

Epistasis detectably alters correlations between genomic sites in a narrow parameter window

Gabriele Pedruzzi¹ and Igor M. Rouzine^{1*}

¹*Sorbonne Université, Institute de Biologie Paris-Seine, Laboratoire de Biologie Computationnelle et Quantitative, LCQB, F-75004 Paris, France*

* To whom correspondence should be addressed: igor.rouzine@sorbonne-universite.fr

Short title: Epistasis and linkage disequilibrium

Abstract

Different genomic sites evolve inter-dependently due to the combined action of epistasis, non-additive contributions of different loci to genome fitness, and physical linkage of different loci due to their common heritage. Both epistasis and linkage, partially compensated by recombination, cause correlations between allele frequencies at the loci (linkage disequilibrium, LD). The interaction and competition between epistasis and linkage are not fully understood, nor is their relative sensitivity to recombination. Modeling an adapting population in the presence of random mutation, natural selection, pairwise epistasis, and random genetic drift, we compare the contributions of epistasis and linkage. For this end, we use a panel of haplotype-based measures of LD and their various combinations calculated for epistatic and non-epistatic pairs separately. We compute the optimal percentages of detected and false positive pairs in a one-time sample of a population of moderate size. We demonstrate that true interacting pairs can be told apart in a sufficiently short genome within a narrow window of time and parameters. Outside of this parameter region, unless the population is extremely large, shared ancestry of individual sequences generates pervasive stochastic LD for non-interacting pairs masking true epistatic associations. In the presence of sufficiently strong recombination, linkage effects decrease faster than those of epistasis, and the detection of epistasis improves. We demonstrate that the epistasis component of locus association can be isolated, at a single time point, by averaging haplotype frequencies over multiple independent populations. These results demonstrate the existence of fundamental restrictions on the protocols for detecting true interactions in DNA sequence sets.

32 Introduction

33

34 Epistasis is inter-dependence of fitness effects of mutations occurring at different loci caused by
35 biological interactions between domains of proteins and between proteins and nucleic acids [1-4]. In
36 biological systems, amino acids in proteins domains interact with each other. The resulting networks of
37 interactions that include direct protein-protein binding and allosteric effects, shape the gene regulation and
38 metabolic networks. Epistasis is a widespread property of biological networks [2, 5-8] and a subject of
39 intense studies. The vital role it plays in the genetic evolution of populations and the heritability of complex
40 traits is well established. The existing estimates indicate that the variation of an inherited trait across a
41 population can only partially be explained by the additive contributions from the relevant alleles. On
42 average, 70% of the inheritance may be due to epistasis or epigenetic effects [9]. Epistasis defines the
43 evolutionary paths and creates fitness valleys, i.e., intermediate genetic variants with reduced fitness [10-12].

44 A crucial biological scenario is a viral population adapting to the abrupt changes in external
45 conditions. Examples include the transmission to a new host, the invasion of a new organ, or the process of
46 immune evasion or the development of drug resistance. Typically, virus adaptation consists of primary
47 mutations followed by a cascade of several compensatory (helper) mutations [13-18]. These mutations help
48 the adapting virus to pass through a fitness valley [11]. During this process, compensatory mutations rescue
49 the replicative fitness of virus while preserving its resistant phenotype [13, 15, 19].

50 However, epistasis is not the only force causing inter-dependence in the evolution of genomic
51 regions. The other dominant factor is the host of linkage effects existing between genomic regions that co-
52 evolve in the same time frame and share the same ancestors [20, 21]. They include Fisher-Muller effect
53 (clonal interference), genetic hitchhiking and genetic background effects, and Hill-Robertson interference
54 between genetic drift and selection [21-23].

55 The other effect of linkage is a genetic association between loci, or linkage disequilibrium (LD).
56 The effects of linkage on the evolution of a long genome in the presence of selection is well understood
57 theoretically [12, 24-31]. The theory shows that linkage significantly slows adaptation many times, enhances
58 accumulation of deleterious mutations, and changes the shape of the phylogenetic tree [32, 33]. The
59 magnitude of linkage effects grows rapidly with the number of loci, L . Recombination partly offsets linkage
60 effects and accelerates evolution [34-40] and competes with epistasis [41]. Epistasis has been shown to be
61 potentially important for the evolution of recombination in a two-locus model [42, 43].

62 One consequence of linkage at large L is the strong interaction between the evolutionary trajectories
63 of different sites that, depending on the case, can be both positive and negative. LD stemming from this
64 interaction is easy to confuse with epistasis effects. Linkage effects become small only in populations that
65 are exponentially large in the number of sites L [25]. Further, working with sequence data from real
66 populations, it is often unclear how to discriminate the effects of shared ancestry from those of epistasis, and
67 which of the two evolutionary forces dominates in each case (for a comprehensive review, see [1, 44, 45]).
68 Therefore, despite of a considerable theoretical and experimental effort, detecting epistasis from genomic

69 data remains a challenge.

70 In the present work, we offer an evolutionary explanation for the observed difficulty of the detection
71 of epistasis from one-time data set. The idea is to generate mock data using a Monte-Carlo model of
72 evolution and then try to discriminate between effects of linkage and epistasis. We use a panel of six
73 pairwise LD measures to compare their distributions between epistatic and random pairs in a broad range of
74 model parameters. We also use 3D and 2D maps of all possible combinations of LD measures and employ an
75 optimization algorithm based on *a priori* knowledge to estimate the best, theoretically possible identification
76 of epistatic pairs. As a result, we delineate the region of time and model parameters where the epistatic pairs
77 can be detected against the linkage background. Finally, we investigate the role of recombination and the
78 effects of averaging over multiple independently-evolving populations.

79

80 **Results**

81

82 **Computer simulation of evolution**

83 We consider a haploid population of N genomic sequences comprised of L sites, where $L \gg 1$, and either a
84 favorable or deleterious allele is present at each site. Evolution of the population between discrete
85 generations is simulated using a Wright-Fisher model including the evolutionary factors of random mutation
86 with the rate μ per site, random genetic drift, and natural selection, as described in *Methods*. Natural
87 selection includes positive (antagonistic) epistatic interaction between selected pairs of deleterious alleles. A
88 simple case of genomes with uniform selection coefficient s_0 and uniform epistatic strength, E , is considered.
89 We also assume that epistatic pairs are isolated, i.e., that each genomic site interacts with only one site. The
90 initial population is randomized as it is done in virus passage experiments, with an average allelic frequency
91 f_0 . In most of our work, we initially neglect the factor of recombination and primarily focus on asexual
92 evolution, but lift this restriction in the end and explore broad parameter ranges. We aim to simulate the
93 detection of epistatic pairs and identify the best conditions for detection theoretically.

94

95 **Measures of linkage disequilibrium (LD)**

96 Various haplotype-based measures based on known haplotype frequencies have been proposed to
97 characterize the allelic association between loci. We will list four measures, as follows.

98 *Lewontin's measure.* A classical measure of statistical correlation between alleles at different loci
99 has a form [46]

100

$$101 \quad D' = \frac{D}{D_{max}}, \quad D = f_{ij} - f_i f_j \quad (1)$$

$$102 \quad D_{max} = \begin{cases} \max\{-f_i f_j, -(1-f_i)(1-f_j)\}, & D < 0 \\ \min\{f_i(1-f_j), (1-f_i)f_j\}, & D > 0 \end{cases}$$

103

104 Here f_{ij} is the average frequency of a bi-allelic haplotype of loci i and j , and D_{max} is a normalization
105 coefficient making sure that $D' \in [0, 1]$.

106 *Pearson's correlation coefficient.* An alternative is the correlation coefficient between pairs of loci r ,
107 expressed as [47]

$$108 \quad r = D / \sqrt{f_i(1 - f_i)f_j(1 - f_j)} \quad (2)$$

109
110 *Kimura-Wu measure.* More recently, Wu and colleagues proposed another statistical marker of
111 linkage disequilibrium, which, for binary alleles, has a form [48]

$$112 \quad WU = \log \frac{f_{11}f_{00}}{f_{01}f_{10}} \quad (3)$$

114
115 which represents the logarithm of the Z-measure proposed much earlier by Kimura [49].

116 *Universal footprint of epistasis.* In our recent work [50], we introduced another bi-allelic measure of
117 LD

$$118 \quad E = 1 - \frac{\log(f_{11}/f_{00})}{\log(f_{01}f_{10}/f_{00}^2)} \quad (4)$$

120 The advantage of this measure with respect to previous three is that it has a direct meaning in terms of fitness.
121 For isolated interacting pairs, it represents the degree of mutual compensation of two deleterious mutations
122 when frequencies in Eq. 4 are ensemble-averaged (see *Methods* below). Here the value $E = 0$ corresponds to
123 the absence of compensation (epistasis), and $E = 1$ to full mutual compensation of the two mutations. Note
124 the singularity in Eq. 4 at $f_{10}f_{01} = f_{00}^2$; we checked that it does not affect our results.

