showing stronger purifying selection (individual p_iN/p_iS are gsA = 0.039,
146 gsB = 0.057, gsC = 0.037, gsD = 0.07, gsE = 0.051). This observation is in line with the theoretical
147 expectation of a positive effect of recombination on the overall efficacy of natural selection. We
148 also observed an excess of rare nonsynonymous relative to synonymous variants (**Fig. 1e**),
149 consistent with the segregation of nonsynonymous variants under weak purifying selection (15).
150 Rare nonsynonymous variants are often deleterious ($s \sim 1/N_e$) (16, 17). Because deleterious
151 variants contribute substantially to polymorphism but rarely to divergence (18, 19), their presence
152 in the genomes, if not controlled for, will lead to an underestimation of α (20).

153 We used GRAPES (10) to estimate the distribution of fitness effects (DFE) (21) and the
154 proportion of adaptive evolution (α) from polymorphism and divergence data while accounting for
155 the presence of deleterious mutations. This approach uses the site frequency distribution of both
156 synonymous and nonsynonymous SFS counts to estimate the DFE while accounting for the effect
157 of demography. The significant amount of shared polymorphism among species (**Table S3**)
158 makes it difficult to reliably call ancestral and derived states (22, 23). Accordingly, we chose to

159 estimate the DFE and α using the folded SFSs (10). To determine the model of the DFE that best
160 fit our data, we used a variety of DFE distribution models (**Table S4**). The DFE models we tested
161 differ by the classes of mutations (deleterious, beneficial, and neutral) included in each DFE
162 model and how fitness effects are distributed within these classes. When using Akaike's
163 Information Criterion (AIC) to select the best DFE model, we found the GammaZero model overall
164 provides the best fit to the SFS data (**Fig. S4**). This model assumes the existence of weakly
165 deleterious nonsynonymous mutations, modeled as a continuous Gamma distribution (10).

166 The proportion of adaptive evolution was first computed between all combinations of
167 "mirror" species ($\alpha_{species1\ species2}$), in which "species 2" is used as outgroup (divergence) for
168 "species 1" (polymorphism), and vice-versa ($\alpha_{species2\ species1}$). This yielded twenty combinations
169 in total. Because "mirror" species share an identical history of divergence, their α estimates can
170 be considered as "biological replicates" (10) (**Table 1**). Except for the comparison between gsA
171 and gsB, in which differences between $\alpha_{gsA\ gsB}$ and $\alpha_{gsB\ gsA}$ exceeded 0.1, the overall
172 discrepancy in the values estimated between mirror species does not exceed 0.1. Using each
173 species' focal polymorphism data, we calculated four α estimates by comparing it to the
174 divergence counts of the remaining species (**Table 1**). The most recombining species (gsC) is
175 observed to have the highest α across all outgroups used, while the least recombining species
176 (gsD) had the lowest α in $\frac{3}{4}$ of the cases.

177 We then investigated whether intraspecies differences in recombination rate affect the
178 amount of adaptive evolution (α). For each species, we split genes into three recombination
179 classes based on their average r^2 values and computed α for each class using the GammaZero
180 model (**Fig. S5, Table S5**). Because we only kept genes with at least ten informative sites, the
181 number of genes evaluated across species was different (**Fig. 1c**). For most species comparisons
182 (gsA, gsB, gsC, and gsE), there is a decrease in the proportion of adaptive evolution with a
183 reduction in recombination (increase in r^2). Except for cases in which we used gsD and gsB
184 polymorphisms to estimate α , all the other species pairwise comparisons led to at least one
185 significant difference (based on non-overlapping CI's) between recombination classes. We further
186 assessed the significance of the pattern reported here by permuting—200 times—across
187 recombination classes (see Methods). Except for simulations in which gsD and gsB
188 polymorphisms were used, all the other simulations led to significant differences (p-value \leq
189 0.025) among the two most extreme classes of recombination (**Fig. S6**).

190 The parameter α can also be viewed as the relative proportion between the rate of amino
191 acid changes fixed by positive selection (ω_a) and the rate of non-adaptive amino acid changes

192 (ω_{na}): $\alpha = \omega_a / (\omega_a + \omega_{na})$. Thus, an increase in α with recombination could be due to either an
193 increase in the rate of adaptive substitutions, a decrease in the rate of non-adaptive substitutions,
194 or both. Figure 3 shows that ω_a increases with recombination rate whereas ω_{na} decreases with
195 recombination rate for most combinations and that the quantitative effects are almost equal in
196 magnitude. Thus, classes of genes evolving under higher recombination rates exhibited lower
197 rates of non-adaptive substitution and increased rates of fixation of adaptive variation. This is
198 exactly as predicted from selective interference theory (1, 3, 4).

199 To evaluate the robustness of these results, we computed an alternative measure of
200 recombination (R/θ). We then made new recombination classes and evaluated its correlation with
201 α . R/θ measures the importance of recombination (R), relative to mutation (θ) across sequences
202 (24, 25). Although the per gene estimates of R/θ are less variable than that of r^2 (**Fig. S7**), these
203 two measures are not independent (Pearson correlation ranged from 0.21 to 0.44) (**Fig. S7**). For
204 most species comparisons, the trend between α and recombination is still consistent: the higher
205 R/θ , the higher α is (**Fig. S8**).

