

CLEAR: Composition of Likelihoods for Evolve And Resequence Experiments

Arya Iranmehr¹, Ali Akbari¹, Christian Schlötterer², and Vineet Bafna³

¹Electrical and Computer Engineering, University of California, San Diego, La Jolla, CA, USA.

²Institut für Populationsgenetik, Vetmeduni, Vienna, Austria.

³Computer Science and Engineering, University of California, San Diego, La Jolla, CA, USA.

Abstract

Experimental evolution (EE) studies are powerful tools for observing molecular evolution “in-action” from populations sampled in controlled and natural environments. The advent of next generation sequencing technologies has made whole-genome and whole-population sampling possible, even for eukaryotic organisms with large genomes, and allowed us to locate the genes and variants responsible for genetic adaptation. While many computational tests have been developed for detecting regions under selection, they are mainly designed for static (single time) data, and work best when the favored allele is close to fixation. Conversely, EE studies provide samples over multiple time points, often at early stages of selective sweep.

While more predictive than static data analysis, a majority of the EE studies are constrained by the limited time span since onset of selection, depending upon the generation time of the organism. This constraint curbs the power of adaptation studies, as the population can only be evolved-and-resequenced for a small number of generations relative to the fixation-time of the favored allele. Moreover, coverage in pooled sequencing experiments varies across replicates and time points for every variant.

In this article, we directly address these issues while developing tools for identifying selective sweep in pool-sequenced EE of sexual organisms and propose Composition of Likelihoods for Evolve-And-Resequence experiments (CLEAR). Extensive simulations show that CLEAR achieves higher power in detecting and localizing selection over a wide range of parameters. In contrast to existing methods, the CLEAR statistic is robust to variation of coverage. CLEAR also provides robust estimates of model parameters, including selection strength and overdominance, as byproduct of the statistical testing, while being orders of magnitude faster. Finally, we applied the CLEAR statistic to data from a study of *D. melanogaster* adaptation to alternating temperatures. We identified selection in many genes, including Heat Shock Proteins. The genes were enriched in “response to heat”, “cold acclimation” and “defense response to bacterium”, and other relevant biological processes.

1 Introduction

Natural selection is a key force in evolution, and a mechanism by which populations can adapt to external ‘selection’ constraints. Examples of adaptation abound in the natural world [26], including for example, classic examples like lactose tolerance in Northern Europeans [11], human adaptation to high altitudes [70, 86], but also drug resistance in pests [18], HIV [29], cancer [34, 87], malarial parasite [5, 55], and other antibiotic resistance [71]. In these examples, understanding the genetic basis of adaptation can provide actionable information, underscoring the importance of the problem.

34 Experimental evolution refers to the study of the evolutionary processes of a model organism
35 in a controlled [9, 12, 36, 46, 47, 58, 59] or natural [7, 10, 19, 20, 51, 64, 85] environment. Recent
36 advances in whole genome sequencing have enabled us to sequence populations at a reasonable cost
37 even for large genomes. Perhaps more important for experimental evolution studies, we can now
38 evolve and re-sequence multiple replicates of a population to obtain *longitudinal time-series data*, in
39 order to investigate the dynamics of evolution at molecular level. Although constraints such as small
40 sizes, limited timescales, and oversimplified laboratory environments may limit the interpretation of
41 EE results, these studies are increasingly being used to test a wide range of hypotheses [43] and have
42 been shown to be more predictive than static data analysis [14, 21, 68]. In particular, longitudinal
43 EE data is being used to estimate model parameters including population size [42, 61, 75, 80, 81, 83],
44 strength of selection [13, 38, 39, 50, 53, 72, 75], allele age [50] recombination rate [75], mutation
45 rate [8, 75], quantitative trait loci [6] and for tests of neutrality hypotheses [10, 15, 28, 75].

46 While objectives, designs and organisms of EE studies can be entirely different [8, 69], here we
47 restrict our attention to the adaptive evolution of multi-cellular sexual organisms. For simplicity,
48 we assume fixed population size, and for the most part, positive single locus selection (only one
49 favored mutation). This regime has been considered earlier, typically with *D. melanogaster* as the
50 model organism of choice, to identify adaptive genes in longevity and aging [15, 65] (600 genera-
51 tions), courtship song [78] (100 generations), hypoxia tolerance [88] (200 generations), adaptation
52 to new laboratory environments [31, 58] (59 generations), egg size [41] (40 generations), C virus
53 resistance [52] (20 generations), and dark-fly [40] (49 generations).

54 The task of identifying genetic adaptation can be addressed at different levels of specificity. At
55 the coarsest level, identification could simply refer to deciding whether some genomic region (or a
56 gene) is under selection or not. In the following, we refer to this task as *detection*. In contrast,
57 the task of *site-identification* corresponds to the process of finding the favored mutation/allele
58 at nucleotide level. Finally, *estimation of model parameters*, such as strength of selection and
59 overdominance at the site, can provide a comprehensive description of the selection process.

60 A wide range of computational methods [79] have been developed to detect regions under
61 positive selection. A majority of the existing methods focus on static data analysis; analysis of a
62 single sample of the population at a specific time, either during the sweep, or subsequent to fixation
63 of the favored allele. Static analysis is focused on reduction in genetic diversity [27, 32, 66, 74] shift
64 in allele-frequencies, prevalence of long haplotypes [67, 79], population differentiation [15, 35, 37]
65 in multiple-population data and others. Many existing methods use the Site Frequency Spectrum
66 (SFS, see S1 Fig) to identify departure from neutrality. Classical examples including Tajima's
67 *D* [74], Fay and Wu's *H* [27], Composite Likelihood Ratio [57], were all shown to be weighted
68 linear combinations of the SFS values [1]. While successful, these methods are prone to both, false
69 negatives [54], and also false-discoveries due to confounding factors such as demography, including
70 bottleneck and population expansions, and ascertainment bias [3, 54, 56, 62, 63]. Nevertheless,
71 SFS based tests continue to be used successfully, often in combination with other tests [3, 79]. One
72 of the contributions of this paper is the extension of SFS based methods to analyze time-series
73 data, and the identification of selection regimes where these methods perform well.

