

Unpublished manuscript. Not for distribution.

Long-range Hill-Robertson effect in adapting populations with recombination and standing variation

Igor M. Rouzine^{1*}

¹Sechenov Institute of Evolutionary Physiology and Biochemistry

Russian Academy of Sciences, Saint-Petersburg 194223

*Correspondence should be addressed to igor.rouzine@iephb.ru

Abstract

In sexual populations, closely-situated genes have linked evolutionary fates, while genes spaced far in genome are commonly thought to evolve independently due to recombination. In the case where evolution depends essentially on supply of new mutations, this assumption has been confirmed by mathematical modeling. Here I examine it in the case of pre-existing genetic variation, where mutation is not important. A haploid population with N genomes, L loci, a fixed selection coefficient, and a small initial frequency of beneficial alleles f_0 is simulated by a Monte-Carlo algorithm. The results demonstrate the existence of extremely strong linkage effects, including clonal interference and genetic background effects, that depend neither on the distance between loci nor on the average number of recombination crossovers. When the number of loci, L , is larger than $4\log^2(Nf_0)$, beneficial alleles become extinct at most loci. The substitution rate varies broadly between loci, with the fastest rate exceeding the one-locus model prediction. All observables and the transition to the independent-locus limit are controlled by single composite parameter $\log^2(Nf_0)/L$. The potential link between these findings and the emergence of new Variants of Concern of SARS CoV-2 is discussed.

Introduction

A typical species is heterozygous at millions of genomic sites, loci. The average difference between an individual's genome and the consensus genome is estimated at 20 million base pairs, or 0.6% of the total of 3.2 billion base pairs (1). The invention of the new

Unpublished manuscript. Not for distribution.

34 methods of full-genome DNA sequencing caused the emergence of the field of genomics
35 and proteomics dedicated to the quantitative aspects of genetic diversity and gene
36 expression at a large number of loci (2-7). To describe and visualize the genetic
37 complexity, various computational methods have been developed including
38 phylogenetics, the principle-components analysis, the cluster analysis. Among them,
39 mathematical modeling of evolution stands out as a tool of a high predictive power.
40 Modeling allows to connect, in the most direct and reproducible fashion, the assumptions
41 about the dominant factors of evolution to the predictions for the observable parameters
42 of genetic diversity and evolutionary dynamics.

43 The assumptions and simplifications of models vary broadly depending on the
44 systems studied and the questions asked. Two distinct groups of models and methods
45 have been applied to animal populations and microbial populations. The classical one-
46 locus and two-locus models that neglect interaction with the other loci in genome (8-10)
47 dominate the way in which many evolutionary biologists think about the evolution of
48 higher organisms. In contrast, monocellular eukaryotes, viruses, and bacteria that are
49 characterized by an extremely high genetic diversity and ultrarapid evolution, are often
50 described by asexual or partly sexual population models that include explicitly large
51 numbers of interacting loci. Analysis of the evolutionary dynamics of multi-locus models
52 is more complex than one-locus and two-locus models and relies either on Monte-Carlo
53 simulation (11-17) or the advanced mathematical methods of statistical physics (11, 18-
54 39). The heavy mathematical artillery is required, because the evolution of many
55 different loci is inter-dependent (40). There are two kinds of interference effects. One
56 kind, not considered in this article, is epistasis arising from biological interaction of
57 different loci, including protein-protein interactions or interactions gene regulation
58 network (29, 31-33, 41-52). The second type, which is the focus of the present article, is
59 the effects originating from the common ancestry of different loci, including Hill-
60 Robertson effect, clonal interference, background selection and hitchhiking (8, 25, 53-
61 55). Linkage effects also slow down adaptation (11, 21, 23, 26), increase accumulation of
62 deleterious alleles (11, 21), and change the statistical shape of genealogical tree (24, 27,
63 56).

64 In sexually reproducing organisms and organisms with frequent recombination
65 such as some viruses, linkage effects are partly compensated by recombination between
66 parental genomes. A fundamental fact of genetics discovered by Morgan is that frequent

Unpublished manuscript. Not for distribution.

67 recombination destroys allelic associations, so that alleles at far-spaced loci segregate
68 independently. Conventional wisdom tells us that all the other linkage effects between
69 far-situated loci must vanish as well. A model of long-term sexual evolution limited by
70 rare mutation seemed to confirm this expectation (57). Assuming that genome consists
71 from independently-evolving blocks and applying the phylogenetic theory of asexual
72 evolution to each block, the authors constructed a scaling argument expressing the length
73 of each block, the lead of the traveling wave, and the average coalescent time in terms of
74 the average adaptation rate. The analytic predictions have been confirmed numerically
75 for two particular models of
76 population in the presence of
77 natural selection and
78 mutation.

79 In the present work, I
80 investigate linkage effects in
81 a different biological
82 scenario, when natural
83 selection and recombination
84 act on pre-existing beneficial
85 alleles, and new mutations
86 can be neglected. This model
87 is appropriate in the case
88 when selection pressure
89 changes its sign at a large
90 number of loci. For example,
91 a population migrates to a
92 new environment, or a virus
93 is subjected to the immune

94 response or a replication inhibitor treatment. In this case, weakly deleterious alleles pre-
95 existing in the mutation-selection balance can become beneficial.

96

97 Results

98 **Model.** Consider a sexually reproducing population comprised of N individual genomes
99 (or $N/2$ diploid genomes without allelic dominance), where each genome has L loci. In

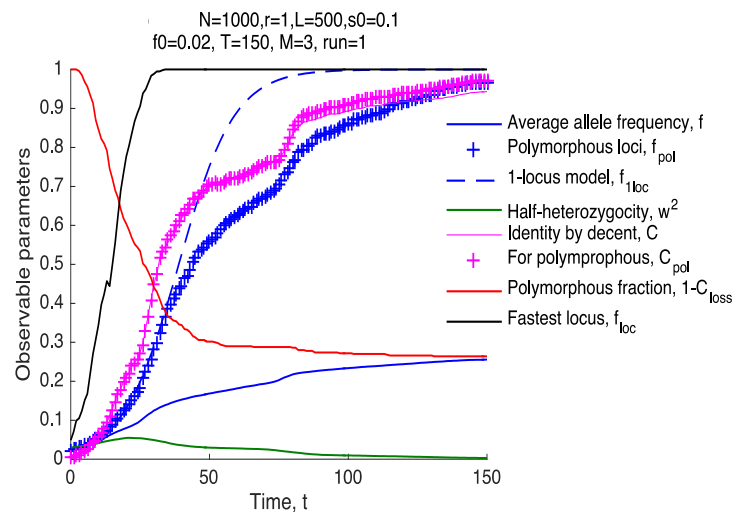


Fig. 1. Dynamics of observables in the model with standing variation and the absence of mutation.

Beneficial alleles become extinct at most loci. X-axis: Time in generations, t . Y-axis: Observable parameters calculated during simulation. The average frequency of beneficial alleles per locus per individual, f , the same value averaged over polymorphous loci only, f_{pol} , the prediction for f of the deterministic one-locus model, f_{1loc} , half-heterozygosity $w^2 = \langle f(1-f) \rangle$, the fraction of homologous pairs of loci with a common initial ancestor, C , the same value for polymorphous loci, C_{pol} , the fraction of polymorphous loci, $1-C_{loss}$, and the largest of allelic frequencies among loci, $\max(f_{loc})$. Parameter values are shown on the top. Parameters are defined in *Methods* and values are shown.

Unpublished manuscript. Not for distribution.