206 Conclusion

207 We have found that five bacterial species within the species complex *R. leguminosarum*
208 display different yet high levels of recombination. The estimates of α ranged between 0.07 and
209 0.39 among species. These estimates are lower than those based on 410 orthologs observed in
210 *E. coli* (0.58, CI=0.45, 0.68) but close to earlier estimates from *S. enterica* (0.34, CI=0.14, 0.50)
211 previously reported (19).

212 Levels of recombination correlate—both across and within species—with higher amounts
213 of adaptive evolution measured either as the rate of adaptive substitutions (ω_a) or as the
214 proportion of amino acid changes which have been fixed by positive selection (α). For instance,
215 the most recombining species (gsC) consistently exhibited the largest α , independent of the
216 outgroup used. Within each species, we also find a positive correlation between intragenomic
217 recombination rate and α , as well as for ω_a . These findings are robust to the measure of
218 recombination (r^2 and R/θ) and the choice of outgroup used for computing divergence.

219 The positive effect of homologous recombination on α we report here is in line with
220 population genetic studies conducted in vertebrates (10, 11) and invertebrates (7, 8, 26–29). It
221 points to recombination being a general facilitator of adaptive evolution across the tree of life—
222 possibly being a selective force for the existence of recombination in prokaryotes in the first place.

223

224 **Material and methods**

225 **Identification of orthologous genes**

226 We previously isolated and sequenced 196 strains from white clover (*Trifolium repens*)
227 root nodules harvested in Denmark, France, and the UK. To identify a set of orthologous genes
228 shared across strains, we followed the methods outlined in Cavassim et al., 2020 (13). Briefly,
229 the strains were previously subjected to whole-genome shotgun sequencing using 2x250bp
230 Illumina paired-end reads (Illumina, USA). Genomes were assembled using SPAdes (v. 3.6.2,
231 (30)) and assembled further, one strain at a time, using a custom Python script (Jigome, available
232 at https://github.com/izabelcavassim/Rhizobium_analysis/tree/master/Jigome).

233 From the assembled genomes (13), we predicted protein-coding sequences using prokka
234 (31) (v1.12); this resulted in a total of 1468264 protein-coding sequences. To predict orthologous
235 genes from these sequences, we used Proteinortho (31, 32) (v5.16b) with default parameters
236 except for enabling the synteny flag. We identified a total of 22115 orthologous, including a total
237 of 17911 orthologous observed in at least two strains (accessory genes), and 4204 orthologous
238 found in all 196 strains (core genes).

239 Genes were then aligned using clustalo (33) (v.1.2.0) in a codon-aware manner. To
240 determine the genetic relationship among all 196 strains, we previously calculated their pairwise
241 average nucleotide identity (ANI) across 305 conserved orthologous gene alignments (13). Under
242 the 95% ANI threshold that delineates species boundaries (34), we demonstrated that these 196
243 *Rhizobium* strains constitute five distinct *R. leguminosarum* species (gsA, gsB, gsC, gsD, and
244 gsE). To ensure that we had a high-quality orthologous dataset for extracting segregating sites,
245 we filtered it further (see below).

246

247 **Confident core orthologous genes**

248 By developing and applying a phylogenetic method to quantify HGT (introgression score),
249 we previously showed that most of the core genes shared among the present species respect the
250 species-tree topology (introgression score = 0) (13). The exceptions are genes sitting in the
251 symbiosis conjugative plasmids, and two chromosomal islands (introgression score > 7). To
252 ensure that we were only analyzing high-quality gene alignments with little evidence of HGT, we
253 imposed some restrictions. We only accepted genes that passed the following criteria: (i) were
254 present in every strain (196 strains), (ii) with a nucleotide diversity (π) below 0.1, (iii) identifiable
255 replicon origin (chromosome and chromids), (iiii) and with an introgression score ≤ 3 (**Fig. S1a**).
256 A total of 3086 out of 4204 core genes were kept, and of these, 2550 genes were found in the

257 chromosome, 288 genes in chromid Rh01, 160 genes in chromid Rh02, and 88 genes in plasmid
258 Rh03.

259

260 **Variant calling**

261 To identify single nucleotide polymorphisms (SNPs) along with our high-quality set of core
262 genes, we evaluated each gene codon-aware alignment using a custom python script
263 https://github.com/izabelcavassim/Popgen_bacteria. For a given core gene alignment and
264 position, we first counted the number of unique nucleotides (A, C, T, G). Only sites containing two
265 unique nucleotides were considered variable sites (bi-allelic SNPs). SNP matrices were then built
266 and encoded as follows: major alleles were encoded as 1 and minor alleles as 0. Gaps were
267 replaced by the site mean. The nucleotide diversity (π), gene length, and the distributions of
268 segregating sites across core genes are described in **Fig. S1 (b-d)**.