74 Relative to the analysis of static samples, fewer tests-of-selection for dynamic time-series data
75 have been proposed. Often, existing tests for static data are adopted for dynamic data with two
76 time-points. Zhu *et al.* [88] used the ratio of the estimated population size of case and control
77 populations to compute test statistic for each genomic region. Burke *et al.* [15] applied Fisher
78 exact test to the last observation of data on case and control populations. Orozco-terWengel
79 *et al.* [58] used the Cochran-Mantel-Haenszel (CMH) test [2] to detect SNPs whose read counts
80 change consistently across all replicates of two time-point data. Turner *et al.* [78] proposed the
81 diffStat statistic to test whether the change in allele frequencies of two populations deviate from the

82 distribution of change in allele frequencies of two drifting populations. Bergland *et al.* [10] applied
83 F_{st} to populations throughout time to signify their differentiation from ancestral (two time-point
84 data) as well as geographically different populations. Jha *et al.* [41] computed test statistic of
85 generalized linear-mixed model directly from read counts. Bollback *et al.* [13] provided a diffusion
86 approximation to the continues Wright Fisher Markov process and estimated s numerically and
87 provided a standard likelihood ratio test under χ^2 distribution.

88 It is only recently that direct tests for analyzing time-series data have been developed, and
89 much of it is based on whole-genome sequencing of pools of individuals (pool-seq) at specific times.
90 Using continuous-time continuous-state Brownian motion process, Feder *et al.* [28] proposed the
91 Frequency Increment Test (FIT). More recently, Topa *et al.* [77] proposed a Gaussian Process (GP)
92 for modeling single-locus time-series pool-seq data. Terhorst *et al.* [75] extended GP to compute
93 joint likelihood of multiple loci under null and alternative hypotheses.

94 A key contribution of our paper is the development of a direct, and significantly faster method,
95 CLEAR, for identifying selection in short-term experimental evolution with pool-seq data. We show
96 for a wide range of parameters that CLEAR provides higher power for detecting selection, is robust
97 to ascertainment bias due to coverage heterogeneity, estimates model parameters consistently, and
98 localizes favored allele more accurately compared to the state-of-the-art methods, while being orders
99 of magnitude faster.

100 2 Materials and Methods

101 **Statistical Procedure.** To find the genes that are responding to the selection pressure, we con-
102 sider a likelihood-based approach [13, 57, 75, 77, 79]. The statistical procedure involves

- 103 (i) **Estimating population size.** The procedure starts by finding maximum likelihood estimate
104 of population size, \hat{N} , over the whole genome.
- 105 (ii) **Estimating selection parameters.** Given \hat{N} , maximizing the likelihood of the time series
106 data w.r.t. selection and overdominance parameters s, h , for each polymorphism.
- 107 (iii) **Computing likelihood statistics.** For each variant, it calculates in order to calculate the
108 log-odds ratio of the likelihood of selection model to the likelihood of neutral evolution/drift
109 model. Likelihood ratios in a genomic region are combined to compute the CLEAR statistic,
110 which is a composite likelihood score for the region being under selection. Extremal analysis
111 will provide candidate regions for further functional analysis.
- 112 (iv) **Hypothesis testing.** The null distribution of likelihood ratio statistics are computed on a
113 set of single locus drift simulations with population size of \hat{N} , and configuration (e.g. sequenc-
114 ing coverage) as closely as possible to the EE in which real data generated. Given the null
115 distribution of statistics, p -value of the observed likelihood ratios are calculated. False Dis-
116 covery Rate (FDR) correction is applied to the p -values to account for multiple testing. The
117 overlapping genes with the variants that satisfy FDR criterion, will be reported for functional
118 analysis or imported to the Gene Set Enrichment Analysis (GSEA).
- 119 (v) **Extra analysis.** of estimated parameter could reveal extra information, e.g., overdominance
120 or fixation time, regarding the ongoing selective sweep.

121 In the subsequent of this section, we outline different steps of the statistical procedure.

122 **2.1 Estimating Population Size**

123 Estimating population sizes from temporal neutral evolution data has been previously studied [4,
 124 13, 75, 83]. Existing methods are well designed for when the allele frequencies are computed from
 125 a finite sample, that is the ascertainment bias is uniform over the genome. However, in the case of
 126 pool-seq data in addition to uniform ascertainment bias, each variant is sampled at different rate,
 127 due to finite sequencing coverage. In addition, majority of the state-of-the-art models [13, 28, 75, 77]
 128 model time and state space as continuous variables. However, in our experiments, we found that
 129 the smooth approximations is inadequate for small populations, low starting frequencies and sparse
 130 sampling (in time) that are typical in experimental evolution (see Results, Fig 7A-C, and Fig 4).
 131 To this end, we use using a discrete-time discrete-state-space Wright-Fisher Markov chain to model
 132 dynamic pool-seq data. In order to find an estimate of population size, we first need to find
 133 likelihood of the population size given neutral pool-seq data.

134 **Likelihood for Neutral Model.** Consider a neutrally evolving diploid population with fixed size
 135 of N individuals where ν_t denotes allele frequency of the derived allele at generation t . Experimental
 136 evolution for R replicates is conducted so that at generations $\mathcal{T} = \{\tau_i : 0 \leq \tau_0 < \tau_1, \dots < \tau_T\}$, n
 137 individuals are chosen for sequencing.

138 At the highest level, the consecutive allele frequencies of the population in a fixed-size Wright-
 139 Fisher model evolves by Binomial sampling

$$\nu_0 \sim \pi, \quad 2N\nu_t | \nu_{t-1} \sim \text{Binomial}(2N, \nu_{t-1}) \quad (1)$$

140 where π is the marginal distribution of allele frequencies in the base population. In general π depends
 141 on demographic history of the founder lines, here we simply assume π is the site frequency
 142 spectrum of fixed sized neutral population S1 Fig.

143 To compute distributions after τ transitions, it is enough to specify the $2N \times 2N$ transition
 144 matrix $P^{(\tau)}$, where $P^{(\tau)}[i, j]$ denotes probability of change in allele frequency from $\frac{i}{2N}$ to $\frac{j}{2N}$ in τ
 145 generations:

$$P^{(1)}[i, j] = \Pr \left(\nu_{t+1} = \frac{j}{2N} \middle| \nu_t = \frac{i}{2N} \right) = \binom{2N}{j} \nu_t^j (1 - \nu_t)^{2N-j}, \quad (2)$$

$$P^{(\tau)} = P^{(\tau-1)} P^{(1)} \quad (3)$$

146 As at each generation n out of N individuals are randomly selected for sequencing, The *sampled*
 147 *allele frequencies*, $\{y_t\}_{t \in \mathcal{T}}$, are also Binomially distributed

$$2ny_t \sim \text{Binomial}(2n, \nu_t) \quad (4)$$

148 We introduce the $2N \times 2n$ sampling matrix Y , where $Y[i, j]$ stores the probability that the sample
 149 allele frequency is $\frac{i}{2n}$ given that the true allele frequency is $\frac{i}{2N}$.