100 the beginning, each locus is assumed to have a fraction f_0 of beneficial alleles, with fitness
 101 benefit s . The value of f_0 is assumed to be in interval $\frac{1}{Ns} \ll f_0 \ll 1$. Next, I assume that a
 102 genome undergoes an average number M of random crossovers with another, randomly

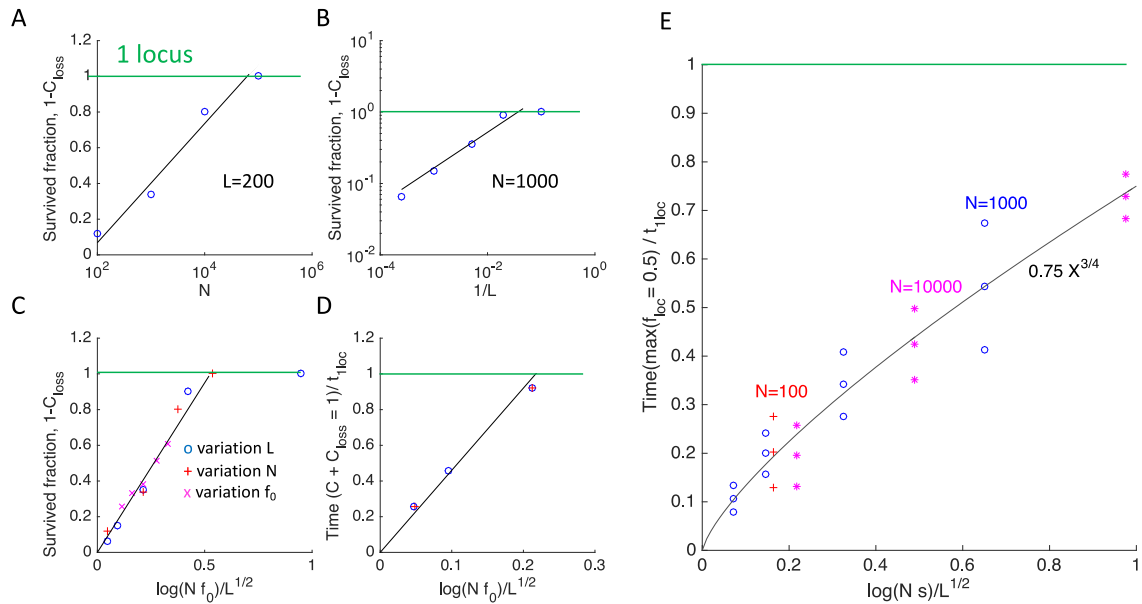


Fig. 2. The observables depend mostly on a single composite parameter.

A-C. The locus fraction where beneficial alleles have survived and completed adaptation, $1 - C_{loss}(\infty)$, is linearly proportional to the natural logarithm of the population size, $\log N$, the inverse square root of the locus number, $1/\sqrt{L}$, and a composite parameter, $\log(Nf_0)/\sqrt{L}$. Colored symbols \circ , $+$, and \times correspond to the variation of model parameters L , N , and f_0 , respectively, where $f_0 > 1/Ns$. The green horizontal line shows the prediction of the one-locus model, $C_{loss} \approx 0$. **D.** The time, t , when the survived-loci fraction, $1 - C_{loss}(t)$, equals the average identity by descent, $C(t)$, [intersection of red and pink curves in Fig. 2] scales linearly with $\log(Nf_0)/\sqrt{L}$ as well. **E.** The time when the allelic frequency at the fastest locus reaches 50%, scales as a power $3/4$ of a similar parameter, $\log(Ns)/\sqrt{L}$. The symbol triplets show the mean and the 95% confidence interval. Colored symbols \circ , $+$, and $*$ show different values of N . The sensitivity to the variation of selection coefficient s , crossover number M , and initial allele frequency f_0 is shown in S1 Fig and S2 Fig. The default parameter values are $N = 1000$, $L = 200$, $f_0 = 0.02$ unless shown otherwise. The other parameters are as in Fig. 2.

103 chosen genome, and one of the two parents is replaced with the recombinant. The
 104 evolution is simulated using a Wright-Fisher process, in which the progeny genomes
 105 replace the parental genome, and the average progeny number is proportional to the
 106 genome fitness. The evolutionary factors included in the model are directional natural
 107 selection, random genetic drift, linkage, and recombination. New mutation and epistasis
 108 are absent. The details of simulation are described in the *Methods* section.

109 **Extinction of beneficial alleles depends on a single composite parameter.** If the
 110 number of loci L is sufficiently large, beneficial alleles at most loci become extinct. The
 111 fraction of remaining polymorphous loci, denoted $1 - C_{loss}(t)$, decreases in time from 1

Unpublished manuscript. Not for distribution.

112 to at a low plateau (Fig. 1A, red line). This result differs from the prediction of the single-
 113 locus model, in which multiple lineages per site are expected to reach fixation at $Nf_0s \gg$
 114 1. In that case, the fixation probability of an allele is s , and the extinction probability is
 115 $1 - s$ (40). The probability of the extinction of all Nf_0 beneficial lineages is given by

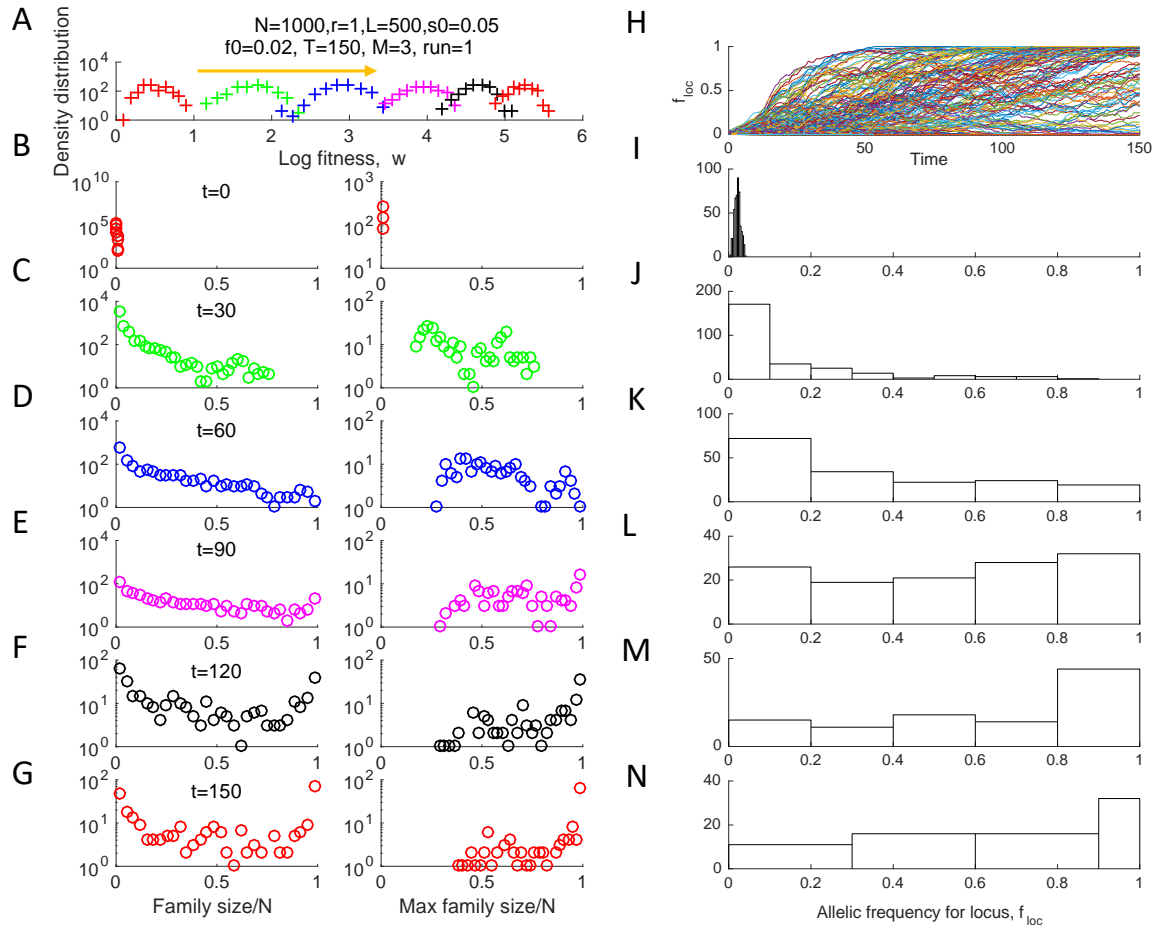


Fig. 3. Traveling fitness wave and nonuniform dynamics of separate loci. A. Distribution density of genomes in fitness at different time points shown in (B-G). B-G. First column: Histograms of the family size defined as the number of sequences with the same initial ancestor at a locus. Second column: Only the largest family per locus is taken into account. H. The average allelic frequency for each separate locus, f_{loc} , as a function of time. I-N. Histograms of f_{loc} across loci at different time points (shown). Parameters are as in Fig. 2.