125 Below we investigate the effect of linkage for interacting and noninteracting pairs of loci using the
126 measures defined in Eqs. 1-4. Also, we employ an optimization algorithm that, exploiting a priori knowledge
127 of the correct epistatic pairs, puts the best possible threshold between the two distributions of LD. We
128 consider different combinations of two or three LD measures to obtain the best detection possible.

129 **LD of epistatic and non-epistatic pairs are distinct in a narrow parameter window**

131 We started by plotting the distribution of six LD measures calculated from Eq. 1 over individual pairs of sites,
132 at different times (Fig. 1). We show separately the distribution for two subsets of pairs: the known epistatic
133 subset (dark shade) and all the pairs (light shade). In the beginning, LD is narrowly distributed around zero,
134 for both epistatic and non-epistatic subsets (Fig. 1, row 1).

135
136 **Fig. 1. LD- and haplotype-based measures of epistasis identify a narrow time window of epistasis detectability.** We compared
137 the time-dependent distribution of 6 markers of LD shown in 6 columns. Each column show the profile of the distribution of a
138 measure of epistasis: D_{11} , D_{01} (Eq. 1), r_{11} , r_{01} (Eq. 2), WU (Eq. 3) and UFE (Eq. 4). Different rows correspond to different time
139 points: $t = 1$, $t = 5$, $t = 10$, $t = 25$ and $t = 50$. The shaded regions correspond to the density distributions for all possible pairwise
140 interactions (lighter color) and the known epistatic pairs (darker shade). The shaded areas are normalized distributions reflecting the
141 fact that epistatic pairs represent a tiny fraction of the all possible pairs in a genome. The fluctuations of non-epistatic pairs increasing
142 in time overlap onto the distributions of epistatic pairs. Parameters: $N = 2 \cdot 10^4$, $s_0 = 0.1$, $L = 50$, E in the range $[0, 1]$, $\mu L = 7 \cdot 10^{-2}$.

143 Each odd site interacts with its neighbor on the right (1-2, 3-4, 5-6, ...) with epistatic strength $E = 0.75$. Initially, sequences were
144 random with average allelic frequency set to $f = 0.4$. The negative control result in the absence of epistasis ($E = 0$) is presented on
145 Supplementary Fig. S1.
146

147 Subsequent time points (Fig 1, rows 2 and 3) show progressive separation of the two distributions. In the
148 course of further evolution (Fig. 1, rows 4 and 5), the distribution of randomly-chosen pairs, which was
149 initially narrow and concentrated near the origin $E = 0$, gradually expands and overlaps with the small
150 epistatic distribution (Fig 1). This effect implies that non-epistatic pairs of sites, due to the stochastic nature
151 of evolution, produce large LD of random sign. In this case, it is impossible to tell apart epistatic pairs from
152 any of these measures of LD.
153

154 **Results are robust to the choice of an LD measure or their combination**

155 Next, we checked whether combinations of LDs used together can improve detection. We have calculated all
156 possible combination of six LD measures in Eq. 2 and tried to separate interacting and non-interacting pairs
157 using 3D and 2D scatter plots. A representative example is shown in Fig. 2, for $E=0$, and for $E=0.75$ at two
158 time points. Other possible combinations of 2 and 3 measures are summarized in Table S1 in *Supplement*.
159

160 **Fig. 2. The optimization algorithm to identify ideal conditions for detection of epistasis is exemplified through the 3D scatter**
161 **plot of three different measures of LD.** Left: A representative example of three-dimensional scatter plots of three statistics, UFE,
162 D'_{01} and r_{11} , plotted for all individual pairs of sites (blue circles) and for known epistatic pairs (red circles). Right and middle: two-
163 dimensional projections. The three rows correspond to the absence of epistasis ($E = 0$, top), and two time points in the presence of
164 epistasis, within the detection window and outside (middle and bottom). All possible combinations of two and three measures have
165 been tested and summarized in Table S1. At intermediate time $t=10$, a distinct cloud of epistatic pairs (red dots) cluster together
166 outside the overall distribution of all pairs and, hence are detectable. At long times, substantial overlap with non-interacting pairs bias
167 contaminates detection. To optimize detection, we define a detection threshold for each of the detection variables (UFE, WU and the
168 four haplotypes) and adopted an optimization algorithm that minimizes the following quantity "DET + a FPOS", where a is a fitting
169 parameter, DET represent the detection percentage, and FPSO is the percentage of false positive, based on prior knowledge of the
170 identity of true epistatic pairs. Parameters as in Fig. 1
171

172 We wrote an optimization algorithm which separates the cloud of interacting pairs from the cloud of non-
173 interacting pairs in the best possible way, using *a priori* knowledge about the identity of pairs (Fig. 2). We
174 adjusted the threshold to optimize the difference between the detection rate and the false positive rate. This
175 method, employing the principle of machine learning, does not give any substantial improvement on the
176 detection window (See Supplementary Table 1). For a real data sets, *a priori* knowledge about interacting
177 pairs is usually unavailable, so that the detection of epistasis in a single population at one time point will be
178 even worse than our prediction.
179

180 **Parameter sensitivity analysis confirms the narrow window of detection**

181 *Selection coefficient.* Next, we investigated how the window of detection changes with model parameters.
182 We calculated the detection rate and the false positive rate for the six measures of LD at different values of

183 selection coefficient, s_0 (Fig. 3). For each measure, the results show an inverse scaling of the detection time
184 window on s_0 . Note that the window closes at very small s_0 , where evolution is almost selectively neutral,
185 and epistasis is never detectable.

186

187 **Fig. 3. Detection of epistasis is confined in a time window whose width is controlled by the mean selection coefficient.**
188 Percentile of detection and false discovery as a function of time is averaged over 25 random simulation runs per each value of s_0 , the
189 constant selection coefficient for each allele in the sub-population. The detection of epistatic pairs for a panel of measures of LD,
190 namely, D_{11} , D_{01} (Eq. 1), r_{11} , r_{01} (Eq. 2), WU (Eq. 3) and UFE (Eq. 4). Results from a detection protocol that maximizes the
191 difference between the detection percentile and the false-positive fractions by tuning the detection threshold, show the same trend for
192 all measures considered. At time $\sim 1.5/s_0$ generations, we observe the beginning of a transition which completely blurs the detection
193 of epistatic interaction at time $\sim 2.5/s_0$. The initial allelic frequency $f_0 = 0.45$, s_0 is shown, the other parameters are as in Fig. 1.
194

195

196 *Distributed selection coefficient.* Next, we conducted a sensitivity analysis with respect to the other
197 model parameters (Fig. S5). Firstly, we lifted the simplifying assumption of a constant selection coefficient, s
198 = s_0 , and allowed variation of s among sites according to a half-Gaussian distribution. We obtain a similar
199 dependence of the window width on the average selection coefficient (Fig. S5), although with a higher false
200 positive rate within the detection window than for the case with constant s .

201 *Length of the genome.* We found out, that sequence length L limits the detectability of epistasis
202 substantially (Fig. S5). An increase of the sequence length or a reduction of the population size leads to
203 narrowing and, eventually, disappearance of the detection window. These results limit the applicability of
204 these methods to short sequences. Indeed, the number of all possible locus pairs increases with genome
205 length L proportionally to L^2 , and the number of epistatic pairs increases only as L , so that the task of finding
206 "the ruby in the rubbish" becomes harder at larger L [1, 44, 45].

207 *Population size.* We observed a very slow (logarithmic) expansion of the detection window with
208 population size N (Fig. S5). This is consistent with the results of asexual evolution models, which predict a
209 very slow logarithmic dependence on N for all the evolutionary observables, including evolution speed,
210 genetic diversity, and the average time to most recent ancestor [25-31, 35-37, 39, 40, 51]. Only in very large
211 populations whose size increases exponentially genome length L , linkage effects become small [25]. In these,
212 astronomically large populations, epistasis would be easily detectable.

213 *Initial standing variation.* We have observed a detection window in time only at the initial
214 frequencies of deleterious alleles above 10% (Fig. S5). At smaller frequencies, detection lapses. We can
215 conclude that detection of epistasis in a single population studied is possible in a narrow parameter range.
216

217 **Recombination improves detection.** Until now, we have assumed a completely asexual evolution. In our
218 next step, we investigated the role of recombination, parametrised by the average number of crossovers per
219 genome, M , and by the probability of outcrossing per genome, r . We obtained that intermediate
220 recombination rates rescue the detection of epistasis by disrupting linkage and yet preserving the epistasis
221 contribution to LD. At our default parameter set (Fig. 1), we observed a significant reduction of linkage
222 fluctuations starting from $r = 20\%$ and $M = 5$ (Fig. 4). The results show that LD effects of linkage are much
223 more resistant to recombination than, for example, the evolution speed, which increases substantially already

224 at tiny values of r [34-40]. We found out also that extremely high levels of recombination decrease LD for
225 epistatic pairs as well, thus rendering epistasis undetectable. Thus, there exists a narrow window of
226 recombination rates where epistasis can be observed outside of the detection window for time and other
227 parameters described above.