150 We denote the pool-seq data for that variant as $\{x_t = \langle c_t, d_t \rangle\}_{t \in \mathcal{T}}$ where d_t, c_t represent the
 151 read depth, and the read count of the derived allele, respectively, at time τ_t . Let $\{\lambda_t\}_{t \in \mathcal{T}}$ be the
 152 sequencing coverage at different generations, then, the observed data are sampled according to

$$d_t \sim \text{Poiss}(\lambda_t), \quad c_t \sim \text{Binomial}(d_t, y_t) \quad (5)$$

153 where the emission probabilities for a observed tuple $x_t = \langle d_t, c_t \rangle$ is

$$\mathbf{e}_i(x_t) = \binom{d_t}{c_t} \left(\frac{i}{2n} \right)^{c_t} \left(1 - \frac{i}{2n} \right)^{d_t - c_t}. \quad (6)$$

154

155 For $1 \leq t \leq T$, let $\alpha_{t,i}$ denote the probability of emitting x_1, x_2, \dots, x_t and reaching state i at
 156 τ_t . Then, $\alpha_{t,i}$ can be computed using the forward-procedure [22]:

$$\begin{aligned}\alpha_{t,i} &= \left(\sum_{1 \leq j \leq 2N} \alpha_{t-1,j} P^{(\delta_t)}[j, i] \right) Y \mathbf{e}_i(x_t) \\ \alpha_t &= \text{Diag}(\alpha_{t-1}) P^{(\delta_t)} Y \mathbf{e}(x_t)\end{aligned}\quad (7)$$

157 where $\delta_t = \tau_t - \tau_{t-1}$. The joint likelihood of the observed data from R independent observations is
 158 given by

$$\Pr(\{\mathbf{x}^{(r)}\} | N, n) = \mathcal{L}(N | \{\mathbf{x}^{(r)}\}, n) = \prod_{r=1}^R \mathcal{L}(N | \mathbf{x}^{(r)}, n) = \prod_{r=1}^R \sum_i \alpha_{T,i}^{(r)}. \quad (8)$$

159 where $\mathbf{x} = \{x_t\}_{t \in \mathcal{T}}$. The graphical model and the generative process for which data is being
 160 generated is depicted in Fig 1.

161 Finally, the last step is to find the \hat{N} in which maximizes the likelihood of the all variants in
 162 whole genome:

$$\hat{N} = \arg \max_N \prod_i^M \mathcal{L}(N | \{\mathbf{x}_i^{(r)}\}) \quad (9)$$

163 2.2 Estimating Selection Parameters

164 **Likelihood for Selection Model.** Assume that the site is evolving under selection constraints $s \in$
 165 \mathbb{R} , $h \in \mathbb{R}_+$, where s and h denote selection strength and overdominance parameters, respectively.
 166 By definition, the relative fitness values of genotypes $0|0$, $0|1$ and $1|1$ are given by $w_{00} = 1$,
 167 $w_{01} = 1 + hs$ and $w_{11} = 1 + s$. Recall that ν_t denotes the frequency of the site at time $\tau_t \in \mathcal{T}$.
 168 Then, ν_{t+} , the frequency at time $\tau_t + 1$ (one generation ahead), can be estimated using:

$$\begin{aligned}\hat{\nu}_{t+} &= \mathbb{E}[\nu_{t+} | s, h, \nu_t] = \frac{w_{11}\nu_t^2 + w_{01}\nu_t(1 - \nu_t)}{w_{11}\nu_t^2 + 2w_{01}\nu_t(1 - \nu_t) + w_{00}(1 - \nu_t)^2} \\ &= \nu_t + \frac{s(h + (1 - 2h)\nu_t)\nu_t(1 - \nu_t)}{1 + s\nu_t(2h + (1 - 2h)\nu_t)}.\end{aligned}\quad (10)$$

169 The machinery for computing likelihood of the selection parameters is identical to that of population
 170 size, except for transition matrices. Hence, here we only introduce the definition transition matrix
 171 $Q_{s,h}^{(\tau)}$ of the selection model. Let $Q_{s,h}^{(\tau)}[i, j]$ denote the probability of transition from $\frac{i}{2N}$ to $\frac{j}{2N}$ in τ
 172 generations, then (See [24], Pg. 24, Eqn. 1.58-1.59):

$$Q_{s,h}^{(1)}[i, j] = \Pr \left(\nu_{t+} = \frac{j}{2N} \middle| \nu_t = \frac{i}{2N}; s, h, N \right) = \binom{2N}{j} \hat{\nu}_{t+}^j (1 - \hat{\nu}_{t+})^{2N-j} \quad (11)$$

$$Q_{s,h}^{(\tau)} = Q_{s,h}^{(\tau-1)} Q_{s,h}^{(1)} \quad (12)$$

173 The maximum likelihood estimates are given by

$$\hat{s}, \hat{h} = \arg \max_{s, h} \prod_i^M \mathcal{L}(s, h | \{\mathbf{x}_i^{(r)}\}, \hat{N}) \quad (13)$$

174 The parameters in Eqs. 9, 13 are optimized using grid search. By broadcasting and vectorizing
 175 the grid search operations across all variants, the genome scan on millions of polymorphisms can
 176 be done in significantly smaller time than iterating a numerical optimization routine for each
 177 variant (see Results and Fig 8).

178 **2.3 Empirical Likelihood Ratio Statistics**

179 Similar to Eq. ??, let \hat{s}, \hat{h} denote the parameters that maximize likelihood. The corresponding
 180 likelihood ratio statistic for each variant of pool-seq data is given by

$$H = -2 \log \left(\frac{\mathcal{L}(\hat{s}, \hat{h} | \{\mathbf{x}^{(r)}\}, \hat{N})}{\mathcal{L}(0, 0.5 | \{\mathbf{x}^{(r)}\}, \hat{N})} \right) \quad (14)$$

181 **Composite Likelihood Ratio.** In general, the favored allele can be in linkage disequilibrium
 182 with some of its surrounding variation. The linked-loci hitchhike and share similar dynamics with
 183 the favored allele [23]. Some models such as multi-locus Gaussian process [75] incorporate these
 184 associations by modeling linkage and recombination explicitly. However, these approaches are com-
 185putationally expensive. Moreover, linkage computations are infeasible without haplotype resolved
 186 data, which pool-seq does not provide. Instead, we work with a simpler Composite Likelihood
 187 Ratio (CLR) [57, 84] computation to combine the individual scores of all variants into a composite
 188 score.