269 **Transition transversion rate bias (kappa) and synonymous and** 270 **nonsynonymous counts**

271 Because transitions are more often synonymous at third codon positions than are
272 transversions, to correctly identify the expected number of synonymous (Lps) and
273 nonsynonymous counts (Lpn), we first estimated the average transition/transversion rate bias
274 (kappa) (35) across species. To this end, we followed the methods described in (36) and used
275 two classes of sites: fourfold-degenerate sites at the third codon positions and nondegenerate
276 sites. Mutations at the fourfold-degenerate sites are synonymous, and therefore kappa at those
277 sites should reflect only the mutational bias. All mutations at nondegenerate sites are
278 nonsynonymous and were also used to estimate kappa. We computed an average kappa by
279 combining these two classes based on equations 8, 9, 10, and 11 of Yang and Nielsen 2000 (36).
280 These equations have been implemented within the CodonSeq class in Biopython (37) (private
281 function `_count_site_YN00()`), and these private functions were adapted to our dataset.

282 To estimate a common kappa for each gene alignment (including all species and strains),
283 we averaged estimates from pairwise analyses across 50 randomly chosen strains. The kappa
284 distribution has a mean of 5.6 and a median of 5.20 (**Fig. S2a**), we used the median to compute
285 the expected number of synonymous and nonsynonymous sites. To this end, we followed the
286 methods described by Ina (1995) (35) and modified by Yang and Nielsen 2000 (36)—also
287 implemented within Biopython. A total of 284742 synonymous, 49298 nonsynonymous sites were
288 counted (**Fig. S2b-c**).

289

290 **Divergence sites and shared polymorphisms**

291 For each pair of species (a focal and an outgroup), we evaluated their variable sites and
292 computed the number of shared synonymous (pS) and nonsynonymous (pN) polymorphisms.
293 Given a bi-allelic SNP (0 and 1), we considered shared polymorphic sites as sites for which both
294 alleles (0,1) were segregating in both species (**Table S3**). We restricted the estimates of
295 divergence to those sites for which we had variable sites across species. We classified
296 synonymous (d_S) and nonsynonymous divergent sites (d_N) as those sites in which we observed
297 fixed differences between a focal species and an outgroup.

298

299 **Calculating the folded Site Frequency Spectrum**

300 One can infer the distribution of fitness effects from Site Frequency Spectrum (SFS) data
301 (20). Because of the amount of shared polymorphism among the present species (**Table S3**), it
302 becomes problematic to confidently distinguish ancestral from derived polymorphisms (38).
303 Therefore, we chose to estimate the DFE using a method that uses the folded site frequency
304 spectrums (SFS) of synonymous and nonsynonymous sites (10). To this end, we built the folded
305 synonymous and nonsynonymous site frequency spectrums by tabulating the observed counts of
306 the minor allele frequencies. The synonymous and nonsynonymous SFSs, and the divergence
307 counts, were then used to estimate the DFE and the proportion of adaptive substitutions (α)
308 across pairs of species.

309

310 **Calculating the strength of purifying selection**

311 The strength of purifying selection was measured as the ratio of nucleotide diversity at
312 nonsynonymous (pi_N) and synonymous sites (pi_S). For each gene and class of polymorphisms
313 (synonymous and nonsynonymous) nucleotide diversity was computed as: $\pi = \sum_1^m (2pq)/Lp$, in
314 which p and q are the allele frequencies, and Lp is the expected number of synonymous (Lps) or
315 nonsynonymous positions (Lpn) along the gene. We use the median of the pi_N/pi_S distribution
316 among genes as a proxy for the strength of purifying selection per species.

317

318 **Calculating the significance levels between recombination classes**

319 To test whether differences among recombination classes were statistically significant
320 across species comparisons, we conducted a non-parametric test by shuffling genes among
321 recombination classes (200 permutations) and recording the amplitude of differences between α

322 estimates ($\Delta_\alpha = \max_\alpha - \min_\alpha$). We calculated a p-value by comparing the observed Δ_α against
323 the simulated Δ_α distribution.

324

325 **Estimation of adaptive substitutions (α)**

326 Fitted parameters of the DFE were used to compute the expected d_N/d_S under the
327 different models, which was compared to the observed d_N/d_S to estimate the adaptive substitution
328 rate (ω_A); non-adaptive substitution rate (ω_{NA}), and the proportion of adaptive substitutions (α)
329 with $\omega_A = \alpha d_N/d_S$ and $\omega_{NA} = (1 - \alpha) d_N/d_S$.

330 To account for potential departures of the SFS from demographic equilibrium (assuming
331 the Wright-Fisher model)—possibly driven by changes in the effective population size or by
332 population structure—the method uses nuisance parameters to correct for these SFS distortions
333 (39). The different DFE models were compared using the Akaike Information Criterion (AIC) (40).

334

335 **Recombination rate estimates**

336 To estimate the recombination rate per gene per species, we used two approaches: one
337 based on the degree of association (or linkage disequilibrium) between pairs of alleles in a sample
338 of haplotypes (r^2), and another, ClonalFrameML (R/θ) (24, 25), which relies on the maximum
339 likelihood inference to detect recombination events that disrupt a clonal pattern of inheritance in
340 bacterial genomes.

341 **(1) Intragenic linkage disequilibrium**

342 Intragenic linkage disequilibrium (LD) measures the correlation between pairs of alleles
343 with genomic distance. Here we used Pearson's r^2 correlation measure.