228

229 **Fig. 4. Variation of the time window of detection with recombination.** Percentile of detection and false discovery as a function of
230 time is averaged over 25 random simulations (runs) in a broad range of parameters values. The detection rare and false positive rate
231 of epistatic pairs with UFE at different values of s , randomly drawn from a half-Gaussian distribution of deleterious alleles. The
232 presence of moderate recombination characterized by outcrossing rate r and the average number of cross-overs, M , broadens the
233 detection window. We observe similar results for all the statistics considered in this study (data not shown). The default parameter set
234 is $E = 0.75$, with the other parameters as in Fig. 1.

235

236

237 **Population divergence creates strong linkage effects**

238 In order to understand the reason behind the strong linkage effects masking epistasis, we investigated the
239 time-dependent changes of the phylogenetic tree using a hierarchical clustering algorithm (Fig. 5a-d). The
240 initial, randomized population display a star-shaped phylogeny, characterized by the same mean distance
241 between all sequences and the most common sequence (Fig. 5). With time, the phylogenetic tree grows
242 branches of increasingly related sequences (Fig. 5c, d). As simulation continues (Fig. 5d), the tree becomes
243 more lopsided, while recent mutations create short branches at the bottom. At the same time, we observe that
244 the tree has a decreasing number of ancestors. Eventually, the tree evolves into Bolthausen-Sznitman
245 coalescent (BSC) with a single common ancestor, previously predicted for the stationary regime of traveling
246 wave [25, 29, 37] (Fig. 5).

247

248 **Fig. 5. Evolution of genealogy within a single, well-mixed population and comparative representation of multiple,
249 independently evolving population.** (a-d) Phylogenetic structure of a single population comprising a sample of 500 genomes at four
250 different times: $t = 0, 10, 20,$ and 30 generations. Mean genetic distance between genomes decreases in time, and the structure of the
251 tree changes from a star-like shape towards a monophyletic tree (BS coalescent), with a single common ancestor. The right panel
252 shows the reconstructed phylogenetic tree of three populations, independently evolved from the same initial random seed. At a glance,
253 it is possible to determine that the three populations do not share much sequence homology and segregate into different,
254 phylogenetically distinct clades. $N = 20000$ genomes, initial average allelic frequency $f_0 = 0.40$, other parameters as in Fig. 1.

255

256 Emergence of this phylogeny is coincident with the increase in the fluctuations of LD of non
257 interacting pairs (Fig. 1). The reason for strong random LD is stochastic divergence of the population from
258 the initial state, as illustrated by clustering of three independently evolved populations (Fig. 5, right). The
259 distance between the trees obtained in separate runs increases linearly in time due to fixed beneficial
260 mutations at randomly chosen sites. Haplotype configurations of the common ancestor of the population are
261 inherited by all members of the population, with some small variation determined by the time to the most
262 recent common ancestor. Thus, the stochastic divergence of individual populations creates strong LD with a
263 random sign.

264

265

266 **The use of multiple populations defeats LD fluctuations and rescues epistatic signature**

267 Because the linkage fluctuations arise due to stochastic divergence of the founder, the common ancestor, the
268 natural idea is to use multiple populations to average over possible founder sequences. To test this idea, we
269 evolved independently multiple populations at the same initial conditions and averaged the haplotype
270 frequencies used in LD markers (Eqs. 1-4) over populations, for each pair of sites, separately. We found out
271 that including a sufficient number of independent populations results in a substantial reduction of the noise
272 and indefinite expansion of the window of detection (Fig. 6). Qualitatively similar results are obtained for all
273 LD markers.

274

275 **Fig. 6. Detection of epistasis is rescued by simultaneous analysis of multiple independently-evolved populations.** Left 4 plots:
276 Percentile of detection (top) and false discovery (bottom) as a function of time are presented for UFE and WU measures. Number of
277 replicate Monte-Carlo runs is shown. The haplotype frequencies are averaged over runs, which represent independently-evolved
278 populations. At time $\sim 1.5/s_0$, we observe the beginning of a transition which completely blurs the detection of epistatic interaction for
279 a single replicate (blue line), however, already 5 replicates are sufficient to significantly extend the detection window up to $\sim 2.5/s_0$,
280 and a higher number of replicates completely eliminate false-positive pairs, while maintaining the average detection above 80%.
281 Parameters: $E = 0.75$, $N = 20000$, the others as in Fig 1. Right: Two-dimensional color maps for UFE measure of LD, which
282 summarize the results of a similar analysis for two population sizes: $N = 100$ (middle plot) and $N = 1000$ (right plot). Y-axis: Number
283 of independent populations. X-axis: time of evolution. Color shows the percentage of detection with the detection threshold of
284 interacting pairs chosen to give the false discovery rate below 20%.
285

286

287 DISCUSSION

288 In the present work, using a Monte-Carlo simulation of a haploid population, we calculated the distributions
289 of six measures of linkage disequilibrium and their combinations for epistatic and random locus pairs. We
290 demonstrated that, in a single asexual population, the footprints of epistatic pairs are readable only in a
291 narrow time interval between $0.2/s_0$ and $1.5/s_0$ generations. During later adaptation, the distribution of
292 linkage disequilibrium for non-interacting pairs broadens and engulfs the distribution for epistatic pairs.
293 These results indicate that, long before the onset of the steady state, linkage effects dominate over the effects
294 of epistasis. This phenomenon is predicted in a broad parameter region and for all the LD statistics,
295 suggesting that, in the context of inherited linkage fluctuations, all statistics based on pairwise linkage
296 disequilibrium are equal.

297 To gain insight into the evolutionary origin of these fluctuations, we investigated phylogenetic trees
298 of the entire population at different time points to observe that the shape of the tree strongly correlates with
299 the magnitude of linkage fluctuations. The shape of the phylogenetic tree changes in time from the initially
300 star-shaped genealogy to a Bolthausen-Sznitman (BS) coalescent [32, 33] previously analyzed in great detail
301 for adapting asexual populations [25, 36, 37]. Once BS genealogy is established, individual sequences share
302 a high degree of interrelatedness due to fixed beneficial mutations at randomly chosen sites. The presence of
303 the BS coalescent is coincident with strong co-inheritance linkage fluctuations. The stochastic nature of their
304 common ancestor sequence, divergent in time from common ancestors in other independent populations (Fig.
305 5) is the cause of the strong fluctuations of LD.

306 We have also directly quantitated the detection of epistatic pairs against the background of random
307 linkage effects. We evaluated the sensitivity of the width of the detection window with respect to several

308 input parameters, such as the mean selection coefficient, the size of the population, the sequence length, and
309 initial genetic variation, and the role of recombination. We observed that the window is proportional to the
310 inverse average selection coefficient, $1/s_0$, but a very small s_0 abolishes any chance of detection, so that the
311 best detection is attained in the case of moderately weak selection. The detection window exists only for
312 sufficiently small genomes. The presence of recombination has the effect of compensating the linkage
313 component and thus significantly improving the detection of epistasis. Yet, very frequent recombination
314 disrupts epistatic effects.

315 To isolate the epistatic component from co-inheritance effects, we performed simulations over
316 several independently-evolved populations and averaged the haplotype frequencies over these runs. The
317 results predict the number of independent population required to attain significant expansion of the detection
318 window (Fig. 6). Thus, the averaging over multiple independently-evolved populations filters out linkage
319 effects leaving a clear footprint of epistasis in a much broader parameter range. However one should note
320 that the multiple-population sampling was conducted under the ideal conditions, in which every population
321 evolved independently for the same time with the same parameter set, and represented the same fraction of
322 the total sample. Unequal sampling or heterogeneous representation in real data sets may create additional
323 problems.

324 Our model adopts several simplifying assumptions. (i) Deleterious alleles are assigned selection
325 coefficient constant in time. (ii) We considered constant and fixed epistatic strength for all pairs. (iii) We
326 focused on a simple topology of epistatic network. While these are reasonable assumptions to describe the
327 problem of linkage fluctuations in biological systems, a real scenario with mixed sign epistasis and complex
328 topology might pose additional challenges for the accurate detection of epistasis.

329 In summary, we offer an evolutionary reason for the fluctuations of epistatic estimates in the
330 existing sequence sets. Linkage due to stochastic divergence of the common ancestor of a population from
331 the origin is responsible for the high false-positive rates of epistasis detection in a single population. We
332 demonstrated how the use of multiple independently-evolving populations, or the use of time series when
333 available, allows us to average out strong linkage effects and rescue the detectability of epistasis.