189 Consider a genomic region L to be a collection of segregating sites with little or no recombination
 190 between sites and the favored allele. This scenario holds when the starting frequency of the favored
 191 allele is not high and the region is small. Let H_ℓ denote the likelihood ratio score based for each
 192 site ℓ in L . The CLR is computed by averaging scores of all the variants withing the testing region

$$\mathcal{H} = \frac{1}{|L|} \sum_{\ell \in L} H_\ell. \quad (15)$$

193 **2.4 Hypothesis Testing**

194 ***p*-value Computation.** By Wilks theorem [82], the likelihood ratio statistic (Eq. 14) is asymp-
 195tically distributed according to χ_k^2 , where k is the difference between dimensions of parameters
 196 of alternative and null model. Feder *et al.* [28] showed that the empirical distribution of statistic
 197 on simulations with \hat{N} provide more accurate *p*-values than χ^2 when the number of independent
 198 samples (replicates) is small. However, the simulations for composite statistics require initial hap-
 199lotypes of the founder population, local recombination rates and massive computational resources.
 200 Hence we use χ^2 for testing composite statistic \mathcal{H} and simulations for polymorphism statistic.

201 Specifically, for composite statistics, we make a simplifying assumption that is all the statistics
 202 H_ℓ are independent, then

$$H_\ell \sim \chi_2^2 \implies \sum_{\ell \in L} H_\ell \sim \chi_{2L}^2 \iff \mathcal{H} \sim \text{Gamma}(\alpha = L, \theta = 2/L) \quad (16)$$

203 To create single locus simulations, we repeat the generative process for $N = \hat{N}$ outlined in
 204 the Fig 1 for every variant in the real data, with only difference that ν_0 and $\{\mathbf{d}^{(r)}\}$ are given from
 205 the real data.

206 **Correcting for multiple testing**

207 **Enrichment Analysis**

Table 1: Overdominance parameter values and their implications.

Value	Condition
$h = 0$	recessive adaptive allele
$h = 0.5$	directional selection
$h = 1$	dominant adaptive allele
$h > 1$	overdominance

208 2.5 Extra Analysis

209 **Site-identification.** Once a genomic region is classified to be under selection by \mathcal{H} or \mathcal{M} statistic,
210 individual variant scores (M or H) in a region are ranked to predict the favored site. In general,
211 identifying the favored site in pool-seq data is difficult [76], due to extensive span of hitchhikers in
212 an ongoing sweep (see S1.2 Text for more detail). In our analysis of the *D. melanogaster* EE data,
213 we identify a set of “candidate” variants whose scores exceed a False Discovery Rate threshold
214 based on the distribution of CLEAR scores on negative controls.

215 **Overdominance.** The value of the overdominance parameter can reveal if the favored allele is
216 overdominant, repressive or dominant [33] (See Table 1, S6 Fig and S7 Fig). The test statistic for
217 such hypothesis test is

$$D = -2 \log \left(\frac{\mathcal{L}(\bar{s}, \hat{h} | \{\mathbf{x}^{(r)}\}, \hat{N})}{\mathcal{L}(\bar{s}, 0.5 | \{\mathbf{x}^{(r)}\}, \hat{N})} \right), \text{ where } \bar{s} = \arg \max_s \prod_i^M \mathcal{L}(s, 0.5 | \{\mathbf{x}_i^{(r)}\}, \hat{N}). \quad (17)$$

218 **Precomputing Transition Matrices.** CLEAR requires a one-time computation of matrices $Q_{s,h}^{(\tau)}$
219 for the entire range of s, h values. Precomputation of 909 transition matrices for $s \in \{-0.5, -0.49, \dots, 0.5\}$
220 and $h \in \{0, 0.25, \dots, 2\}$ took less than 15 minutes (≈ 1 second per matrix) on a desktop computer
221 with a Core i7 CPU and 16GB of RAM.

222 2.6 Simulations

223 We performed extensive simulations using parameters that have been used for *D. melanogaster*
224 experimental evolution [45]. See also Fig 3 for illustration. To implement real world pool-seq
225 experimental evolution, we conducted simulations as follows:

226 I. **Creating initial founder line haplotypes.** Using `msms` [25], we created neutral popula-
227 tions for F founding haplotypes with command `./msms <F> 1 -t <2μLNe> -r <2rNeL>`
228 `<L>`, where $F = 200$ is number of founder lines, $N_e = 10^6$ is effective population size,
229 $r = 2 \times 10^{-8}$ is recombination rate, $\mu = 2 \times 10^{-9}$ is mutation rate and $L = 50K$ is the
230 window size in base pairs which gives $\theta = 2\mu N_e L = 200$ and $\rho = 2N_e r L = 2000$.

231 II. **Creating initial diploid population.** To simulate experimental evolution of diploid organ-
232 isms, initial haplotypes were first cloned to create F diploid homozygotes. Next, each diploid
233 individual was cloned N/F times to yield diploid population of size N .

234 III. **Forward Simulation.** We used forward simulations for evolving populations under selection.
235 We also consider selection regimes which the favored allele is chosen from standing variation
236 (not *de novo* mutations). Given initial diploid population, position of the site under selection,
237 selection strength s , number of replicates $R = 3$, recombination rate $r = 2 \times 10^{-8}$ and

238 sampling times $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$, simuPop [60] was used to perform forward simulation
239 and compute allele frequencies for all of the R replicates. For hard sweep (respectively, soft
240 sweep) simulations we randomly chose a site with initial frequency of $\nu_0 = 0.005$ (respectively,
241 $\nu_0 = 0.1$) to be the favored allele.

242 **IV. Sequencing Simulation.** Give allele frequency trajectories we sampled depth of each site
243 identically and independently from $\text{Poisson}(\lambda)$, where $\lambda \in \{30, 100, 300, \infty\}$ is the coverage
244 for the experiment. Once depth d is drawn for the site with frequency ν , the number of reads
245 c carrying the derived allele are sampled according to $\text{Binomial}(d, \nu)$. For experiments with
246 finite depth the tuple (c, d) is the input data for each site. Infinite depth experiments refer
247 to the case, where the true allele frequency is provided and Markov and HMM likelihood
248 computations give identical results.