116 $C_{loss}(\infty) = (1 - s)^{Nf_0} \approx e^{-Nf_0s}$, which is exponentially small.

117 Varying model parameters in simulation, we found out empirically out that the
 118 fraction of loci with non-extinct alleles, $1 - C_{loss}(\infty)$, depends mostly on a single
 119 composite parameter (Fig. 2A-C)

$$1 - C_{loss} = \begin{cases} 2.0 \frac{\log(Nf_0)}{\sqrt{L}} & 1 \ll \log(Nf_0) < 0.5\sqrt{L} \\ 1 & \log(Nf_0) > 0.5\sqrt{L} \end{cases} \quad (1)$$

Unpublished manuscript. Not for distribution.

121 Note the critical point, $\log(Nf_0) = 0.5\sqrt{L}$. If the population size is too large or the
122 number of loci is too small, no significant loss of polymorphism is predicted.

123 **The fastest adaptation rate among loci is much faster than in a single-locus**
124 **model.** Because most loci fail to complete adaptation, the average frequency of beneficial
125 alleles per locus, $f_{av}(t)$, saturates far below 1 (Fig. 1A, blue line). The dependence of
126 average heterozygosity on time, $2w^2(t)$, is decreased accordingly (Fig. 1A, green). The
127 allele frequency averaged over remaining polymorphic sites, $f_{pol}(t)$, increases in the
128 same general time range as the one-locus prediction. The time of half-fixation of
129 polymorphous sites, t_{50} , is very close to the deterministic one-locus prediction, $t_{50} \approx t_{1loc}$
130 (Fig. 1B)

$$131 \quad t_{1loc} = \frac{1}{s} \log \frac{1}{f_0} \quad (2)$$

132 In the range of parameters $s = 0.025 - 0.2, L = 200 - 2000, N = 1000 - 10,000$, the
133 relative difference between t_{50} and t_{1loc} is between -0.11 and 0.14. Compared to the 1-
134 locus model prediction (blue dashed line in Fig. 1), the dependence $f(t)$ experiences a
135 delay in the late phases of adaptation and has a noticeable random oscillation component
136 (Fig. 1, blue +).

137 The speed of adaptation is extremely broadly distributed among loci with non-extinct
138 alleles (Fig. 3H). At some loci, alleles accumulate much faster than predicted by the one-
139 locus model (Fig 1, black line). The half-time of adaptation of the fastest locus, $\max(t_{loc})$,
140 is much shorter than t_{1loc} and increases as power $3/4$ of composite parameter $\frac{\log(Ns)}{\sqrt{L}}$ (Fig.
141 2E) (compare with Eq. 1). The broad variation between loci is created by random
142 recombination events, which bring together different numbers of favorable alleles, and
143 natural selection, which favors the best. As a result, the distribution of genomes in fitness
144 forms a traveling wave well-known for both asexual and sexual populations (40) (Fig.
145 3A).

146 The fitness classes of the traveling wave have a complex lineage structure that
147 varies between loci. For a given locus, a lineage is determined as the set of individuals
148 that have the same initial ancestor. The lineages all initially consists from a single
149 individual (Fig. 3B), but their sizes grow in time, at different rates for different loci, and
150 become distributed in a very broad range (Fig. 3C-G). The size distribution shifts in time
151 towards larger lineages eventually occupying almost the entire population. If we take into
152 account only the largest lineage for each locus, their size distribution looks similar but

Unpublished manuscript. Not for distribution.

153 has a low cutoff increasing in time (Fig. 3B-G, column 2). The largest lineages grow to a
 154 half of the population at a much earlier time than t_{1loc} in Eq. 2.

155 **Phylogenetic time scale depends only on the same composite parameter.**

156 Another quantity affected by linkage effects is the identity by descent, C , defined as the
 157 probability of a homologous locus pair to have the same initial ancestor. The average
 158 identity by descent averaged over all loci and over only polymorphous loci is almost the

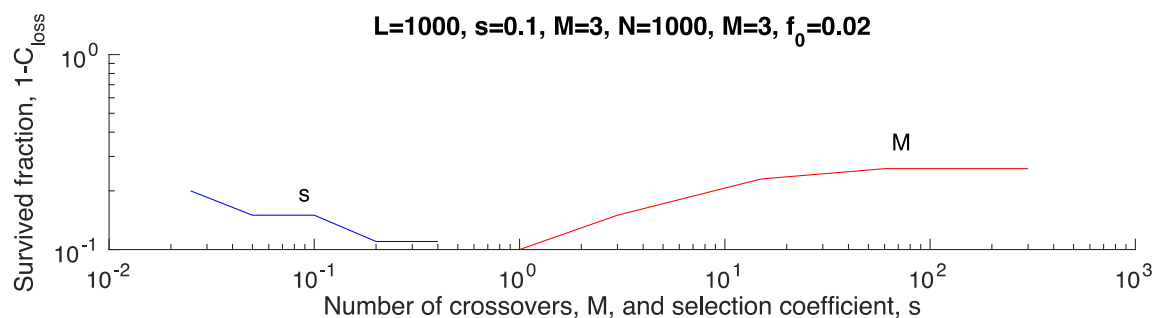


Fig. 4. Weak sensitivity of the fraction of loci that complete adaptation to the selection coefficient and the average crossover number. The default parameter values are shown on the top.

159 same (magenta line and magenta +, Fig. 1A). This result differs from the single-locus
 160 model, where common ancestry is rare, $C(t) < f_{pol}^2(t)$, because each of the pair of loci
 161 must fall into the same growing lineage to have the same ancestor, and the size of each
 162 lineage relative to the population size is smaller than $f_{pol}(t)$. In contrast, in our case, C is
 163 larger than $f_{pol}(t)$, which is larger than $f_{pol}^2(t)$. At the time point T_2 where $C = 1 - C_{loss}$,
 164 both values are both close to a half in a broad parameter range, $C(T_2) \approx C_{loss}(T_2) \approx 0.5$.
 165 The dependence of T_2 on model parameters can be fit by the formula

$$166 \quad T_2 \approx t_{1loc} \frac{5.0 \log(Nf_0)}{\sqrt{L}} \quad (3)$$

167 In other words, T_2 is proportional to the same composite parameter that controls the
 168 fraction of fixed loci, $1 - C_{loss}$, Eq. 1 (Fig. 3D). Time T_2 defined by Eq. 3 represents a
 169 proxy time scale of the phylogenetic tree. Although, at this time point, a population does
 170 not have a single ancestor for an average locus as yet, T_2 approximates the time to the
 171 most recent common ancestor by an order of magnitude.

172 **Weak dependence of all observables on the average number of recombination**
 173 **crossovers.** The above results in Figs 1 to 3 are weakly sensitive to the average crossover
 174 number, M . In its entire range of between 1 and L , the fraction of loci that do not lose
 175 alleles, $1 - C_{loss}(\infty)$, varies only by the factor of ~ 2 (Fig. 4).

Unpublished manuscript. Not for distribution.

176 **The absence of long-range linkage disequilibrium.** No linkage disequilibrium is
177 predicted in the long range. Pearson's correlator between allelic frequencies at two loci
178 defined as

$$179 \quad r^2(l_{12}) = \frac{\langle (f_1 - \langle f \rangle)(f_2 - \langle f \rangle) \rangle}{\langle (f_1 - \langle f \rangle)^2 \rangle}$$

180 decreases rapidly with the distance between loci, l_{12} , and the characteristic distance of
181 the decrease shrinks with time (Fig S1). In other words, alleles at far loci segregate
182 independently, as they should in the presence of recombination.