334

335 **METHODS**

336 We consider a haploid population of N binary sequences, where each genome site (nucleotide position)
337 numbered by $i=1, 2, \dots, L$ is either $K_i=0$ or $K_i=1$. We assume that the genome is long, $L \gg 1$. Evolution of
338 the population in discrete time measured in generations is simulated using a standard Wright-Fisher model,
339 which includes the factors of random mutation with rate μL per genome, natural selection, and random
340 genetic drift. Recombination is assumed to be absent. Once per generation, each genome is replaced by a
341 random number of its progeny which obeys multinomial distribution. The total population stays constant
342 with the use of the broken-stick algorithm. To include natural selection, we calculate fitness (average
343 progeny number) e^W of sequence K_i as given by [50]

344

345
$$W = \sum_{i=1}^L s_i K_i + \sum_{i<j}^L S_{ij} K_i K_j \quad (5)$$

346

347
$$S_{ij} = E_{ij}(|s_i| + |s_j|)T_{ij} \quad (6)$$

348

349 The first term in Eq. 5 stands for the additive contribution of single mutations to fitness with selection
350 coefficients s_i . The second term in Eq. 5 describes pairwise interactions of sites with magnitudes S_{ij} , which
351 are given by Eq. 6. Coefficient E_{ij} represents the relative strength of epistatic interaction between sites i and
352 j , while the binary elements of the matrix \mathbf{T} indicate the interacting pairs by $T_{ij} = 1$ and the other pairs by T_{ij}
353 $= 0$. An example of positive epistasis is the compensation of two deleterious mutations inside protein
354 segments that bind each other. Note that $E_{ij} = 1$ corresponds to full mutual compensation of deleterious
355 mutants at sites i and j . We consider the simplest interaction topology of interacting neighbors, as given by
356 $T_{2i,2i+1} = 1$ and 0 for all other pairs.

357

358 **Conflict of Interest**

359 None declared

360 **Acknowledgments**

361 We thank Martin Weight for insightful comments.

362 **Funding**

363 This work was supported by Agence Nationale de la Recherche grant J16R389 to IMR.

364

365 **REFERENCES**

366

367 1. Cordell HJ. Epistasis: what it means, what it doesn't mean, and statistical methods to detect it in
368 humans. *Hum Mol Genet.* 2002;11(20):2463-8. Epub 2002/09/28. PubMed PMID: 12351582.

369 2. Combarros O, Cortina-Borja M, Smith AD, Lehmann DJ. Epistasis in sporadic Alzheimer's disease.
370 *Neurobiol Aging.* 2009;30(9):1333-49. Epub 2008/01/22. doi: 10.1016/j.neurobiolaging.2007.11.027.
371 PubMed PMID: 18206267.

372 3. Combarros O, van Duijn CM, Hammond N, Belbin O, Arias-Vasquez A, Cortina-Borja M, et al.
373 Replication by the Epistasis Project of the interaction between the genes for IL-6 and IL-10 in the risk of
374 Alzheimer's disease. *J Neuroinflammation.* 2009;6:22. Epub 2009/08/25. doi: 10.1186/1742-2094-6-22.
375 PubMed PMID: 19698145; PubMed Central PMCID: PMC2744667.

376 4. Bullock JM, Medway C, Cortina-Borja M, Turton JC, Prince JA, Ibrahim-Verbaas CA, et al.
377 Discovery by the Epistasis Project of an epistatic interaction between the GSTM3 gene and the

- 378 HHEX/IDE/KIF11 locus in the risk of Alzheimer's disease. *Neurobiol Aging*. 2013;34(4):1309 e1-7. Epub
379 2012/10/06. doi: 10.1016/j.neurobiolaging.2012.08.010. PubMed PMID: 23036584.
- 380 5. McKinney BA, Pajewski NM. Six Degrees of Epistasis: Statistical Network Models for GWAS.
381 *Front Genet*. 2011;2:109. Epub 2012/02/04. doi: 10.3389/fgene.2011.00109. PubMed PMID: 22303403;
382 PubMed Central PMCID: PMC3261632.
- 383 6. Steen KV. Travelling the world of gene-gene interactions. *Brief Bioinform*. 2012;13(1):1-19. Epub
384 2011/03/29. doi: 10.1093/bib/bbr012. PubMed PMID: 21441561.
- 385 7. Ritchie MD. Using biological knowledge to uncover the mystery in the search for epistasis in
386 genome-wide association studies. *Ann Hum Genet*. 2011;75(1):172-82. Epub 2010/12/17. doi:
387 10.1111/j.1469-1809.2010.00630.x. PubMed PMID: 21158748; PubMed Central PMCID:
388 PMC3092784.
- 389 8. Zhang Y, Jiang B, Zhu J, Liu JS. Bayesian models for detecting epistatic interactions from genetic
390 data. *Ann Hum Genet*. 2011;75(1):183-93. Epub 2010/11/26. doi: 10.1111/j.1469-1809.2010.00621.x.
391 PubMed PMID: 21091453.
- 392 9. Zuk O, Hechter E, Sunyaev SR, Lander ES. The mystery of missing heritability: Genetic interactions
393 create phantom heritability. *Proc Natl Acad Sci U S A*. 2012;109(4):1193-8. Epub 2012/01/10. doi:
394 10.1073/pnas.1119675109. PubMed PMID: 22223662; PubMed Central PMCID: PMC3268279.
- 395 10. Desai MM, Fisher DS, Murray AW. The speed of evolution and maintenance of variation in asexual
396 populations. *Curr Biol*. 2007;17(5):385-94. PubMed PMID: 17331728.
- 397 11. Weissman DB, Desai MM, Fisher DS, Feldman MW. The rate at which asexual populations cross
398 fitness valleys. *Theor Popul Biol*. 2009;75(4):286-300. Epub 2009/03/17. doi: 10.1016/j.tpb.2009.02.006.
399 PubMed PMID: 19285994; PubMed Central PMCID: PMC2992471.
- 400 12. Gerrish PJ, Lenski RE. The fate of competing beneficial mutations in an asexual population.
401 *Genetica*. 1998;102-103(1-6):127-44. Epub 1998/08/28. PubMed PMID: 9720276.
- 402 13. Gonzalez-Ortega E, Ballana E, Badia R, Clotet B, Este JA. Compensatory mutations rescue the virus
403 replicative capacity of VIRIP-resistant HIV-1. *Antiviral Res*. 2011;92(3):479-83. Epub 2011/10/27. doi:
404 10.1016/j.antiviral.2011.10.010. PubMed PMID: 22027647.
- 405 14. Handel A, Regoes RR, Antia R. The role of compensatory mutations in the emergence of drug
406 resistance. *PLoS Comput Biol*. 2006;2(10):e137. Epub 2006/10/17. doi: 10.1371/journal.pcbi.0020137.
407 PubMed PMID: 17040124; PubMed Central PMCID: PMC1599768.
- 408 15. Levin BR, Perrot V, Walker N. Compensatory mutations, antibiotic resistance and the population
409 genetics of adaptive evolution in bacteria. *Genetics*. 2000;154(3):985-97. Epub 2000/04/11. PubMed PMID:
410 10757748; PubMed Central PMCID: PMC1460977.
- 411 16. Nijhuis M, Schuurman R, de Jong D, Erickson J, Gustchina E, Albert J, et al. Increased fitness of
412 drug resistant HIV-1 protease as a result of acquisition of compensatory mutations during suboptimal
413 therapy. *AIDS*. 1999;13(17):2349-59. Epub 1999/12/22. PubMed PMID: 10597776.