249 3 Results

250 **Modeling Allele Frequency Trajectories in Finite Populations.** We first tested the good-
251 ness of fit of the discrete-time Markov chain versus continuous-time Brownian motion (Gaussian
252 approximation) in modeling allele frequency trajectories in finite populations, under different sam-
253 pling schemes and starting frequencies. For this purpose, we conducted 100K simulations with
254 two time samples $\mathcal{T} = \{0, \tau\}$ where $\tau \in \{1, 10, 100\}$ is the parameter controlling the density of
255 sampling in time. In addition, we repeated simulations for different values of starting frequency
256 $\nu_0 \in \{0.005, 0.1\}$ (i.e., hard and soft sweep) and selection strength $s \in \{0, 0.1\}$ (i.e., neutral and
257 selection). Then, given initial frequency ν_0 , we computed the expected distribution of the frequency
258 of the next sample ν_τ under two models and compared them with empirical distributions calculated
259 from simulated data. Fig 4A-F shows that Brownian motion is inadequate when ν_0 is far from 0.5,
260 or when sampling times are sparse ($\tau > 1$). If the favored allele arises from standing variation in
261 a neutral population, it is unlikely to have frequency close to 0.5, and the starting frequencies are
262 usually much smaller (see S1 Fig). Moreover, in typical *D. melanogaster* experiments for example,
263 sampling is sparse. Often, the experiment is designed so that $10 \leq \tau \leq 100$ [31, 45, 58, 88].

264 In contrast to the Brownian motion results, Markov chain can provide predictions when the the
265 allele is under selection. In addition Fig 4A-M also shows that Markov chain predictions (Eq. 12)
266 are highly consistent with empirical data for a wide range of simulation parameters.

267 **Detection Power.** We compared the performance of CLEAR against other methods for detect-
268 ing selection. For each method we calculated detection power as the percentage of true-positives
269 identified with false-positive rate ≤ 0.05 (S9 Fig). For each configuration (specified with values for
270 selection coefficient s , starting allele frequency ν_0 and coverage λ), power of each method is eval-
271 uated over 2000 distinct simulations, half of which modeled neutral evolution and the rest modeled
272 positive selection.

273 Before comparing against other methods, we first evaluated the use of CLEAR with different
274 percentile-cutoffs π (Eq. 15) in computing composite statistics of a region. For each configuration,
275 we computed average Power for $s \in \{0.025, 0.05, 0.075, 0.1\}$, using $\mathcal{H}_\pi, \mathcal{H}_\pi^+$. We computed the
276 optimal value of π using a line-search. Fig 5 reveals several important trade-off between π , initial
277 frequency, and coverage.

- 278 • \mathcal{H}_π^+ consistently achieves a high power for $\pi = 0$, and in the absence of knowledge of the
279 selection regime or the ancestral allele, \mathcal{H}_0^+ is a powerful statistic to use.

- In every scenario tested, $\max_{\pi}\{\text{Power}(\mathcal{H}_{\pi})\} \geq \max_{\pi}\{\text{Power}(\mathcal{H}_{\pi}^+)\}$, suggesting that it is beneficial to make predictions based on \mathcal{H}_{π} , if the selection regime is well-understood and ancestral allele is known.
- In soft sweep, relative to hard sweep, it helps to choose a higher value of the cut-off π . This is consistent with the fact that LD between the favored site and other sites is generally lower for soft sweep. For instance, in soft sweep with infinite coverage (Fig 5F), optimum is gained at $\pi = 100$, equivalent to considering the score of the highest scoring site as \mathcal{H} statistic of the region.
- When coverage is low (Fig 5A,D), it helps to accumulate evidence from multiple sites, and the best results are achieved for lower values of π .

Finally, we compared the power of CLEAR with Gaussian process (GP) [75], FIT [28], and CMH [2] statistics. As CMH only takes read count data, here we used $\lambda = 300$ to implement infinite coverage scenario. All methods other than CLEAR and CMH convert read counts to allele frequencies prior to computing the test statistic. CLEAR shows the highest power in all cases and the power stays relatively high even for low coverage (Fig 7 and S1 Table). In particular, the difference in performance of CLEAR with other methods is pronounced when starting frequency is low. Starting frequency is at its minimum in selection on an allele from *de novo* mutations and is likely to be low if selection is on an allele from standing variation (S1 Fig). The advantage of CLEAR stems from the fact that favored allele with low starting frequency might be missed by low coverage sequencing. In this case, incorporating the signal from linked sites becomes increasingly important. We note that methods using only two time points, such as CMH, do relatively well for high selection values and high coverage. However, the use of time-series data can increase detection power in low coverage experiments or when starting frequency is low. Moreover, time-series data provide means for estimating selection parameters s, h (see below). Finally, as CLEAR is robust to change of coverage, our results (Fig 7B,C) suggest that taking many samples with lower coverage is preferable to sparse sampling with higher coverage.

SFS for Detection in Natural Samples. We did not show the SFS based statistics in Fig 7 as they did not perform better than random. In majority of controlled experimental evolution studies, the population is restricted set of F founder lines, where $F \ll N_e$ (Fig 3B) and inbred during the experiment. This creates a severe bottleneck, confounding SFS. S10 Fig demonstrates the effect of experimental evolution on different SFS statistics under neutral evolution for 1000 simulations. A second problem with using SFS for experimental evolution is that the sampling starts right after the onset of experimentally induced selection, and the favored allele may not reach high enough frequency to modify the site frequency spectrum (Fig 2).

However, in experiments involving naturally evolving populations, even if the span of the time-series is small, the onset of selection might occur many generations prior to sampling. To test performance of SFS-based statistics in natural evolution, using `mssms`, we conducted 200 (100 neutral and 100 sweep) forward simulations for different values of s , $N_e = 10K$ and $N = 200$. The start of sampling was chosen randomly after onset of selection in two distinct scenarios. Let $t_{\nu=x}(s, N_e)$ denote the expected time (in generations) required to reach carrier frequency x in a hard sweep and $U[a, b]$ denote discrete uniform distribution in the interval $[a, b]$. First we considered the case when start of sampling is chosen throughout the whole sweep. i.e., $\tau_0 \sim U[1, t_{\nu=1}(s, N_e)]$ (Fig 9A). Next, we considered sampling start time chosen nearer to fixation of the favored allele, i.e., $\tau_0 \sim U[t_{\nu=0.9}(s, N_e), t_{\nu=1}(s, N_e)]$ (Fig 9B). In both scenarios, sampling was done over 6 time points within 50 generations of τ_0 (Fig 3A). We compared CLEAR, GP, FIT with both static and dynamic SFS based statistics of SFSelect and Tajima's D. Fig 9A shows that SFS based statistics are outperformed by other methods. However, when sampling is performed close to fixation, i.e.,

327 when the favored allele has frequency of 0.9 or higher, SFS based statistics perform considerably
328 better than GP, FIT and CLEAR (Fig 9B). Moreover, dynamic SFS statistics provide higher power
329 than static SFS statistics, demonstrating that in the use of time-series SFS based statistics is
330 advantageous.