344 Each gene genotype matrix (containing a minimal set of ten single nucleotide
345 polymorphisms (SNPs)) was first normalized as follow: let N denote the total number of
346 individuals, and M the total number of SNPs, the full gene genotype matrix (X) has dimensions
347 $N \times M$ with genotypes encoded as 0's and 1's for the N haploid individuals. Each column S_i ($i =$
348 $1, \dots, M$) of the X matrix is a vector of SNP information of size N . The first step of the calculation
349 was to apply a Z-score normalization to each SNP vector by subtracting the vector by its mean
350 and dividing it by its standard deviation ($\frac{S_i - \bar{S}_i}{\sqrt{\text{var}(S_i)}}$), resulting in a vector with mean 0 and variance

351 1. The linkage disequilibrium was then calculated as a function of distance d (maximum 1000
352 base pairs apart) and was computed as the average LD of pairs of SNPs d base pairs away from
353 each other. The calculations were done in the following way:

354
$$Cor(X_i, X_j) = \frac{Cov(X_i, X_j)}{\sqrt{var(X_i)var(X_j)}}$$

355
$$r^2 = Cor(X_i, X_j)^2$$

356

357 In which $j > i$ and X_i is composed of the genotypes of all individuals of a given species for SNP
358 position i in the genotype matrix. X_j is formed of the genotypes of all individuals of the same
359 species for position j in the genotype matrix, and $d = j - i$ with $d \leq 1000$ base pairs. Results
360 were then summarized into bins of 100 base pairs apart; for each bin, a mean r^2 was computed
361 and then averaged to a singular r^2 value.

362 (2) ClonalFrameML

363 To estimate the changes in the clonal phylogeny by recombination (R), relative to mutation
364 (θ) (R/θ), we used the software ClonalFrameML (24, 25). For each species, we first concatenated
365 all core gene alignments (3086 genes) to build the starting phylogenetic species-tree using a
366 maximum likelihood approach (Raxml-ng (41, 42)). We then input each phylogenetic tree within
367 each gene alignment to estimate R/θ .

368

369 Data sharing plans

370 Code generated for this study can be found at
371 https://github.com/izabelcavassim/Popgen_bacteria

372

373 Funding information

374 This work was funded by grant no. 4105-00007A from Innovation Fund Denmark (S.U.A. and
375 M.H.S.).

376

377 Acknowledgments

378 The authors would like to thank industrial partners DLF Trifolium, SEGES, and Legume
379 Technology Ltd. for their contribution to the field trials. The authors would also like to thank J.
380 Peter W. Young, Paula Tataru, Bjarni Vilhjálmsson, and Marjolaine Rousselle for their helpful
381 discussions about this work.

382
383
384

385 References

386

- 387 1. N. H. Barton, The reduction in fixation probability caused by substitutions at linked loci. *Genetical*
388 *Research* **64**, 199–208 (1994).
- 389 2. W. G. Hill, A. Robertson, The effect of linkage on limits to artificial selection. *Genetical Research* **8**,
390 269–294 (1966).
- 391 3. J. Felsenstein, The evolutionary advantage of recombination. *Genetics* **78**, 737–756 (1974).
- 392 4. G. A. McVean, B. Charlesworth, The effects of Hill-Robertson interference between weakly selected
393 mutations on patterns of molecular evolution and variation. *Genetics* **155**, 929–944 (2000).
- 394 5. B. Charlesworth, A. J. Betancourt, V. B. Kaiser, I. Gordo, Genetic recombination and molecular
395 evolution. *Cold Spring Harb. Symp. Quant. Biol.* **74**, 177–186 (2009).
- 396 6. D. Castellano, M. Coronado-Zamora, J. L. Campos, A. Barbadilla, A. Eyre-Walker, Adaptive
397 Evolution Is Substantially Impeded by Hill-Robertson Interference in *Drosophila*. *Mol. Biol. Evol.* **33**,
398 442–455 (2016).
- 399 7. J. L. Campos, D. L. Halligan, P. R. Haddrill, B. Charlesworth, The relation between recombination
400 rate and patterns of molecular evolution and variation in *Drosophila melanogaster*. *Mol. Biol. Evol.*
401 **31**, 1010–1028 (2014).
- 402 8. J. Grandaubert, J. Y. Dutheil, E. H. Stukenbrock, The genomic determinants of adaptive evolution in
403 a fungal pathogen. *Evol Lett* **3**, 299–312 (2019).
- 404 9. A. F. Moutinho, F. F. Trancoso, J. Y. Dutheil, The Impact of Protein Architecture on Adaptive
405 Evolution. *Mol. Biol. Evol.* **36**, 2013–2028 (2019).
- 406 10. N. Galtier, Adaptive Protein Evolution in Animals and the Effective Population Size Hypothesis. *PLoS*
407 *Genet.* **12**, e1005774 (2016).
- 408 11. A. F. Moutinho, T. Bataillon, J. Y. Dutheil, Variation of the adaptive substitution rate between species
409 and within genomes. *Evolutionary Ecology* **34**, 315–338 (2020).
- 410 12. X. Didelot, M. C. J. Maiden, Impact of recombination on bacterial evolution. *Trends Microbiol.* **18**,
411 315–322 (2010).
- 412 13. M. I. A. Cavassim, *et al.*, Symbiosis genes show a unique pattern of introgression and selection
413 within a *Rhizobium leguminosarum* species complex. *Microbial Genomics* **6** (2020).
- 414 14. Harrison, *et al.*, Introducing the bacterial “chromid”: not a chromosome, not a plasmid. *Trends*
415 *Microbiol.* **18**, 141–148 (2010).
- 416 15. T. Ohta, Role of very slightly deleterious mutations in molecular evolution and polymorphism. *Theor.*
417 *Popul. Biol.* **10**, 254–275 (1976).
- 418 16. A. L. Hughes, *et al.*, Widespread purifying selection at polymorphic sites in human protein-coding
419 loci. *Proc. Natl. Acad. Sci. U. S. A.* **100**, 15754–15757 (2003).
- 420 17. A. L. Hughes, Evidence for abundant slightly deleterious polymorphisms in bacterial populations.
421 *Genetics* **169**, 533–538 (2005).
- 422 18. J. C. Fay, G. J. Wyckoff, C. I. Wu, Positive and negative selection on the human genome. *Genetics*
423 **158**, 1227–1234 (2001).