- 414 17. Noviello CM, Lopez CS, Kukull B, McNett H, Still A, Eccles J, et al. Second-site compensatory
415 mutations of HIV-1 capsid mutations. *J Virol.* 2011;85(10):4730-8. Epub 2011/03/04. doi:
416 10.1128/JVI.00099-11. PubMed PMID: 21367891; PubMed Central PMCID: PMC3126181.
- 417 18. Piana S, Carloni P, Rothlisberger U. Drug resistance in HIV-1 protease: Flexibility-assisted
418 mechanism of compensatory mutations. *Protein Sci.* 2002;11(10):2393-402. Epub 2002/09/19. doi:
419 10.1110/ps.0206702. PubMed PMID: 12237461; PubMed Central PMCID: PMC3126181.
- 420 19. Cong ME, Heneine W, Garcia-Lerma JG. The fitness cost of mutations associated with human
421 immunodeficiency virus type 1 drug resistance is modulated by mutational interactions. *J Virol.*
422 2007;81(6):3037-41. Epub 2006/12/29. doi: 10.1128/JVI.02712-06. PubMed PMID: 17192300; PubMed
423 Central PMCID: PMC1865994.
- 424 20. Barton NH. Linkage and the limits to natural selection. *Genetics.* 1995;140(2):821-41. Epub
425 1995/06/01. PubMed PMID: 7498757; PubMed Central PMCID: PMC1206655.
- 426 21. Hill WG, Robertson A. The effect of linkage on limits to artificial selection. *Genet Res.*
427 1966;8(3):269-94. Epub 1966/12/01. PubMed PMID: 5980116.
- 428 22. Felsenstein J. The evolutionary advantage of recombination. *Genetics.* 1974;78(2):737-56. Epub
429 1974/10/01. PubMed PMID: 4448362; PubMed Central PMCID: PMC1213231.
- 430 23. McVean GA, Charlesworth B. The effects of Hill-Robertson interference between weakly selected
431 mutations on patterns of molecular evolution and variation. *Genetics.* 2000;155(2):929-44. PubMed PMID:
432 10835411; PubMed Central PMCID: PMC1461092.
- 433 24. Tsimring LS, Levine H, Kessler DA. RNA virus evolution via a fitness-space model. *Phys Rev Lett.*
434 1996;76(23):4440-3. doi: 10.1103/PhysRevLett.76.4440. PubMed PMID: 10061290.
- 435 25. Rouzine IM, Wakeley J, Coffin JM. The solitary wave of asexual evolution. *Proc Natl Acad Sci U S*
436 *A.* 2003;100(2):587-92. Epub 2003/01/15. doi: 10.1073/pnas.242719299. PubMed PMID: 12525686;
437 PubMed Central PMCID: PMC141040.
- 438 26. Brunet E, Rouzine IM, Wilke CO. The stochastic edge in adaptive evolution. *Genetics.*
439 2008;179(1):603-20. Epub 2008/05/22. doi: 179/1/603 [pii]
440 10.1534/genetics.107.079319. PubMed PMID: 18493075; PubMed Central PMCID: PMC2390637.
- 441 27. Rouzine IM, Brunet E, Wilke CO. The traveling-wave approach to asexual evolution: Muller's
442 ratchet and speed of adaptation. *Theor Popul Biol.* 2008;73(1):24-46. doi: 10.1016/j.tpb.2007.10.004.
443 PubMed PMID: 18023832; PubMed Central PMCID: PMC2246079.
- 444 28. Desai MM, Fisher DS. Beneficial mutation selection balance and the effect of linkage on positive
445 selection. *Genetics.* 2007;176(3):1759-98. Epub 2007/05/08. doi: 10.1534/genetics.106.067678. PubMed
446 PMID: 17483432; PubMed Central PMCID: PMC1931526.
- 447 29. Hallatschek O. The noisy edge of traveling waves. *Proc Natl Acad Sci U S A.* 2011;108(5):1783-7.
448 Epub 2010/12/29. doi: 10.1073/pnas.1013529108. PubMed PMID: 21187435; PubMed Central PMCID:
449 PMC3033244.

- 450 30. Good BH, Rouzine IM, Balick DJ, Hallatschek O, Desai MM. Distribution of fixed beneficial
451 mutations and the rate of adaptation in asexual populations. *Proc Natl Acad Sci U S A*. 2012;109(13):4950-5.
452 Epub 2012/03/01. doi: 10.1073/pnas.1119910109. PubMed PMID: 22371564; PubMed Central PMCID:
453 PMCPMC3323973.
- 454 31. Good BH, Desai MM. The impact of macroscopic epistasis on long-term evolutionary dynamics.
455 *Genetics*. 2015;199(1):177-90. Epub 2014/11/15. doi: 10.1534/genetics.114.172460. PubMed PMID:
456 25395665; PubMed Central PMCID: PMCPMC4286683.
- 457 32. Neher RA, Hallatschek O. Genealogies of rapidly adapting populations. *Proc Natl Acad Sci U S A*.
458 2013;110(2):437-42. Epub 2012/12/28. doi: 10.1073/pnas.1213113110. PubMed PMID: 23269838; PubMed
459 Central PMCID: PMCPMC3545819.
- 460 33. Walczak AM, Nicolaisen LE, Plotkin JB, Desai MM. The structure of genealogies in the presence of
461 purifying selection: a fitness-class coalescent. *Genetics*. 2012;190(2):753-79. doi:
462 10.1534/genetics.111.134544. PubMed PMID: 22135349; PubMed Central PMCID: PMCPMC3276618.
- 463 34. Batorsky R, Kearney MF, Palmer SE, Maldarelli F, Rouzine IM, Coffin JM. Estimate of effective
464 recombination rate and average selection coefficient for HIV in chronic infection. *Proc Natl Acad Sci U S A*.
465 2011;108(14):5661-6. Epub 2011/03/26. doi: 1102036108 [pii]
466 10.1073/pnas.1102036108. PubMed PMID: 21436045; PubMed Central PMCID: PMC3078368.
- 467 35. Gheorghiu-Svirshchevski S, Rouzine IM, Coffin JM. Increasing sequence correlation limits the
468 efficiency of recombination in a multisite evolution model. *Mol Biol Evol*. 2007;24(2):574-86. PubMed
469 PMID: 17138627.
- 470 36. Rouzine IM, Coffin JM. Evolution of human immunodeficiency virus under selection and weak
471 recombination. *Genetics*. 2005;170(1):7-18. PubMed PMID: 15744057.
- 472 37. Rouzine IM, Coffin JM. Multi-site adaptation in the presence of infrequent recombination. *Theor*
473 *Popul Biol*. 2010;77(3):189-204. PubMed PMID: 20149814.
- 474 38. Xiao Y, Rouzine IM, Bianco S, Acevedo A, Goldstein EF, Farkov M, et al. RNA recombination
475 enhances adaptability and is required for virus spread and virulence. *Cell Host Microbe*. 2016;19(4):493-503.
476 doi: 10.1016/j.chom.2016.03.009. PubMed PMID: 27078068.
- 477 39. Neher RA, Leitner T. Recombination rate and selection strength in HIV intra-patient evolution.
478 *PLoS Comput Biol*. 2010;6(1):e1000660. Epub 2010/02/04. doi: 10.1371/journal.pcbi.1000660. PubMed
479 PMID: 20126527; PubMed Central PMCID: PMC2813257.
- 480 40. Neher RA, Shraiman BI, Fisher DS. Rate of adaptation in large sexual populations. *Genetics*.
481 2010;184(2):467-81. PubMed PMID: 19948891.
- 482 41. Neher RA, Shraiman BI. Competition between recombination and epistasis can cause a transition
483 from allele to genotype selection. *Proc Natl Acad Sci U S A*. 2009;106(16):6866-71. Epub 2009/04/16. doi:
484 10.1073/pnas.0812560106. PubMed PMID: 19366665; PubMed Central PMCID: PMCPMC2672512.
- 485 42. Barton NH. A general model for the evolution of recombination. *Genet Res*. 1995;65(2):123-45.
486 PubMed PMID: 7605514.

- 487 43. Kouyos RD, Otto SP, Bonhoeffer S. Effect of varying epistasis on the evolution of recombination.
488 Genetics. 2006;173(2):589-97. doi: 10.1534/genetics.105.053108. PubMed PMID: 16547114; PubMed
489 Central PMCID: PMCPMC1526506.
- 490 44. Cordell HJ. Detecting gene-gene interactions that underlie human diseases. Nature reviews Genetics.
491 2009;10(6):392-404. doi: 10.1038/nrg2579. PubMed PMID: PMC2872761.
- 492 45. Wei W-H, Hemani G, Haley CS. Detecting epistasis in human complex traits. Nature Reviews
493 Genetics. 2014;15:722. doi: 10.1038/nrg3747
494 <https://www.nature.com/articles/nrg3747#supplementary-information>.
- 495 46. Lewontin RC. On measures of gametic disequilibrium. Genetics. 1988;120(3):849-52. PubMed
496 PMID: 3224810; PubMed Central PMCID: PMCPMC1203562.
- 497 47. Balding DJ. A tutorial on statistical methods for population association studies. Nat Rev Genet.
498 2006;7(10):781-91. doi: 10.1038/nrg1916. PubMed PMID: 16983374.
- 499 48. Wu X, Dong H, Luo L, Zhu Y, Peng G, Reveille JD, et al. A novel statistic for genome-wide
500 interaction analysis. PLoS Genet. 2010;6(9):e1001131. doi: 10.1371/journal.pgen.1001131. PubMed PMID:
501 20885795; PubMed Central PMCID: PMCPMC2944798.
- 502 49. Kimura M. Attainment of Quasi Linkage Equilibrium When Gene Frequencies Are Changing by
503 Natural Selection. Genetics. 1965;52(5):875-90. PubMed PMID: 17248281; PubMed Central PMCID:
504 PMCPMC1210959.
- 505 50. Pedruzzi G, Barlukova A, Rouzine IM. Evolutionary footprint of epistasis. PLoS Comput Biol.
506 2018;14(9):e1006426. doi: 10.1371/journal.pcbi.1006426. PubMed PMID: 30222748; PubMed Central
507 PMCID: PMCPMC6177197.
- 508 51. Desai MM, Weissman D, Feldman MW. Evolution can favor antagonistic epistasis. Genetics.
509 2007;177(2):1001-10. Epub 2007/08/28. doi: 10.1534/genetics.107.075812. PubMed PMID: 17720923;
510 PubMed Central PMCID: PMCPMC2034608.
511