331 **Site-identification.** In general, localizing the favored variant, using pool-seq data is a nontrivial
332 task [76]. We used the simple approach of ranking each site in a region detected as being under
333 selection. The sites were ranked according to the likelihood ratio scores (Eqns. ??, ??). For each
334 setting of ν_0 and s , we conducted 1000 simulations and computed the rank of the favored mutation
335 in each simulation. The cumulative distribution of the rank of the favored allele in 1000 simulation
336 for each setting (Fig 10) shows that CLEAR outperforms other statistics. We also compared each
337 method to see how often it ranked the favored site in as the top ranked site (Table 2A-B), among
338 the top 10 ranked sites (Table 2C-D), and among the top 50 (Table 2E-F) ranked sites. In the
339 ≈ 1150 variants tested, CLEAR performed consistently better than other methods in all of these
340 measures.

341 An interesting observation is revisiting the contrast between site-identification and detection [48,
342 76] (see S1.2 Text for more details). When selection coefficient is high, detection is easier (Fig 7A-
343 F), but site-identification is harder due to the high LD between hitchhiking sites and the favored
344 allele (Table 2A-F). Moreover, site-identification is harder in hard sweep scenarios relative to soft
345 sweeps. For example, when coverage $\lambda = 100$ and selection coefficient $s = 0.1$, the detection power
346 is 75% for hard sweep, but 100% for soft sweep (Fig 7B-E). In contrast, the favored site was ranked
347 as the top in 14% of hard sweep cases, compared to and 95% of soft sweep simulations (Table 2A-B).

Table 2: Percentage of simulations which favored allele appears in top of the ranking.

Hard Sweep					Soft Sweep				
(A)					(B)				
s	CMH	FIT	GP	CLEAR	s	CMH	FIT	GP	CLEAR
0.025	3	0	0	2	0.025	9	1	4	11
0.05	5	0	0	5	0.05	53	10	51	63
0.075	11	0	0	10	0.075	87	31	92	91
0.1	15	0	0	14	0.1	93	54	95	95
(C)					(D)				
s	CMH	FIT	GP	CLEAR	s	CMH	FIT	GP	CLEAR
0.025	21	3	0	15	0.025	34	8	27	44
0.05	39	2	0	28	0.05	86	41	86	92
0.075	57	4	0	49	0.075	100	81	100	100
0.1	78	4	0	71	0.1	100	97	100	100
(E)					(F)				
s	CMH	FIT	GP	CLEAR	s	CMH	FIT	GP	CLEAR
0.025	51	19	0	43	0.025	60	22	62	75
0.05	70	18	0	66	0.05	96	69	97	99
0.075	85	20	2	81	0.075	100	97	100	100
0.1	94	26	6	93	0.1	100	100	100	100

Percentage of simulations in which the favored allele is ranked first (A-B); appears in top 10 (C-D); or, appears in top 50 (E-F). In soft sweep simulations (B,D,F), the ranks are consistently better than hard sweep simulations (A,C,E). This can be attributed to lower LD between the hitchhikers (false positives) and favored allele in soft sweep scenarios.

348 Estimating Parameters. CLEAR computes the selection parameters \hat{s} and \hat{h} as a byproduct of
 349 the hypothesis testing. We computed bias of selection fitness ($s - \hat{s}$) and overdominance ($h - \hat{h}$)
 350 for of CLEAR and GP in each setting. The distribution of the error (bias) for $100\times$ coverage is
 351 presented in Fig 11 for different configurations. S11 Fig and S12 Fig provide the distribution of
 352 estimation errors for $30\times$, and infinite coverage, respectively. For hard sweep, CLEAR provides
 353 estimates of s with lower variance of bias (Fig 11A). In soft sweep, GP and CLEAR both provide
 354 unbiased estimates with low variance (Fig 11B). Fig 11C-D shows that CLEAR provides unbiased
 355 estimates of h as well.

356 Running Time. As CLEAR does not compute exact likelihood of a region (i.e., does not explicitly
 357 model linkage between sites), the complexity of scanning a genome is linear in number of polymor-
 358 phisms. Calculating score of each variant requires $\mathcal{O}(TR)$ and $\mathcal{O}(TRN^2)$ computation for \mathcal{M} , and
 359 \mathcal{H} , respectively. However, most of the operations are can be vectorized for all replicates to make
 360 the effective running time for each variant. We conducted 1000 simulations and measured running
 361 times for computing site statistics M , H , FIT, CMH and GP with different number of linked-loci.
 362 Our analysis reveals (Fig 8) that CLEAR is orders of magnitude faster than GP, and comparable
 363 to FIT. While slower than CMH on the time per variant, the actual running times are comparable
 364 after vectorization and broadcasting over variants (see below).

365 These times can have a practical consequence. For instance, to run GP in the single locus
 366 mode on the entire pool-seq data of the *D. melanogaster* genome from a small sample ($\approx 1.6M$
 367 variant sites), it would take 1444 CPU-hours (≈ 1 CPU-month). In contrast, after vectorizing and

368 broadcasting operations for all variants operations using `numba` package, CLEAR took 75 minutes
369 to perform an scan, including precomputation, while the fastest method, CMH, took 17 minutes.

370 3.1 Analysis of a *D. melanogaster* Adaptation to Cold and Hot Temperatures

371 We applied CLEAR to the data from a study of *D. melanogaster* adaptation to alternating temper-
372 atures [31, 58], where 3 replicate samples were chosen from a population of *D. melanogaster* for
373 59 generations under alternating 12-hour cycles of hot (28°C) and cold (18°C) temperatures and
374 sequenced. In this dataset, sequencing coverage is different across replicates and generations (see
375 S2 Fig of [75]) which makes variant depths highly heterogeneous (S5 Fig and S4 Fig). We computed
376 the \mathcal{H}^+ statistic for sliding windows of 30Kbp with steps of 10Kbp over the whole genome. After
377 filtering out heterochromatic, centromeric and telomeric regions[17, 30], and applying a local false
378 discovery rate ≤ 0.01 (Methods), we identified 89 intervals (Fig 12) containing 968 genes (S6 Table).

379 We found 11 GO Biological Process terms to be enriched with Fisher exact P -value $\leq 10^{-3}$
380 (Table 3). The selected genes include many heat shock proteins in enriched ontologies including
381 ‘cold acclimation’ and ‘response to heat’ (S5 Table). As longer genes contain more variants, the
382 probability of a false variant being selected could increase with the length of the gene. Although the
383 CLEAR statistic for genes does not favor longer genes, we also performed a single variant based GO
384 enrichment using Gowinda [44]. The analysis identified 34 enriched GO terms (S4 Table) associated
385 with Biological Process. 5 of 11 GO terms in the gene level analysis were also among the 34 Gowinda
386 terms (Fisher exact p -val: 10^{-8}) suggesting consistency between variant and gene-based analysis.