- 424 19. J. Charlesworth, A. Eyre-Walker, The rate of adaptive evolution in enteric bacteria. *Mol. Biol. Evol.*
425 **23**, 1348–1356 (2006).
- 426 20. A. Eyre-Walker, P. D. Keightley, Estimating the rate of adaptive molecular evolution in the presence
427 of slightly deleterious mutations and population size change. *Mol. Biol. Evol.* **26**, 2097–2108 (2009).
- 428 21. A. Eyre-Walker, P. D. Keightley, The distribution of fitness effects of new mutations. *Nature Reviews*
429 *Genetics* **8**, 610–618 (2007).
- 430 22. P. Tataru, M. Mollion, S. Glémin, T. Bataillon, Inference of Distribution of Fitness Effects and
431 Proportion of Adaptive Substitutions from Polymorphism Data. *Genetics* **207**, 1103–1119 (2017).
- 432 23. A. Schneider, B. Charlesworth, A. Eyre-Walker, P. D. Keightley, A method for inferring the rate of
433 occurrence and fitness effects of advantageous mutations. *Genetics* **189**, 1427–1437 (2011).
- 434 24. M. Vos, X. Didelot, A comparison of homologous recombination rates in bacteria and archaea. *ISME*
435 *J.* **3**, 199–208 (2009).
- 436 25. X. Didelot, D. J. Wilson, ClonalFrameML: efficient inference of recombination in whole bacterial
437 genomes. *PLoS Comput. Biol.* **11**, e1004041 (2015).
- 438 26. D. C. Presgraves, Recombination enhances protein adaptation in *Drosophila melanogaster*. *Curr.*
439 *Biol.* **15**, 1651–1656 (2005).
- 440 27. A. J. Betancourt, J. J. Welch, B. Charlesworth, Reduced effectiveness of selection caused by a lack
441 of recombination. *Curr. Biol.* **19**, 655–660 (2009).
- 442 28. J. R. Arguello, *et al.*, Recombination Yet Inefficient Selection along the *Drosophila melanogaster*
443 Subgroup's Fourth Chromosome. *Molecular Biology and Evolution* **27**, 848–861 (2010).
- 444 29. T. F. C. Mackay, *et al.*, The *Drosophila melanogaster* Genetic Reference Panel. *Nature* **482**, 173–
445 178 (2012).
- 446 30. A. Bankevich, *et al.*, SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell
447 Sequencing. *Journal of Computational Biology* **19**, 455–477 (2012).
- 448 31. T. Seemann, Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).
- 449 32. M. Lechner, *et al.*, Orthology detection combining clustering and synteny for very large datasets.
450 *PLoS One* **9**, e105015 (2014).
- 451 33. F. Sievers, *et al.*, Fast, scalable generation of high-quality protein multiple sequence alignments
452 using Clustal Omega. *Mol. Syst. Biol.* **7**, 539 (2011).
- 453 34. K. T. Konstantinidis, A. Ramette, J. M. Tiedje, The bacterial species definition in the genomic era.
454 *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **361**, 1929–1940 (2006).
- 455 35. Y. Ina, New methods for estimating the numbers of synonymous and nonsynonymous substitutions.
456 *J. Mol. Evol.* **40**, 190–226 (1995).
- 457 36. Z. Yang, R. Nielsen, Estimating synonymous and nonsynonymous substitution rates under realistic
458 evolutionary models. *Mol. Biol. Evol.* **17**, 32–43 (2000).
- 459 37. P. J. A. Cock, *et al.*, Biopython: freely available Python tools for computational molecular biology and
460 bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- 461 38. R. D. Hernandez, S. H. Williamson, C. D. Bustamante, Context dependence, ancestral
462 misidentification, and spurious signatures of natural selection. *Mol. Biol. Evol.* **24**, 1792–1800 (2007).