183 **Far blocks of genome do not evolve independently.** The above results for the
184 phylogeny time scale differ from that of scaling theory (57). In my notation, their general
185 result for the average time to the most recent common ancestor has the form [(57), Eq.
186 5]

$$187 \quad T_{MRCA} \approx const \frac{M}{v} \log \left(\frac{Nv}{M} \right) \quad (4)$$

188 where v is the average rate of long-term adaptation, defined as the fitness gain per unit
189 time, $const$ is a number on the order of 1, and the logarithm is supposed to be much larger
190 than 1. In my case, the proxy of T_{MRCA} by the order of magnitude is T_2 in Eq. 3, and the
191 adaptation rate is (see Fig. 1A)

$$192 \quad v \approx const \frac{sL(1 - C_{loss})}{t_{50}} \quad (5)$$

193 As already mentioned, the average time to a half-fixation for the loci that do not lose
194 alleles, t_{50} , is always close to one-locus limit t_{1loc} . Substituting Eq. 3 and Eq. 5 into Eq. 4,
195 we get

$$196 \quad \frac{M \log(Nv/M)}{s \log^2(Nf_0)} = const$$

197 which is clearly false, because M , N , and s are independent parameters. Hence, Eq. 4 does
198 not work in the case with pre-existing variation.

Unpublished manuscript. Not for distribution.

199 Note that the analytic argument in (57) was developed and tested for a different
 200 scenario, when the sexual evolution is limited by new mutation events. It was based on
 201 two statements: the assumption that a genome evolves as quasi-independent asexual
 202 blocks, and an expression for the time to the most recent common ancestor in terms of
 203 the average adaptation rate. The expression was based on the basic concept that the time
 204 to most recent common ancestor is the lead of the wave divided by the adaptation rate
 205 and was confirmed for various multi-locus models, both sexual and asexual. Therefore,
 206 it is likely that the quasi-
 207 independence assumption is
 208 the cause of the discrepancy. In
 209 other words, in the case of pre-
 210 existing variation, the genome
 211 does not evolve as a set of
 212 quasi-independent segments.
 213 That conclusion is indirectly
 214 confirmed by the results in Fig.
 215 3 showing that beneficial
 216 alleles can form highly-fit
 217 genomes whose rapid growth
 218 outruns mixing of genomes due
 219 to recombination (Fig. 3). A
 220 recombinant that decreases
 221 fitness is not relevant for
 222 future generations.

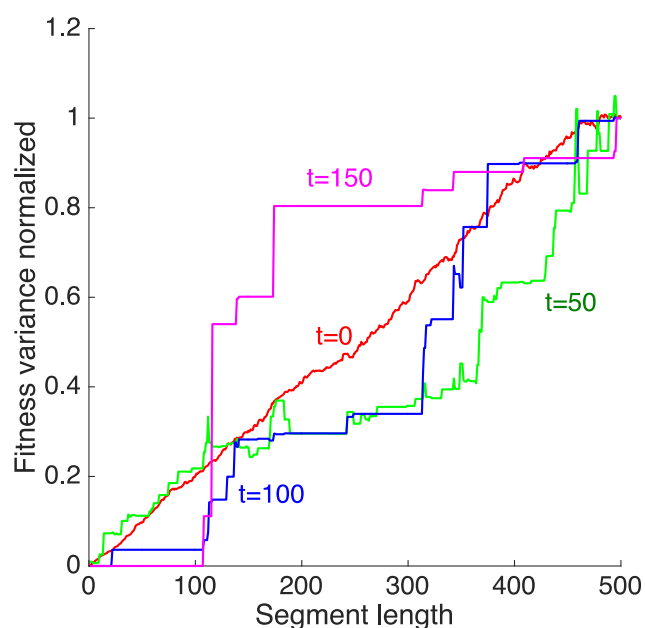


Fig. 5. Non-linear dependence of genome segment variance on segment length. X-axis: the length of a genome segment starting from locus 1. Y-axis: Fitness variation between homologous genomic segments divided by the genome fitness variation at the same moment of time. A single run is shown. Parameters are, as in Fig. 1.

223 Furthermore, within one realization (Monte Carlo run), the fitness variance of a genomic
 224 segment normalized to the genome fitness is not linearly proportional to its length, but
 225 shows a complex step-like dependence (Fig. 5).

226 **Alleles are fixed inter-dependently.** The fixation probability of an allele can be
 227 calculated as

$$228 \quad P_{fix} = \frac{1 - C_{loss}(\infty)}{Nf_0} \quad (6)$$

229 In the parameter interval of interest, this value falls far below the 1-locus prediction,
 230 $P_{fix}^{1loc} = s$ (Fig S2). Probability P_{fix} plateaus on the value of s in the dilute limit of

Unpublished manuscript. Not for distribution.

231 sufficiently small f_0 , which agrees with a previous finding in the case of rare mutations,
232 see the limit $r \gg s$ in (30). Based on simulation, the transition point to the dilute limit
233 f_0^{dilute} decreases with N and L . One can determine the transition point from condition
234 $P_{fix}(f_0^{dilute}) = s$ and Eqs. 1 and 6. Replacing $\log(Nf_0)$ with 1 if it smaller than 1, we
235 obtain

$$237 \quad f_0^{dilute} \approx \frac{2}{Ns\sqrt{L}}, \quad s\sqrt{L} \geq 1 \quad (3)$$

236 This estimate agrees with the simulation results in Fig S2.

238 **Phylogenetic tree and allele surfing.** In addition to calculating the phylogeny time
239 scale (Fig. 2D), we constructed the ancestral trajectory of a locus between individuals in
240 real-time by memorizing the parentage of each individual locus and then tracing its
241 ancestry back in time. Lineage of each locus jumps among individuals randomly due to
242 recombination (Fig. 6A). If we straighten these trajectories and keep only the topology of
243 coalescence and the coalescent times, we arrive at phylogenetic trees for different loci
244 (Fig. 6B-D). As expected, the tree varies strongly across loci, and the early branches are
245 relatively shorter than in the neutral Kingman's coalescent. The average density of
246 coalescent events averaged over 10 runs and normalized to the prediction of the
247 selectively-neutral model (*Methods*) decreases exponentially with time (Fig. 6E, F), as it
248 would also in the one-locus limit. This is because coalescent density is proportional to the
249 inverse effective population size (58), which is the size of the growing variant
250 subpopulation. However, the coalescent density is also much larger than in the one-locus
251 limit and increases with number of loci L . Thus, in agreement with the previous studies
252 for various models, uncompensated linkage in the presence of selection makes
253 phylogenetic trees denser and changes their shape by making early branches shorter (24,
254 27, 56, 59) (Fig. 6E, F).

255 In addition to the trajectory of a locus over specific ancestors (Fig 6A), we can also
256 construct its fitness trajectory, by memorizing the fitness values of its ancestors (Fig 6G).
257 The fitness trajectory comprises alternating straight horizontal segments due to the
258 clonal expansion connected to jumps caused by recombination. The jumps occur in both
259 directions, but more often towards a genetic background with a higher fitness (Fig. 6G).

Unpublished manuscript. Not for distribution.

260 This “allelic surfing” behavior with vertical and horizontal segments was predicted
261 analytically for sexual populations with a small outcrossing rate (30, 60).