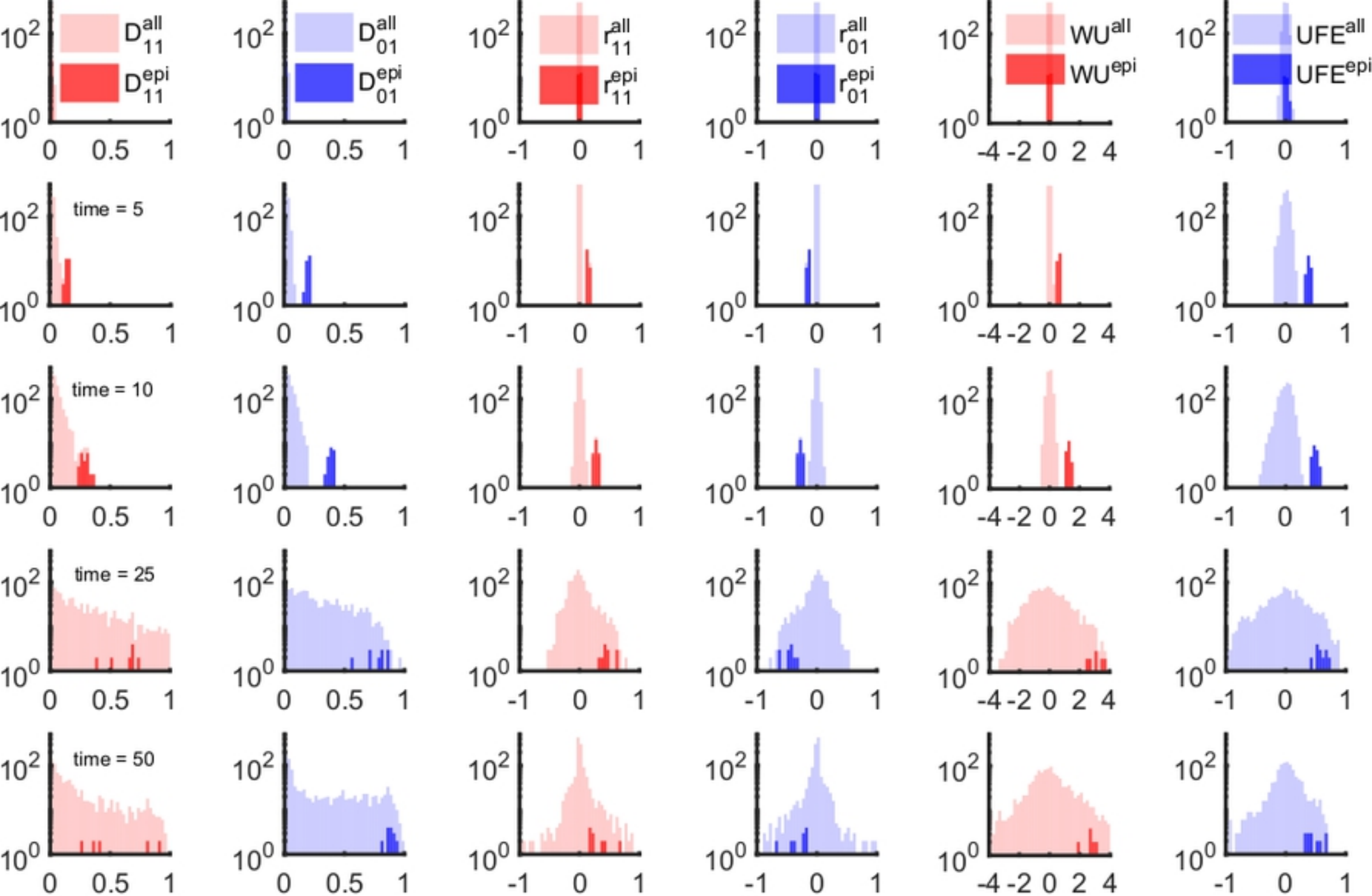


Figure 1

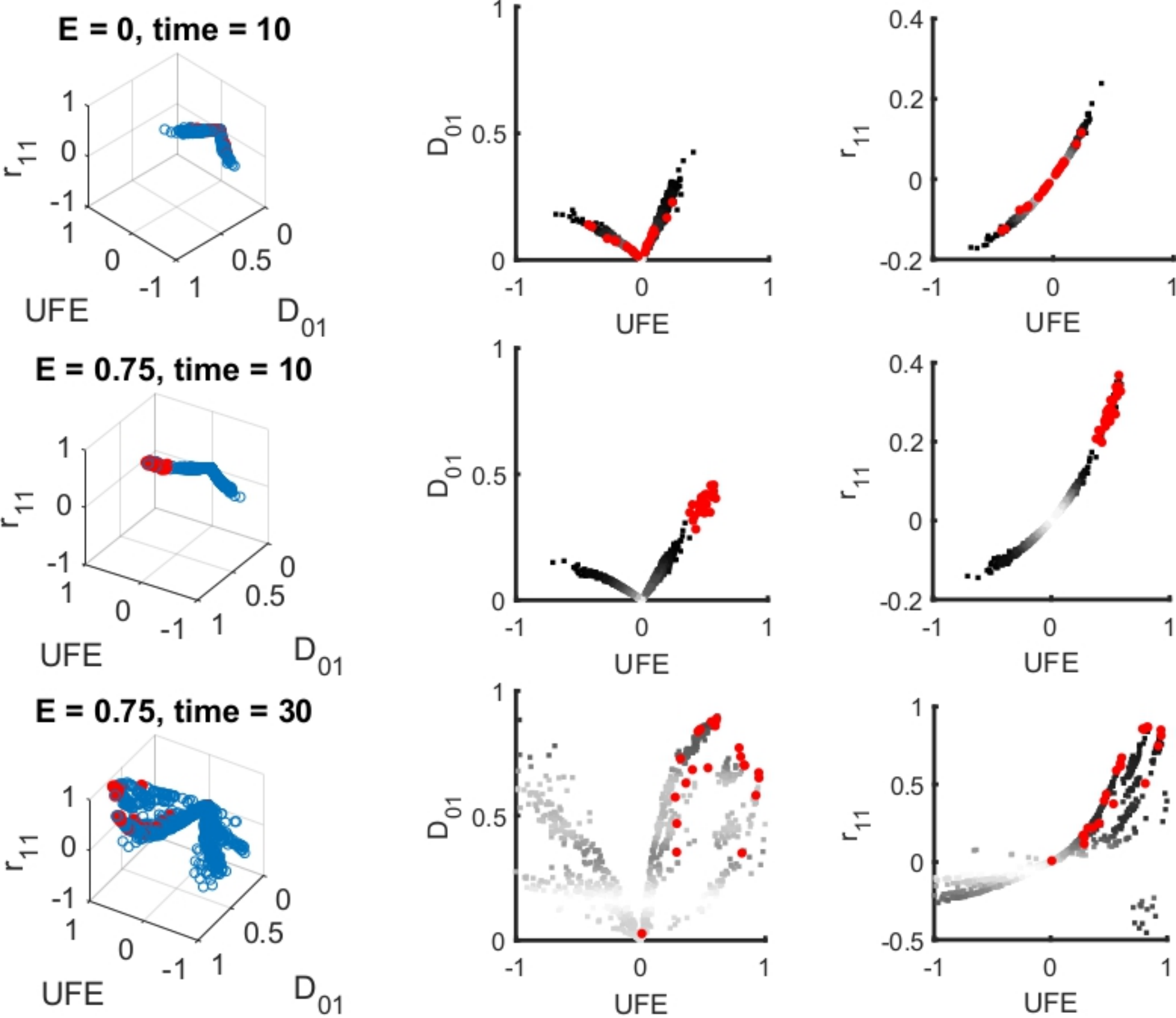


Figure 2

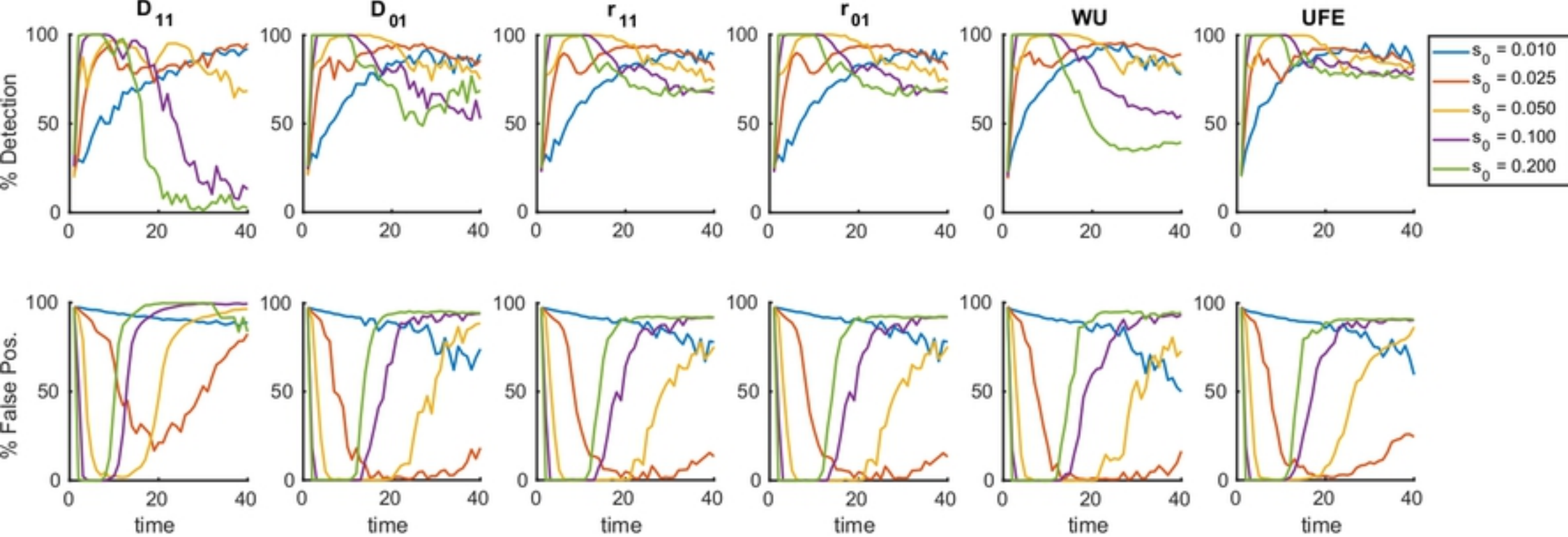


Figure 3

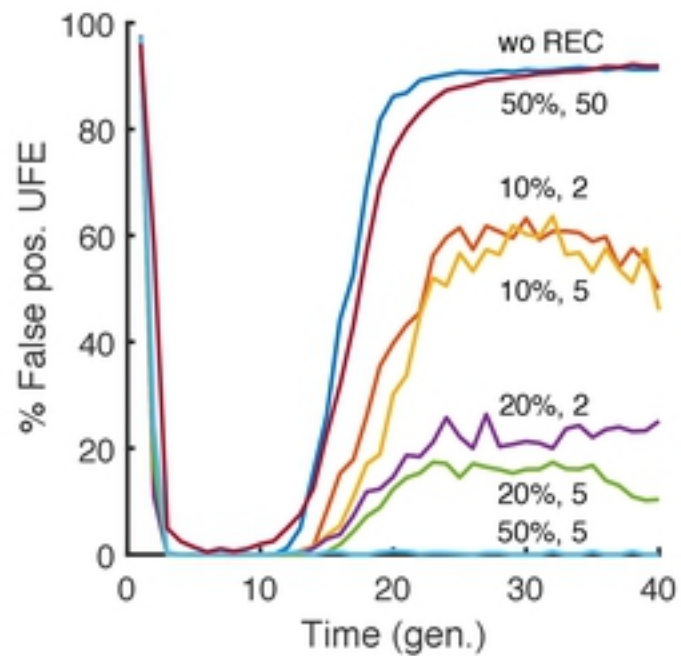
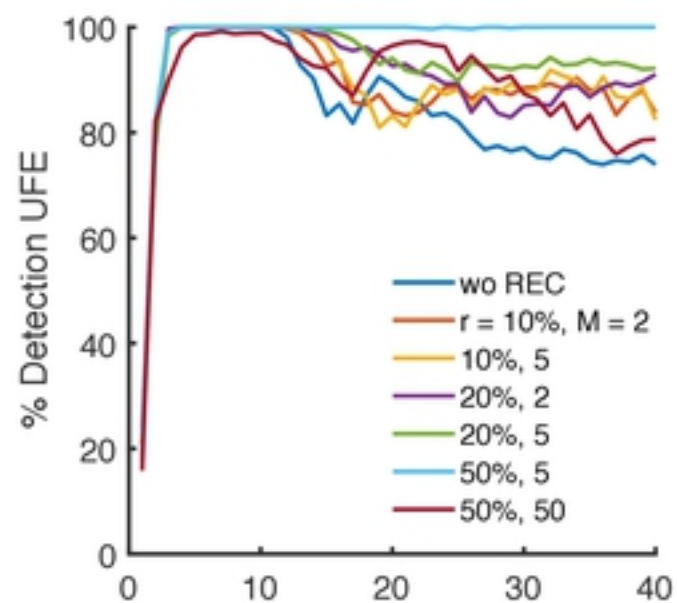