Table 3: **GO (Biological Process) enrichment.**

Rank	GO ID	GO Term	$-\log(p\text{-value})$	Hits	Num of Genes
1	GO:0042742	defense response to bacterium	5.0	15	62
2	GO:0009408	response to heat	4.8	16	71
3	GO:0006719	juvenile hormone catabolic process	4.5	3	4
4	GO:0008363	larval chitin-based cuticle development	4.5	6	14
5	GO:0045664	regulation of neuron differentiation	4.5	3	4
6	GO:0051291	protein heterooligomerization	4.5	3	4
7	GO:0061077	chaperone-mediated protein folding	4.4	4	7
8	GO:0009631	cold acclimation	4.0	4	8
9	GO:0030837	negative regulation of actin filament polymerization	3.4	3	6
10	GO:0042026	protein refolding	3.1	4	11
11	GO:0007552	metamorphosis	3.0	8	36

A Fisher exact test was performed for GO enrichment in genes located in selected regions. All GO terms that contained at least 3 selected genes, and had Fisher exact p -value $\leq 10^{-3}$, are listed above.

387 **4 Discussion**

388 We developed a computational tool, CLEAR, that can detect regions under selection experimental
389 evolution experiments of sexual populations. Using extensive simulations, we show that CLEAR
390 outperforms existing methods in detecting selection, locating the favored allele, and estimating
391 selection parameters. Importantly, we make design choices that make CLEAR very fast in practice,
392 facilitating genome-wide studies.

393 Many factors play a role in adaptation during experimental evolution studies. The statistics
394 used by CLEAR perform well because they account for many of these aspects. CLEAR is not
395 restricted to two-time points, but uses the complete time-series data. Because it uses an exact
396 model, CLEAR achieves robust predictions for all values of the initial frequency. It adjusts for
397 heterogeneous ascertainment bias in finite-depth pooled-seq data to avoid hard filtering variants.
398 It exploits presence of high linkage within a region to compute composite likelihood ratio statistic.
399 Finally, CLEAR uses s, h as model parameters in its likelihood calculation, and provides optimized
400 estimates of these parameters, which can provide extra information such as fixation time, and
401 dominance ([S7 Fig](#)).

402 In our simulations, we found that the power of detection can be severely affected by the sampling
403 schedule as well as initial frequency of the favored allele. In general, while EE studies are powerful,
404 they also pose some challenges that are not adequately considered by other tools. One serious
405 constraint is the sampling time span, the gap between the first and last sampled generations, which
406 depends upon the generation time of the organism. It can be very small relative to the time of
407 fixation of the favored allele. In *D. melanogaster* for example, 30-50 generations are typical [45],
408 although there are some notable exceptions [88]. Therefore, unless the selection coefficient is very
409 strong, the time series data will only capture a ‘partial sweep’. This limitation is more pronounced
410 in controlled experimental evolution, where the sampling often starts at the onset of selection. In
411 particular, in a hard sweep scenario, the initial frequency of the favored allele is low, and may not
412 reach detectable frequency in sequencing, given the sampling time span. Through exact (discrete-
413 time, discrete-frequency) modeling, CLEAR performs better than competing tools even when initial
414 frequency is low and sampling time span is limited.

415 However, even if it were possible to sample over a larger time-span, many methods, especially
416 the ones that compute full likelihoods, would simply not scale to allow computation of evolutionary
417 trajectories over a large time-span. In contrast, CLEAR precomputes the transition matrices, and
418 scales linearly with number of samples, irrespective of the time-span in which they were acquired.

419 Sequence coverage is a practical consideration that is often ignored by other tools. Low se-
420 quencing coverage can lead to incorrect frequency estimates, even for the favored allele, especially
421 when the initial frequency is low. CLEAR uses HMMs to explicitly model variation in sequence
422 coverage. Moreover, it computes the composite likelihood from multiple linked sites, reducing the
423 impact of coverage on any one site, and detects selection even when the favored site is not sampled
424 due to low sequencing depth.

425 In controlled experimental evolution experiments, populations are evolved and inbred. As this
426 scenario involves picking a small number of founders, the effective population size significantly
427 drops from the large number of wild type (e.g., for *D. melanogaster*, $N_e \approx 10^6$) to a small number
428 of founder lines ($F \approx 10^2$). This creates a severe population bottleneck. The bottleneck confounds
429 SFS-based statistics and makes it difficult to fit a model or test a hypothesis ([S10 Fig](#)). Hence,
430 statistical testing based on SFS statistic provides poor performance in controlled experiments where
431 the initial sampling time is close to the onset of selection. However, SFS-based methods perform
432 very well when sampling is started long after the onset of selection (e.g., sampling from natural
433 populations). The larger time gap from the onset of selection provides an opportunity for the site

434 frequency spectrum to shift away from neutrality.

435 The comparison of hard and soft sweep scenarios lead to interesting observations. First, when
436 LD is high in the selected region, as is often the case in a hard sweep, composition of scores
437 significantly improves power of detection. When LD is low, as in soft sweep scenarios, composition of
438 scores does not work as well. However, the favored allele is well established at the onset of selection,
439 and will grow faster compared to the hard sweep scenario under identical selection regimes. This
440 makes it possible to detect selection even in soft sweep scenarios. The situation is a little different
441 with respect to localizing the favored allele. In soft sweep scenarios, the favored allele is not in high
442 LD with nearby variants, and its frequency change is independent of them. Therefore, we obtain
443 better localization results in soft sweep scenarios.

444 There are many directions to improve the analyses presented here. In particular, we plan to
445 focus our attention on other organisms with more complex life cycles, experiments with variable
446 population size and longer sampling-time-spans. As evolve and resequencing experiments continue
447 to grow, deeper insights into adaptation will go hand in hand with improved computational analysis.

448 **Software and Data Availability.** The source code and running scripts for CLEAR are publicly
449 available at <https://github.com/bafnalab/clear>. *D. melanogaster* data originally published [31, 58].
450 The dataset of the *D. melanogaster* study, until generation 37, is obtained from Dryad digital repos-
451 itory (<http://datadryad.org>)under accession DOI: 10.5061/dryad.60k68. Generation 59 of the *D.*
452 *melanogaster* study is accessed from European Sequence Read Archive (<http://www.ebi.ac.uk/ena/>)
453 under the project accession number: PRJEB6340. The dataset containing experimental evolution of
454 Yeast populations [16] is downloaded from <http://wfitch.bio.uci.edu/> tdlong/PapersRawData/BurkeYeast.gz
455 (last accessed 01/24/2017).