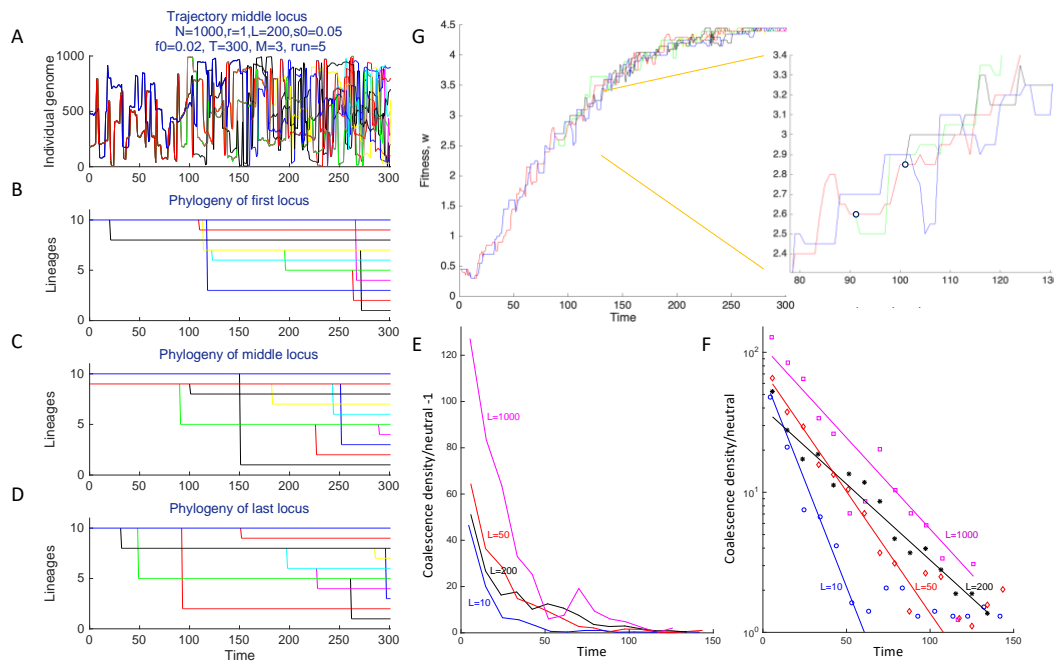


Fig. 6. Phylogenetic tree and ancestral history of separate loci.

A. A reverse-time trajectory of the middle locus ($i = L/2$) in 10 individuals numbered 1, 101, 201, ... 901 at time $t = 300$ obtained by tracing ancestral history. **B-D.** Phylogenetic trees for three loci (first, middle, and last). **E-F.** The time density of coalescent events averaged over 10 simulation runs and normalized to their values predicted by the selectively neutral model (*Methods*). Linear (E) and logarithmic (F) scales are used for Y-axis. **G.** Fitness trajectories for the middle locus in (A, C). Insert: a small segment is shown by the orange square. Parameters are on (A) unless shown otherwise.

262

263 Discussion

264

265 I modeled numerically stochastic sexual evolution of a multi-locus system due to
266 natural selection and pre-existing variation in the form of small numbers of beneficial
267 alleles. Despite of the lack of observable LD for far-situated loci, simulation predicts the
268 existence of string long-range linkage effects encompassing the entire genome. The
269 effects include the extinction of beneficial alleles at most loci due to clonal interference,
270 weak sensitivity of most observables to the average number of crossovers, and a very fast
271 evolution at a fraction of loci. These results are in striking contrast to the previous
272 findings for the long-term evolution driven by mutation, selection, and recombination,
273 where genome was demonstrated to consist from quasi-independent blocks (57). The
274 linkage effects are predicted only for sufficiently long genomes.

Unpublished manuscript. Not for distribution.

275 If the locus number is decreased, or if the population size is increased, a transition to
276 the independent-locus limit is predicted. The predicted dependence of all linkage effects
277 on the population size is logarithmic (Fig. 2). For a genome of 200 loci and $f_0 = 0.02$, $s =$
278 0.1, the transition to the independent-locus regime can be observed already for 100,000
279 individuals. For a longer genome of 1000 loci, however, loci evolve independently only
280 for populations of 10^{12} individuals or larger, which is unrealistic for most species. A
281 human or an animal population has millions of variable loci, of which a sizeable portion
282 is under selection, so that independent-locus models, probably, never work in most
283 animals, except for rare mutations that are under very strong selection pressure.

284 We have investigated the case of a constant selection coefficient, but the results are
285 expected to apply also for a sufficiently fast decaying distribution of selection coefficients,
286 such as a Gaussian distribution. Distributions with long tails may have different
287 properties, where the traveling wave is replaced by pairwise clonal interference (26). The
288 case of an exponential distribution can have a mixed behavior, depending on parameter
289 values (26). The exponential distribution is often observed in experiments on pathogens
290 which fact has been explained in a recent work (34).

291 The results obtained are directly relevant for the viruses that have frequent
292 recombination, such as HIV, polio, or SARS CoV-2. Similar to seasonal human
293 coronaviruses or influenza virus, SARS-CoV-2 is constantly acquiring new mutations in
294 its genome. Evolution is especially fast in receptor Spike protein (61-64). Two major
295 reasons account for the high speed of evolution, as follows. Firstly, Spike has receptor-
296 binding motives that affect transmission, and their evolution leads to the emergence of
297 VOCs with enhanced transmissibility. Secondly, Spike contains epitopes, regions that are
298 very important for the immune response because of their involvement in binding of
299 antibodies that can neutralize virus. Mutations in epitopes are a major factor that limits
300 the virus recognition by the immune system and, hence, the durability of protection (65,
301 66).

302 An important puzzle important for devising future vaccination strategies is the origin
303 of the VOCs produced by large groups of new mutations that emerge all together at once
304 (67, 68). Alternative theories of the emergence of VOCs (69) include reverse zoonosis, the
305 evolution within immunocompromised patients (70-72), and the evolution in population
306 pockets not covered by the genetic surveillance. Still another possibility is the fitness

Unpublished manuscript. Not for distribution.

307 valley effect, a cascade emergence of compensating mutations following a primary
308 mutation inferred for HIV and influenza (33, 73).

309 Based on the present study, we may add yet another possible explanation. While
310 in another respiratory virus, influenza, we observe only rare reassortment of its eight
311 chromosomes, SARS CoV-2, with its single-chromosome genome, has observable
312 crossover recombination (74-78). Hence, the large packages of mutations may emerge
313 due to the combined effects of recombination and natural selection and represent the
314 sequences formed by the fastest loci (Fig. 3). To understand the importance of
315 recombination for SARS CoV-2, we need to know the frequency of co-infected individuals
316 among all the infected, which determines outcrossing probability r , an important input
317 parameter entering the models of sexual populations (13, 19, 30, 57, 60, 79). For fully
318 sexual reproduction considered in the present work, by the definition, $r = 1$. The
319 outcrossing number for SARS-CoV-2 is presently unknown. It could be quite large due to
320 the possibility of a co-infection during superspreading events (80-83). Methods
321 developed previously to quantify recombination from RNA sequence data for HIV could
322 be re-applied to SARS-CoV-2 (13).

323 **Conclusion.** In sexual populations with pre-existing beneficial alleles, in an
324 exponentially broad range of population size, recombination cannot suppress long-range
325 linkage effects, such as the excessive loss of beneficial alleles at most loci, the lack of
326 dependence on the crossover number, and superfast evolution at some loci. These
327 findings may be relevant for interpreting the emergence of new strains of SARS CoV-2.

328

329 **Materials and methods**

330

331 Consider a fully sexual population with L loci comprised of N individual genomes. Each
332 locus has initially Nf_0 alleles, $1/Ns \ll f_0 \ll 1$, with fitness benefit $s \ll 1$. In each
333 generation step, each genome undergoes random crossovers with another, randomly
334 chosen genome, with average crossover number M producing a recombinant genome.
335 One of the two parents is replaced with the recombinant. Genome number j with
336 k_j favorable alleles is replaced with a random number of its copies distributed according
337 to the polynomial distribution implemented by “broken stick” method, as follows. N
338 random points are generated uniformly within the interval $[0, N]$ broken into N

Unpublished manuscript. Not for distribution.