- 463 39. A. Eyre-Walker, M. Woolfit, T. Phelps, The distribution of fitness effects of new deleterious amino
464 acid mutations in humans. *Genetics* **173**, 891–900 (2006).
- 465 40. H. Akaike, Information Theory and an Extension of the Maximum Likelihood Principle. *Springer*
466 *Series in Statistics*, 610–624 (1992).
- 467 41. K. Kobert, T. Flouri, A. Aberer, A. Stamatakis, The Divisible Load Balance Problem and Its
468 Application to Phylogenetic Inference. *Lecture Notes in Computer Science*, 204–216 (2014).
- 469 42. A. Stamatakis, RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
470 phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

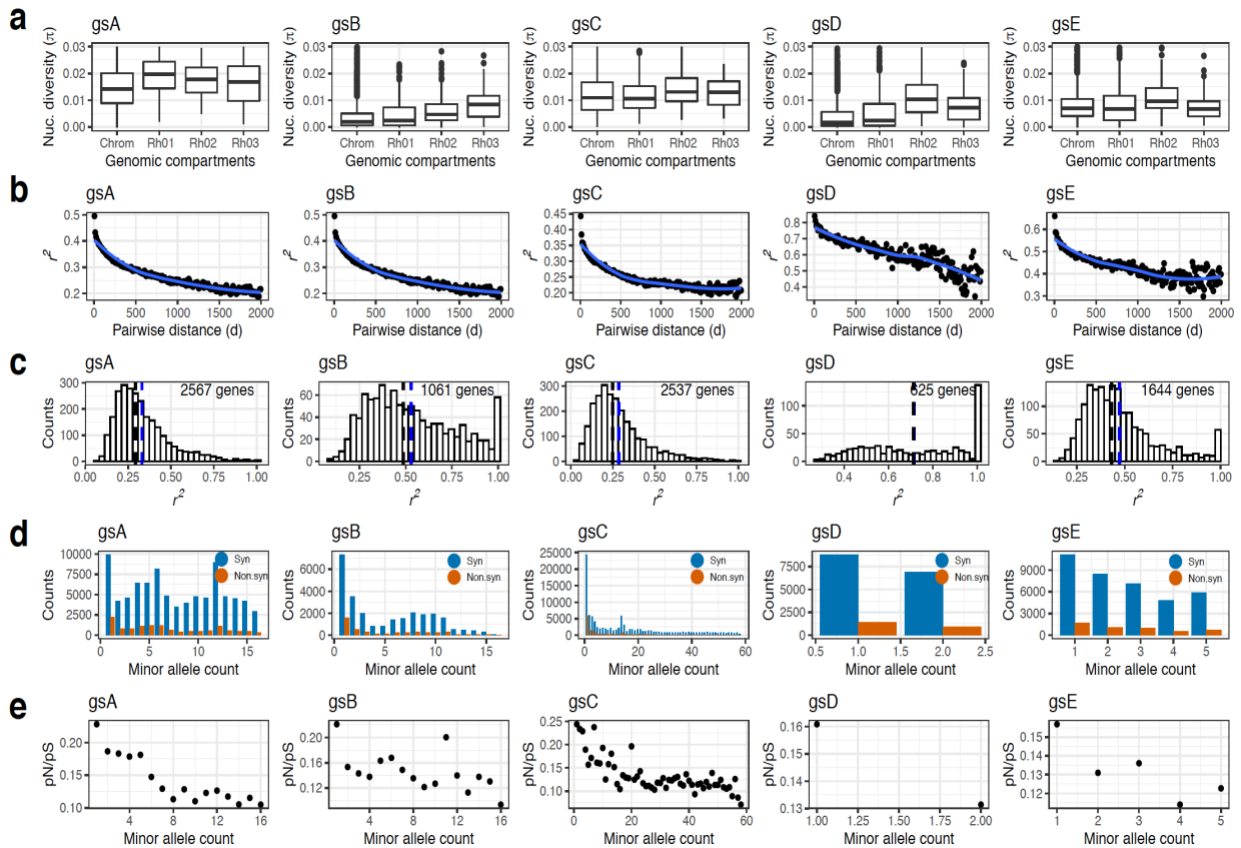
492

493

494

495 **Figures and Tables**

496



497

498

499 **Figure 1. Population genetics parameters across five species. (a)** Nucleotide diversity (π)

500 across 3086 genes distributed along with genomic compartments (chromosome and chromids:

501 Rh01, Rh02, and Rh03). To exclude outliers only genes with $\pi \leq 0.03$ are shown. **(b)** Intragenic

502 linkage disequilibrium measured via the decay of r^2 for all core genes (3086 genes). The curve

503 fitting line (in blue) is from a local regression method (loess). **(c)** Linkage disequilibrium (r^2)

504 distribution across genes. Only genes with at least ten segregating sites were kept (gsA: 2567

505 genes, gsB: 1061, gsC: 2537, gsD: 625, and gsE: 1644). The black and blue dashed lines