Figure 4

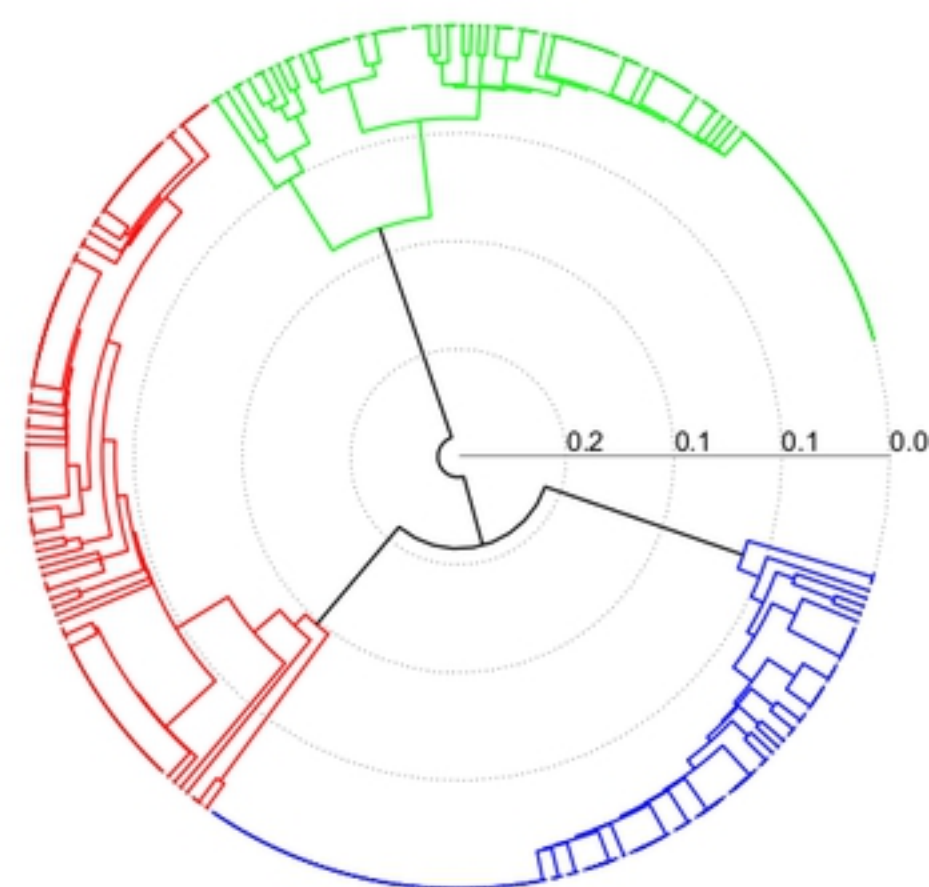
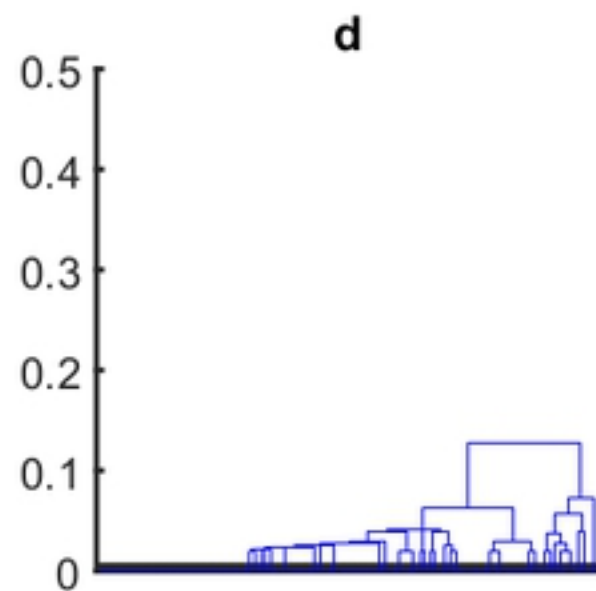
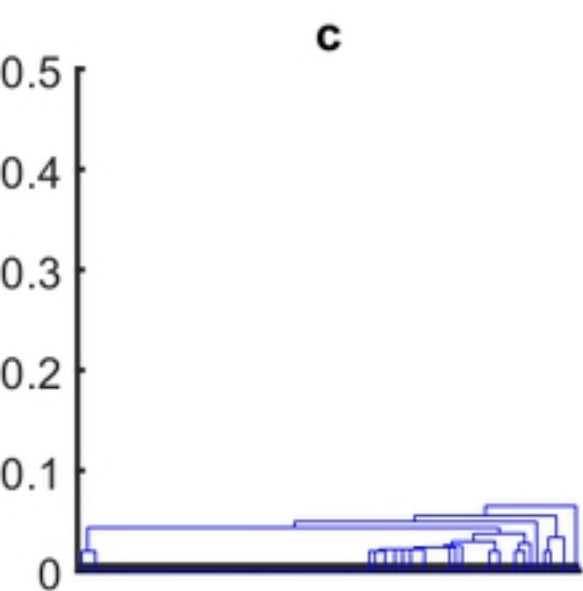
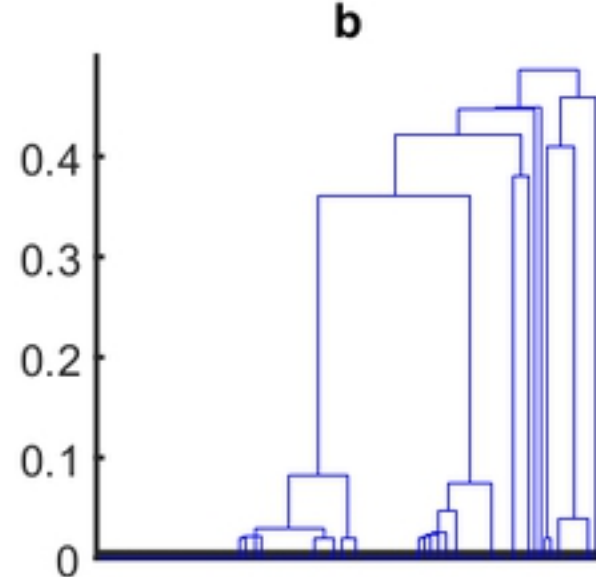
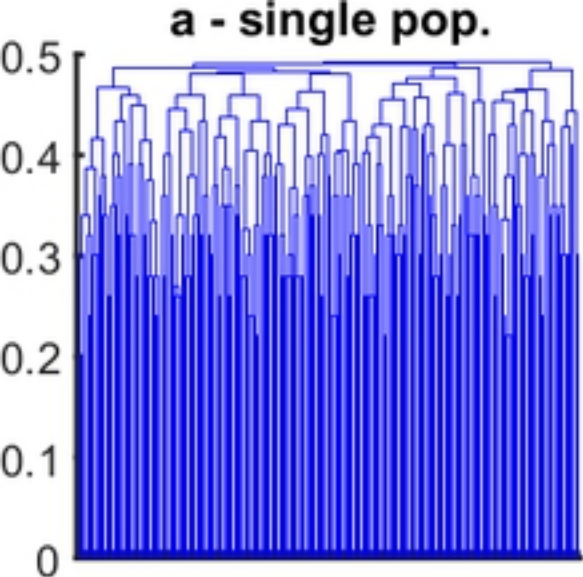