339 segments. The length of segment j is proportional to the fitness of the corresponding
340 genome w_j

$$341 \quad w_j = \frac{\exp(-sk_j)}{\sum_{j=1}^N \exp(-sk_j)} \quad (3)$$

342 The number of random values that fall into segment j are taken to be the number of his
343 progeny in the next generation. Thus, the total number of genomes stays constant. New
344 mutations are neglected, which is shown to be correct in the short-term in the presence
345 of pre-existing genetic variation, both in simulation and experimentally (84). Epistasis is
346 absent. For the modeling studies of epistatic effect, the reader is referred to (31-33, 47-
347 51).

348 Input model parameters are the selection coefficient across loci, $s = s_0$, population
349 size N , outcrossing rate $r = 1$, number of loci L , initial beneficial allele frequency f_0 , total
350 simulation time t , average number of recombination crossovers M , and the seed number
351 of the generator of pseudorandom numbers.

352 Parameter ranges studied are $s = [0.025, 0.4]$, $L = [10, 4000]$, $N = [10^2, 10^5]$, $M =$
353 $[1, 300]$, $f_0 = [0.0001, 0.02]$. The main focus is on the interval of f_0 such that $\frac{1}{Ns} \ll f_0 \ll$
354 1. The transition to dilute limit $Nf_0s \ll 1$ when alleles are fixed independently is shown
355 in Fig. S2.

356

357

358 **Funding:** This research was partly funded by Agence Nationale de la Recherche, France,
359 grant number J16R389 to I.M.R.

360

361 **Acknowledgement:** The study was carried out within the framework of the state
362 assignment of the Federal Agency for Scientific Organizations (FASO Russia: topic no.
363 AAAA-A18-118012290142-9).

364

365 **Competing interests:** The funders had no role in the design of the study; in the collection,
366 analyses, or interpretation of data; in the writing of the manuscript, or in the decision to
367 publish the results.

368

Unpublished manuscript. Not for distribution.

369 **Data and materials availability:** The simulation code is available at
370 <https://github.com/irouzine/Strong-linkage-in-sex>.

371
372

373 **References**

374

- 375 1. G. P. Consortium *et al.*, A global reference for human genetic variation. *Nature* **526**,
376 68-74 (2015).
- 377 2. M. Plesser Duvdevani *et al.*, Whole-genome sequencing reveals complex chromosome
378 rearrangement disrupting NIPBL in infant with Cornelia de Lange syndrome. *Am J Med*
379 *Genet A* **182**, 1143-1151 (2020).
- 380 3. H. Stranneheim *et al.*, Integration of whole genome sequencing into a healthcare
381 setting: high diagnostic rates across multiple clinical entities in 3219 rare disease
382 patients. *Genome Med* **13**, 40 (2021).
- 383 4. V. Chat, R. Ferguson, L. Morales, T. Kirchhoff, Ultra Low-Coverage Whole-Genome
384 Sequencing as an Alternative to Genotyping Arrays in Genome-Wide Association
385 Studies. *Front Genet* **12**, 790445 (2021).
- 386 5. E. J. Muturi, C. Dunlap, D. P. Tchouassi, J. Swanson, Next generation sequencing
387 approach for simultaneous identification of mosquitoes and their blood-meal hosts. *J*
388 *Vector Ecol* **46**, 116-121 (2021).
- 389 6. Q. Wang *et al.*, Using Next-generation Sequencing to Identify Novel Exosomal miRNAs
390 as Biomarkers for Significant Hepatic Fibrosis. *Discov Med* **31**, 147-159 (2021).
- 391 7. N. P. Mthethwa, I. D. Amoah, P. Reddy, F. Bux, S. Kumari, A review on application of
392 next-generation sequencing methods for profiling of protozoan parasites in water:
393 Current methodologies, challenges, and perspectives. *J Microbiol Methods* **187**,
394 106269 (2021).
- 395 8. D. S. Fisher, *The genetical theory of natural selection* (Clarendon Press, Oxford, United
396 Kingdom, 1930).
- 397 9. M. Kimura, *Population genetics, molecular evolution, and the neutral theory. Selected*
398 *papers.* (The University of Chicago Press, Chicago 1994).
- 399 10. I. M. Rouzine, A. Rodrigo, J. M. Coffin, Transition between stochastic evolution and
400 deterministic evolution in the presence of selection: general theory and application to
401 virology. *Microbiol Mol Biol Rev* **65**, 151-185 (2001).
- 402 11. I. M. Rouzine, J. Wakeley, J. M. Coffin, The solitary wave of asexual evolution. *Proc*
403 *Natl Acad Sci U S A* **100**, 587-592 (2003).
- 404 12. S. Gheorghiu-Svirschevski, I. M. Rouzine, J. M. Coffin, Increasing sequence correlation
405 limits the efficiency of recombination in a multisite evolution model. *Mol Biol Evol* **24**,
406 574-586 (2007).
- 407 13. R. Batorsky *et al.*, Estimate of effective recombination rate and average selection
408 coefficient for HIV in chronic infection. *Proc Natl Acad Sci U S A* **108**, 5661-5666 (2011).
- 409 14. T. Bedford, A. Rambaut, M. Pascual, Canalization of the evolutionary trajectory of the
410 human influenza virus. *BMC Biol* **10**, 38 (2012).
- 411 15. T. Bedford *et al.*, Global circulation patterns of seasonal influenza viruses vary with
412 antigenic drift. *Nature* **523**, 217-220 (2015).
- 413 16. T. Bedford *et al.*, Integrating influenza antigenic dynamics with molecular evolution.
414 *Elife* **3**, e01914 (2014).

Unpublished manuscript. Not for distribution.

- 415 17. J. Neidhart, I. G. Szendro, J. Krug, Adaptation in tunably rugged fitness landscapes: the
416 rough Mount Fuji model. *Genetics* **198**, 699-721 (2014).
- 417 18. E. Brunet, B. Derrida, Exactly soluble noisy traveling-wave equation appearing in the
418 problem of directed polymers in a random medium. *Phys Rev E Stat Nonlin Soft Matter*
419 *Phys* **70**, 016106 (2004).
- 420 19. I. M. Rouzine, J. M. Coffin, Evolution of human immunodeficiency virus under selection
421 and weak recombination. *Genetics* **170**, 7-18 (2005).
- 422 20. E. Brunet, B. Derrida, A. H. Mueller, S. Munier, Phenomenological theory giving the full
423 statistics of the position of fluctuating pulled fronts. *Phys Rev E Stat Nonlin Soft Matter*
424 *Phys* **73**, 056126 (2006).
- 425 21. I. M. Rouzine, E. Brunet, C. O. Wilke, The traveling-wave approach to asexual
426 evolution: Muller's ratchet and speed of adaptation. *Theor Popul Biol* **73**, 24-46 (2008).
- 427 22. E. Brunet, I. M. Rouzine, C. O. Wilke, The stochastic edge in adaptive evolution.
428 *Genetics* **179**, 603-620 (2008).
- 429 23. M. M. Desai, D. S. Fisher, Beneficial mutation selection balance and the effect of
430 linkage on positive selection. *Genetics* **176**, 1759-1798 (2007).
- 431 24. M. M. Desai, A. M. Walczak, D. S. Fisher, Genetic diversity and the structure of
432 genealogies in rapidly adapting populations. *Genetics* **193**, 565-585 (2013).
- 433 25. P. J. Gerrish, R. E. Lenski, The fate of competing beneficial mutations in an asexual
434 population. *Genetica* **102-103**, 127-144 (1998).
- 435 26. B. H. Good, I. M. Rouzine, D. J. Balick, O. Hallatschek, M. M. Desai, Distribution of fixed
436 beneficial mutations and the rate of adaptation in asexual populations. *Proc Natl Acad*
437 *Sci U S A* **109**, 4950-4955 (2012).
- 438 27. R. A. Neher, O. Hallatschek, Genealogies of rapidly adapting populations. *Proc Natl*
439 *Acad Sci U S A* **110**, 437-442 (2013).
- 440 28. R. A. Neher, C. A. Russell, B. I. Shraiman, Predicting evolution from the shape of
441 genealogical trees. *Elife* **3** (2014).
- 442 29. R. A. Neher, B. I. Shraiman, Statistical genetics and evolution of quantitative traits.
443 *Reviews of Modern Physics* **83**, 1283 (2011).
- 444 30. R. A. Neher, B. I. Shraiman, D. S. Fisher, Rate of adaptation in large sexual populations.
445 *Genetics* **184**, 467-481 (2010).
- 446 31. G. Pedruzzi, A. Barlukova, I. M. Rouzine, Evolutionary footprint of epistasis. *PLoS*
447 *Comput Biol* **14**, e1006426 (2018).
- 448 32. G. Pedruzzi, I. M. Rouzine, Epistasis detectably alters correlations between genomic
449 sites in a narrow parameter window. *PLoS One* **14**, e0214036 (2019).
- 450 33. G. Pedruzzi, I. M. Rouzine, An evolution-based high-fidelity method of epistasis
451 measurement: Theory and application to influenza. *PLoS Pathog* **17**, e1009669 (2021).
- 452 34. A. Barlukova, I. M. Rouzine, The evolutionary origin of the universal distribution of
453 mutation fitness effect. *PLoS Comput Biol* **17**, e1008822 (2021).
- 454 35. I. M. Rouzine, G. Rozhnova, Antigenic evolution of viruses in host populations. *PLoS*
455 *Pathog* **14**, e1007291 (2018).
- 456 36. S. Schiffels, G. J. Szollosi, V. Mustonen, M. Lassig, Emergent neutrality in adaptive
457 asexual evolution. *Genetics* **189**, 1361-1375 (2011).
- 458 37. J. Marchi, M. Lassig, A. M. Walczak, T. Mora, Antigenic waves of virus-immune
459 coevolution. *Proc Natl Acad Sci U S A* **118** (2021).
- 460 38. L. Yan, R. A. Neher, B. I. Shraiman, Phylodynamic theory of persistence, extinction and
461 speciation of rapidly adapting pathogens. *Elife* **8** (2019).