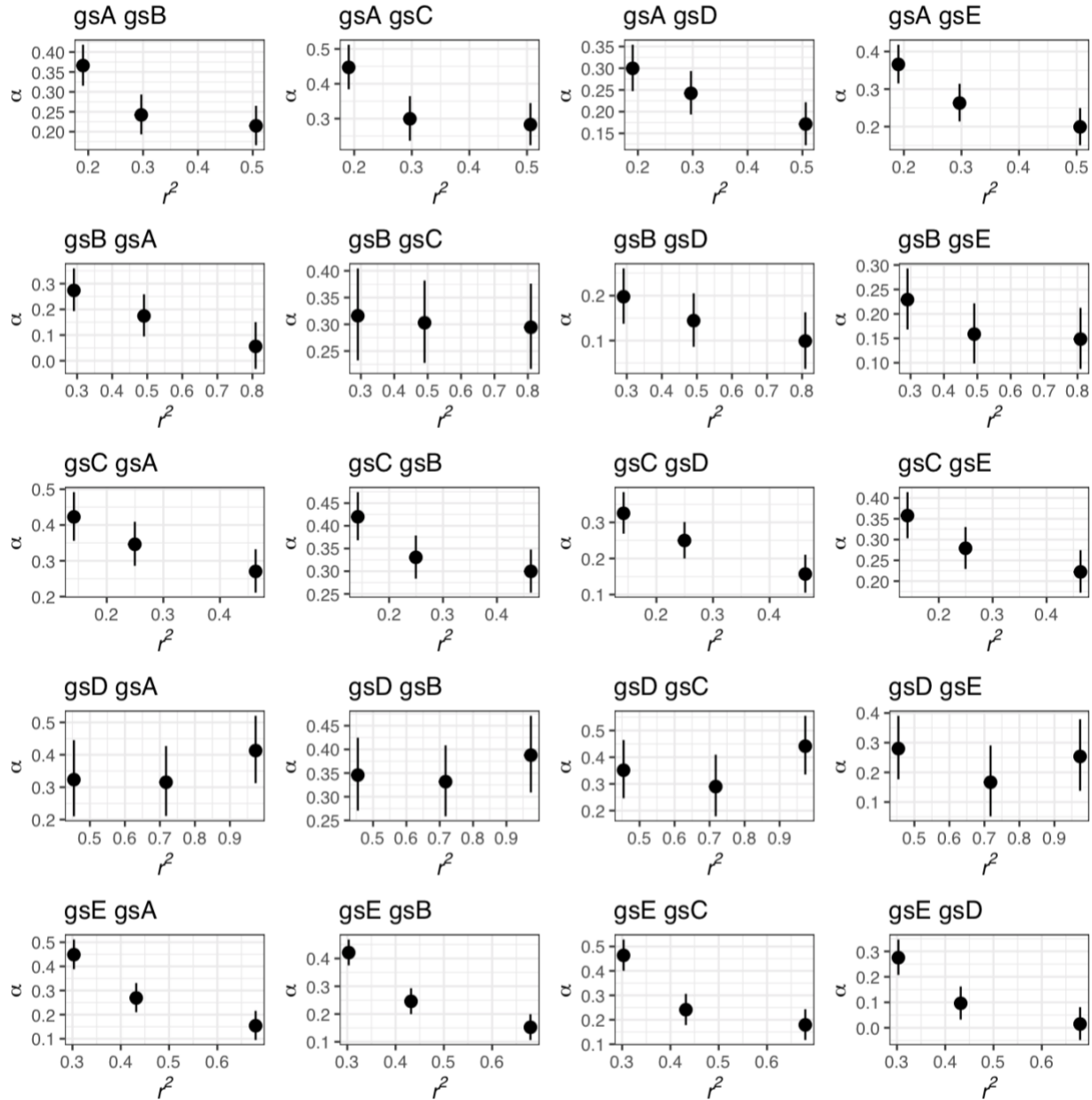
506 correspond to the median and mean r^2 , respectively. **(d)** Site frequency spectrum counts of

507 synonymous and non-synonymous sites by minor allele count based on all core genes (3086

508 genes). **(e)** The ratio of non-synonymous to synonymous polymorphisms by minor allele count.

509

510



511

512

513 **Figure 2. The proportion of adaptive evolution (α) by classes of recombination.** For each
514 pairwise estimates of α the polymorphism data from one species (left in title) is compared against
515 the divergence counts of an outgroup (right in title), and vice-versa. Results are divided into
516 classes of recombination based on r^2 (a measure that is inversely proportional to the level of
517 recombination). The α estimates and their associated confidence intervals were obtained using
518 the best fitting DFE model (GammaZero).