Figure 5

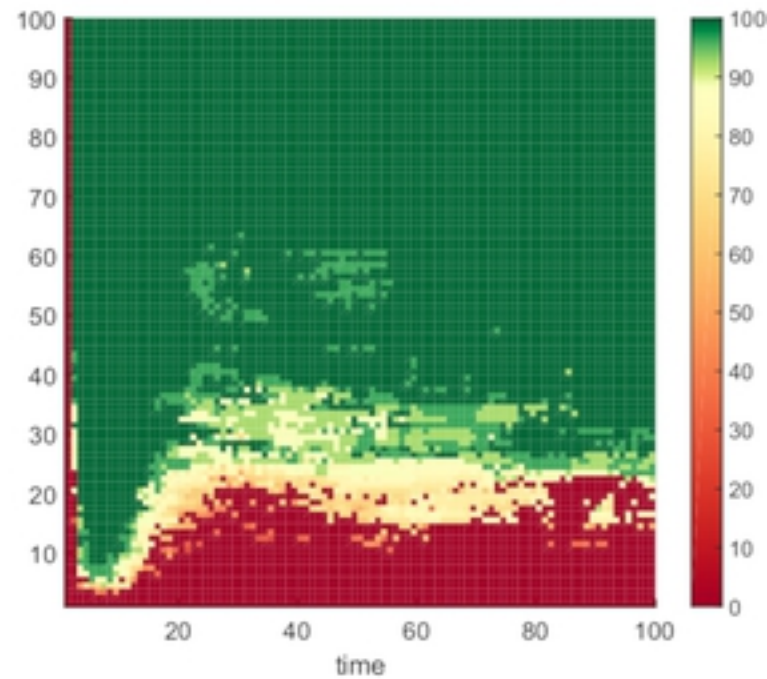
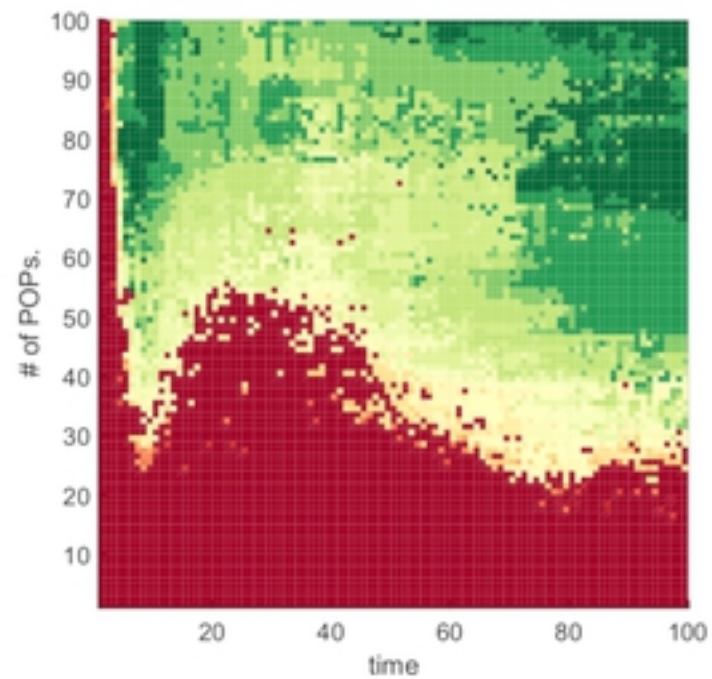
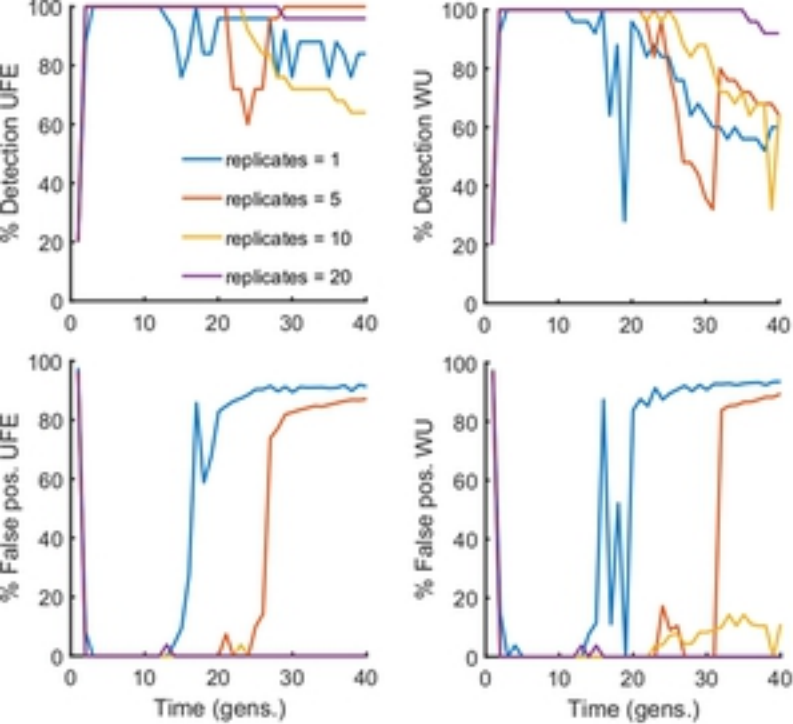


Figure 6