Unpublished manuscript. Not for distribution.

- 462 39. I. M. Rouzine, An Evolutionary Model of Progression to AIDS. *Microorganisms* **8** (2020).
463 40. I. Rouzine, *Mathematical Modeling of Evolution*
464 *Volume 1: One-Locus and Multi-Locus Theory and Recombination*, De Gruyter Series in
465 Mathematics and Life Sciences (De Gruyter, Berlin/Boston, 2021),
466 10.1515/9783110615456, pp. 185.
467 41. F. Morcos *et al.*, Direct-coupling analysis of residue coevolution captures native
468 contacts across many protein families. *Proc Natl Acad Sci U S A* **108**, E1293-1301
469 (2011).
470 42. A. Pick *et al.*, Structure-activity relationships of flavonoids as inhibitors of breast
471 cancer resistance protein (BCRP). *Bioorg Med Chem* **19**, 2090-2102 (2011).
472 43. A. Procaccini, B. Lunt, H. Szurmant, T. Hwa, M. Weigt, Dissecting the specificity of
473 protein-protein interaction in bacterial two-component signaling: orphans and
474 crosstalks. *PLoS One* **6**, e19729 (2011).
475 44. J. Rodriguez-Rivas, G. Croce, M. Muscat, M. Weigt, Epistatic models predict mutable
476 sites in SARS-CoV-2 proteins and epitopes. *Proc Natl Acad Sci U S A* **119** (2022).
477 45. A. Schug, M. Weigt, J. N. Onuchic, T. Hwa, H. Szurmant, High-resolution protein
478 complexes from integrating genomic information with molecular simulation. *Proc Natl*
479 *Acad Sci U S A* **106**, 22124-22129 (2009).
480 46. M. Weigt, R. A. White, H. Szurmant, J. A. Hoch, T. Hwa, Identification of direct residue
481 contacts in protein-protein interaction by message passing. *Proc Natl Acad Sci U S A*
482 **106**, 67-72 (2009).
483 47. M. M. Desai, D. Weissman, M. W. Feldman, Evolution can favor antagonistic epistasis.
484 *Genetics* **177**, 1001-1010 (2007).
485 48. B. H. Good, M. M. Desai, The impact of macroscopic epistasis on long-term
486 evolutionary dynamics. *Genetics* **199**, 177-190 (2015).
487 49. E. R. Jerison, M. M. Desai, Genomic investigations of evolutionary dynamics and
488 epistasis in microbial evolution experiments. *Curr Opin Genet Dev* **35**, 33-39 (2015).
489 50. S. Kryazhimskiy, D. P. Rice, E. R. Jerison, M. M. Desai, Microbial evolution. Global
490 epistasis makes adaptation predictable despite sequence-level stochasticity. *Science*
491 **344**, 1519-1522 (2014).
492 51. J. I. Rojas Echenique, S. Kryazhimskiy, A. N. Nguyen Ba, M. M. Desai, Modular epistasis
493 and the compensatory evolution of gene deletion mutants. *PLoS Genet* **15**, e1007958
494 (2019).
495 52. Y. Wang, R. Lei, A. Nourmohammad, N. C. Wu, Antigenic evolution of human influenza
496 H3N2 neuraminidase is constrained by charge balancing. *Elife* **10** (2021).
497 53. R. E. Lenski, Experimental evolution and the dynamics of adaptation and genome
498 evolution in microbial populations. *ISME J* **11**, 2181-2194 (2017).
499 54. N. Strelkova, M. Lassig, Clonal interference in the evolution of influenza. *Genetics* **192**,
500 671-682 (2012).
501 55. B. H. Good, A. M. Walczak, R. A. Neher, M. M. Desai, Genetic diversity in the
502 interference selection limit. *PLoS Genet* **10**, e1004222 (2014).
503 56. E. Brunet, B. Derrida, A. H. Mueller, S. Munier, Effect of selection on ancestry: an
504 exactly soluble case and its phenomenological generalization. *Phys Rev E Stat Nonlin*
505 *Soft Matter Phys* **76**, 041104 (2007).
506 57. R. A. Neher, T. A. Kessinger, B. I. Shraiman, Coalescence and genetic diversity in sexual
507 populations under selection. *Proc Natl Acad Sci U S A* **110**, 15836-15841 (2013).
508 58. J. F. Kingman, Origins of the coalescent. 1974-1982. *Genetics* **156**, 1461-1463 (2000).

Unpublished manuscript. Not for distribution.