519

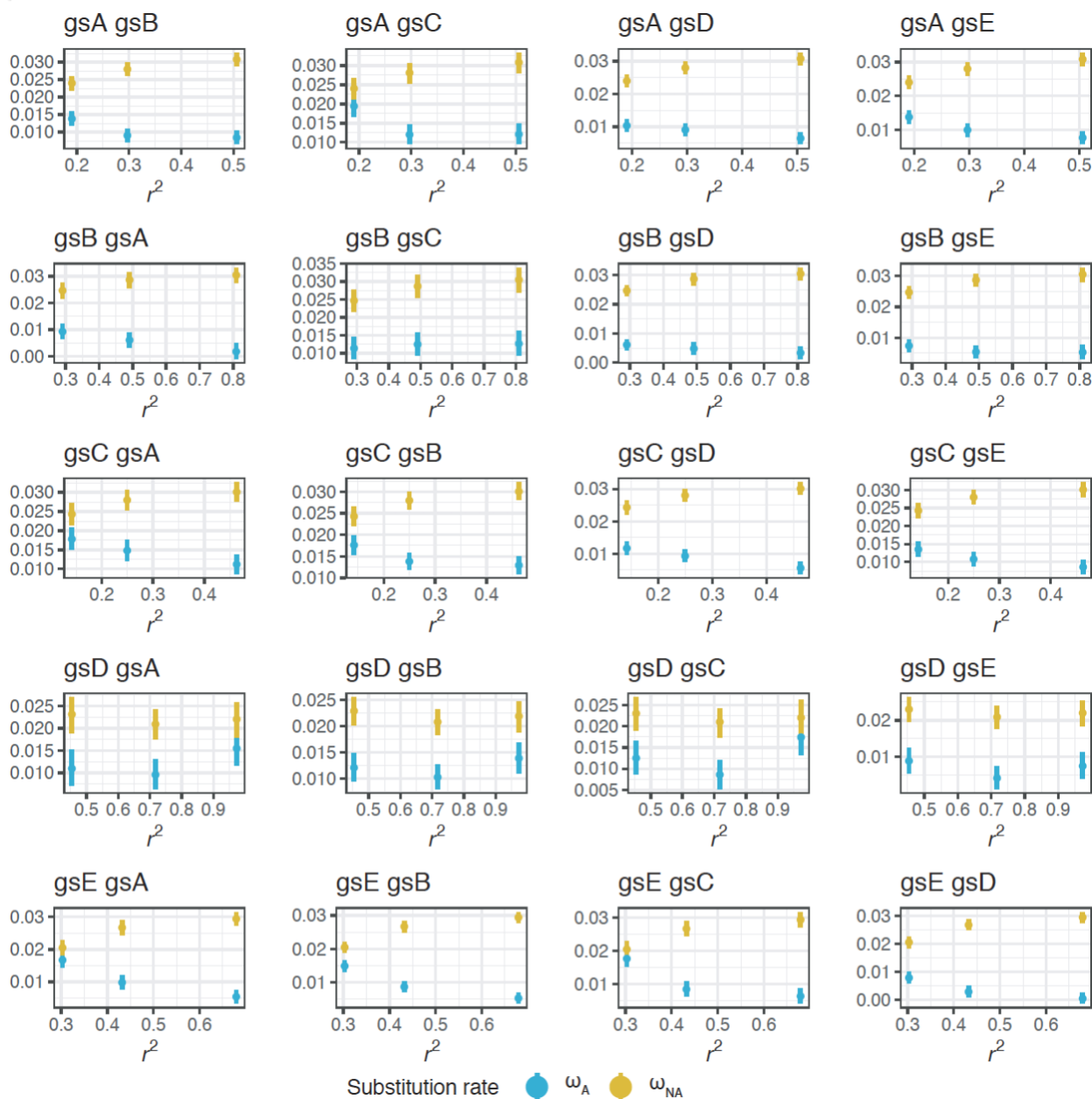
520

521

522

523

524



525

526

527 **Figure 3. The rates of adaptive (ω_A) and non-adaptive (ω_{NA}) evolution by classes of**
 528 **recombination.** For each pairwise estimates of ω_A (in blue) and ω_{NA} (in yellow) the polymorphism
 529 data from one species (left in title) is compared against the divergence counts of an outgroup
 530 (right in title), and vice-versa. Results are divided into classes of recombination based on r^2 (a
 531 measure that is inversely proportional to the level of recombination). An opposite effect of
 532 recombination on ω_A and ω_{NA} is observed in most pairwise comparisons. The estimates (ω_A ,
 533 ω_{NA}) and their associated confidence intervals were obtained using the best fitting DFE model
 534 (GammaZero).

535

536 **Table 1. The proportion of adaptive evolution (α) across pairs of species.** The α estimates
 537 were computed based on the best fitting DFE model (GammaZero) (**Table S4**). For each pairwise
 538 estimate of α ($\alpha_{species1\ species2}$), the polymorphism data from a focal species (rows) is compared
 539 against the divergence counts of an outgroup (columns), and vice-versa ($\alpha_{species2\ species1}$).
 540 Confidence intervals are displayed in brackets (grey) and numbers in parentheses represent the
 541 α ranking (in decreasing order) by outgroup (by column).
 542

Polymorphism (focal)	Divergence (outgroup)				
	gsA	gsB	gsC	gsD	gsE
gsA	-	0.28 [0.26-0.31] (2)	0.35 [0.33-0.39] (1)	0.25 [0.23-0.28] (1)	0.29 [0.27-0.32] (2)
gsB	0.18 [0.16-0.21] (4)	-	0.26 [0.24-0.29] (3)	0.15 [0.13-0.17] (3)	0.17 [0.15-0.19] (3)
gsC	0.36 [0.33-0.39] (1)	0.36 [0.33-0.38] (1)	-	0.25 [0.23-0.28] (1)	0.30 [0.27-0.32] (1)
gsD	0.25 [0.22-0.28] (3)	0.25 [0.23-0.27] (4)	0.25 [0.22-0.28] (4)	-	0.12 [0.10-0.15] (4)
gsE	0.27 [0.24-0.30] (2)	0.25 [0.23-0.27] (4)	0.27 [0.24-0.30] (2)	0.10 [0.07-0.13] (4)	-

543
 544