- 509 59. E. Brunet, B. Derrida, D. Simon, Universal tree structures in directed polymers and
510 models of evolving populations. *Phys Rev E Stat Nonlin Soft Matter Phys* **78**, 061102
511 (2008).
- 512 60. I. M. Rouzine, J. M. Coffin, Highly fit ancestors of a partly sexual haploid population.
513 *Theor Popul Biol* **71**, 239-250 (2007).
- 514 61. J. W. Yewdell, Antigenic drift: Understanding COVID-19. *Immunity* **54**, 2681-2687
515 (2021).
- 516 62. R. T. Eguia *et al.*, A human coronavirus evolves antigenically to escape antibody
517 immunity. *PLoS Pathog* **17**, e1009453 (2021).
- 518 63. N. D. Rochman *et al.*, Ongoing global and regional adaptive evolution of SARS-CoV-2.
519 *Proc Natl Acad Sci U S A* **118** (2021).
- 520 64. W. A. Haynes *et al.*, High-resolution epitope mapping and characterization of SARS-
521 CoV-2 antibodies in large cohorts of subjects with COVID-19. *Commun Biol* **4**, 1317
522 (2021).
- 523 65. A. J. Greaney *et al.*, Comprehensive mapping of mutations in the SARS-CoV-2 receptor-
524 binding domain that affect recognition by polyclonal human plasma antibodies. *Cell*
525 *Host Microbe* **29**, 463-476 e466 (2021).
- 526 66. A. J. Greaney *et al.*, Complete Mapping of Mutations to the SARS-CoV-2 Spike
527 Receptor-Binding Domain that Escape Antibody Recognition. *Cell Host Microbe* **29**, 44-
528 57 e49 (2021).
- 529 67. Anonymous, SARS-CoV-2 variants of concern as of 27 January 2022. *European Centre*
530 *for Disease Prevention and Control* <https://www.ecdc.europa.eu/en/covid-19/variants-concern>
531 (2022).
- 532 68. J. H. Tay, A. F. Porter, W. Wirth, S. Duchene, The Emergence of SARS-CoV-2 Variants
533 of Concern Is Driven by Acceleration of the Substitution Rate. *Mol Biol Evol* **39** (2022).
- 534 69. S. P. Otto *et al.*, The origins and potential future of SARS-CoV-2 variants of concern in
535 the evolving COVID-19 pandemic. *Curr Biol* **31**, R918-R929 (2021).
- 536 70. S. A. Kemp *et al.*, SARS-CoV-2 evolution during treatment of chronic infection. *Nature*
537 **592**, 277-282 (2021).
- 538 71. L. Corey *et al.*, SARS-CoV-2 Variants in Patients with Immunosuppression. *N Engl J Med*
539 **385**, 562-566 (2021).
- 540 72. B. Choi *et al.*, Persistence and Evolution of SARS-CoV-2 in an Immunocompromised
541 Host. *N Engl J Med* **383**, 2291-2293 (2020).
- 542 73. I. M. Rouzine, J. M. Coffin, Search for the mechanism of genetic variation in the pro
543 gene of human immunodeficiency virus. *J Virol* **73**, 8167-8178 (1999).
- 544 74. A. Ignatieva, J. Hein, P. A. Jenkins, Ongoing Recombination in SARS-CoV-2 Revealed
545 Through Genealogical Reconstruction. *Mol Biol Evol* 10.1093/molbev/msac028
546 (2022).
- 547 75. B. Jackson *et al.*, Generation and transmission of interlineage recombinants in the
548 SARS-CoV-2 pandemic. *Cell* **184**, 5179-5188 e5178 (2021).
- 549 76. Y. Turkahia *et al.*, Pandemic-Scale Phylogenomics Reveals Elevated Recombination
550 Rates in the SARS-CoV-2 Spike Region. *bioRxiv* 10.1101/2021.08.04.455157 (2021).
- 551 77. H. Yi, 2019 Novel Coronavirus Is Undergoing Active Recombination. *Clin Infect Dis* **71**,
552 884-887 (2020).
- 553 78. D. VanInsberghe, A. Neish, A. Lowen, K. Koelle, Recombinant SARS-CoV-2 genomes
554 circulated at low levels over the first year of the pandemic. *Virus Evolution* **7**, veab059
555 (2021).

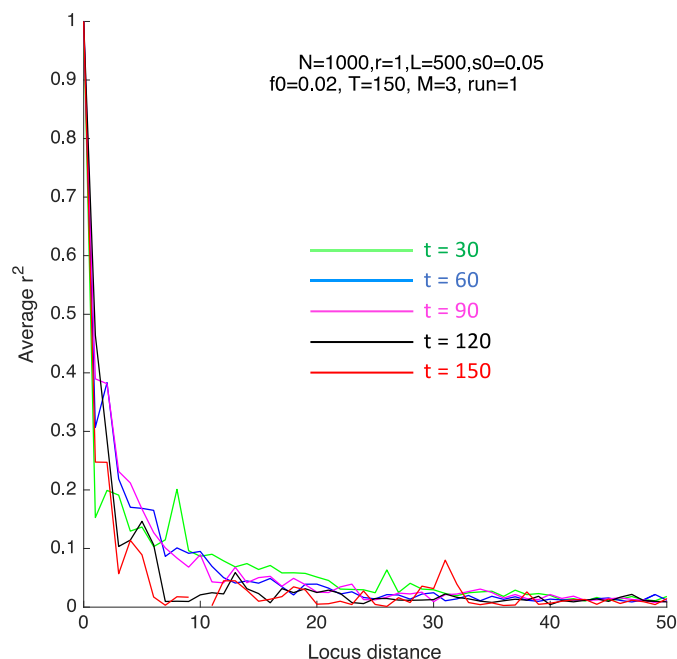
Unpublished manuscript. Not for distribution.

- 556 79. I. M. Rouzine, J. M. Coffin, Multi-site adaptation in the presence of infrequent
557 recombination. *Theor Popul Biol* **77**, 189-204 (2010).
- 558 80. M. S. Y. Lau *et al.*, Characterizing superspreading events and age-specific
559 infectiousness of SARS-CoV-2 transmission in Georgia, USA. *Proc Natl Acad Sci U S A*
560 **117**, 22430-22435 (2020).
- 561 81. D. C. Adam *et al.*, Clustering and superspreading potential of SARS-CoV-2 infections in
562 Hong Kong. *Nat Med* **26**, 1714-1719 (2020).
- 563 82. A. Gomez-Carballa, X. Bello, J. Pardo-Seco, F. Martinon-Torres, A. Salas, Mapping
564 genome variation of SARS-CoV-2 worldwide highlights the impact of COVID-19 super-
565 spreaders. *Genome Res* **30**, 1434-1448 (2020).
- 566 83. Y. Liu, R. M. Eggo, A. J. Kucharski, Secondary attack rate and superspreading events for
567 SARS-CoV-2. *Lancet* **395**, e47 (2020).
- 568 84. R. N. Dutta, I. M. Rouzine, S. D. Smith, C. O. Wilke, I. S. Novella, Rapid adaptive
569 amplification of preexisting variation in an RNA virus. *J Virol* **82**, 4354-4362 (2008).
570
571

Unpublished manuscript. Not for distribution.

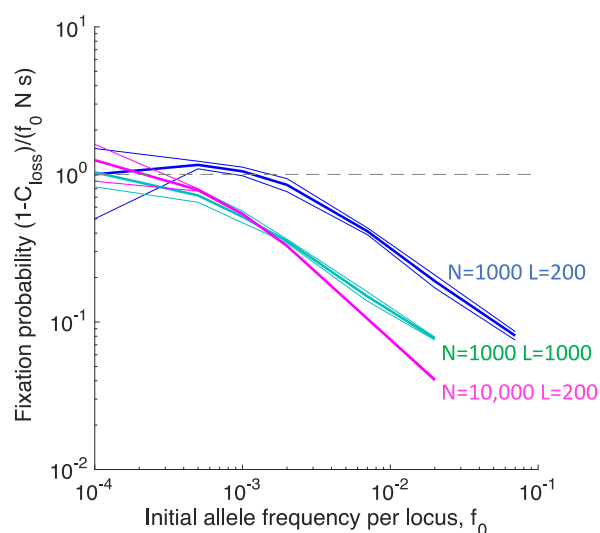
572 Supporting Information

573



574

575 **S1 Fig. Linkage disequilibrium as a function of distance between loci in the genome.** Pearson's
 576 measure r^2 is averaged over pairs of sufficiently heterozygous loci, $2f_{loc}(1 - f_{loc}) > 0.1$. The time points
 577 and parameters (shown) are the same as in Fig. 3. At $t = 0$, linkage disequilibrium is identically zero due
 578 to the initial random distribution of alleles.
 579



580

581 **S2 Fig. Fixation probability per beneficial allele as a function of the initial allelic frequency exhibits**
 582 **transition to the dilute limit of independent alleles with fixation probability s .** Y-axis: The average
 583 fraction of surviving polymorphic loci, $1 - C_{loss}(\infty)$, divided by $f_0 N s$, which is the product of the average
 584 number of beneficial alleles per locus, $f_0 N$, and the allelic fixation probability in the 1-locus model, s . X-axis:
 585 The initial frequency of beneficial alleles, f_0 . The dependence is shown at three combinations of values of
 586 N and L . Three lines of each color show the mean and the mean plus minus the standard deviation between
 587 three simulation runs, i.e., the 67% confidence interval. The independent-locus limit of fixation probability
 588 shown by the dashed horizontal line is reached at small f_0 . Fixed parameters are $M = 3$ and $s = 0.1$.
 589

589