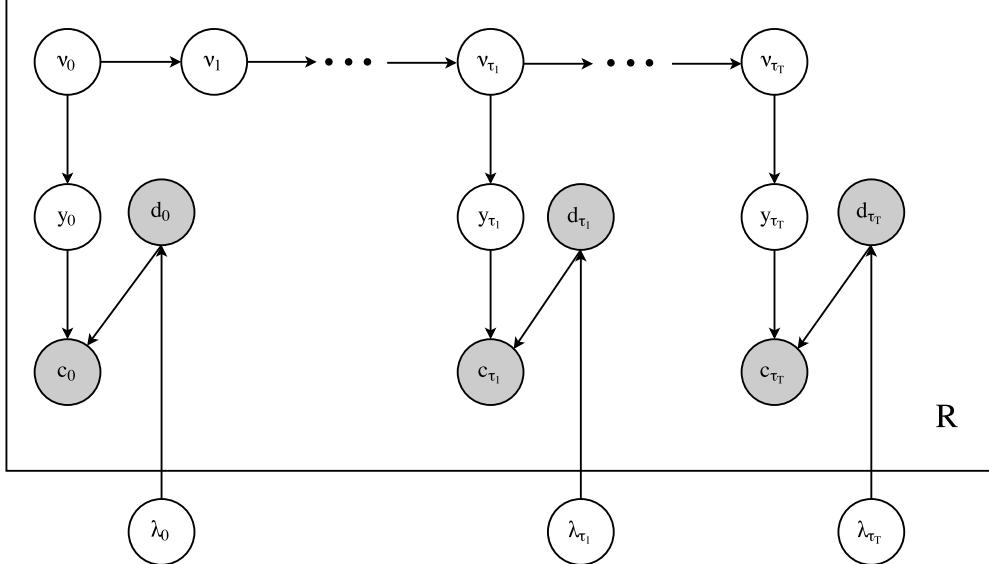
456 Acknowledgments

457 AI, AA, and VB were supported by grants from the NIH (1R01GM114362) and NSF (DBI-1458557
458 and IIS-1318386). CS is supported by the European Research Council grant ArchAdapt.

459 Conflict of interest

460 VB is a co-founder, has an equity interest, and receives income from Digital Proteomics, LLC (DP).
461 The terms of this arrangement have been reviewed and approved by the University of California,
462 San Diego in accordance with its conflict of interest policies. DP was not involved in the research
463 presented here.

464 **Figures**



Generative Process 1: The Generative Process for Dynamic Pool-seq Data.

Input: $N, n, R, \{\lambda_{\tau_0}, \dots, \lambda_{\tau_T}\}, \mathcal{T} = \{\tau_0, \dots, \tau_T\}$

Output: Time-series pool-seq data for R replicates of a single locus $\{\mathbf{c}^{(r)}\}$ and $\{\mathbf{d}^{(r)}\}$.

for $r \leftarrow 1$ **to** R **do**

for $t \leftarrow \tau_0$ **to** τ_T **do**

$2N\nu_t \sim \text{Binomial}(2N, \nu_{t-1})$;

if $t \in \mathcal{T}$ **then**

$d_t^{(r)} \sim \text{Poiss}(\lambda_{\tau_i})$;

$2ny_t \sim \text{Binomial}(2n, \nu_t)$;

$c_t^{(r)} \sim \text{Binomial}(d_t^{(r)}, y_t)$;

end

end

end

Fig 1: Graphical model and generative process of the hidden Markov model for single-locus pool-seq experimental evolution.

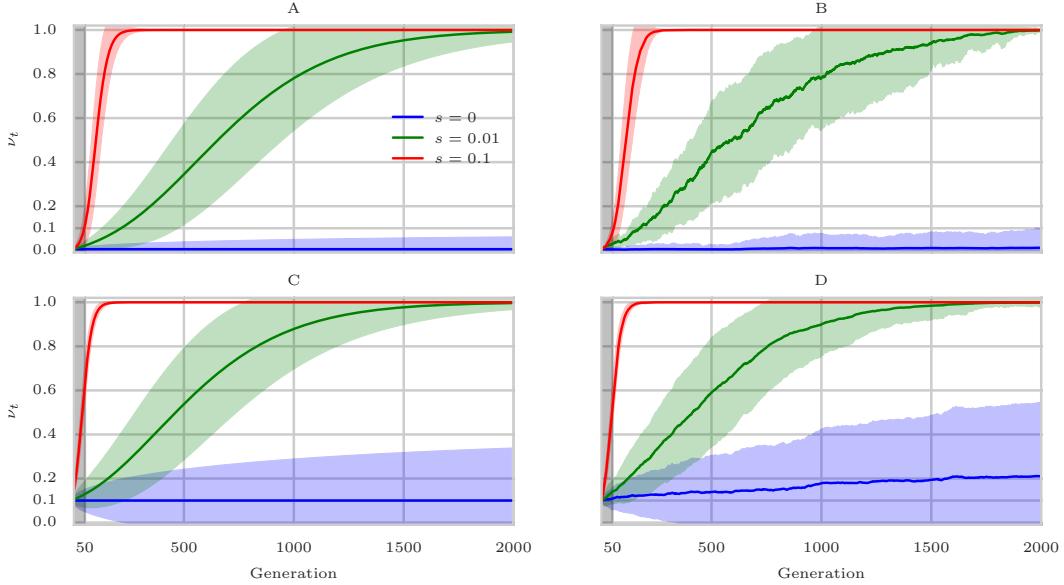


Fig 2: Theoretical and empirical trajectories of favored allele for hard and soft sweep scenarios.

Single-loci diallelic Wright-Fisher Markov chain model (A,C) and empirical (B,D) trajectories of the frequency of the favored allele are computed for 1000 simulations of populations with 1000 diploid individuals. Each curve shows the mean and the 95% confidence interval. Panels A and C depict theoretical calculations of the favored allele frequency under hard (ν_0 is small), and soft sweep due to standing variation (higher ν_0), for a range of values of s . $s = 0$ corresponds to neutral evolution. Similarly, panels B and D show the empirical forward simulations of populations under the same selection regimes, and hard/soft-sweep scenarios. The first 50 generations are shaded in gray to represent the typical sampling span of EE experiments. The plot illustrates the difficulty of EE experiments in having to predict selection at a very early stage of the sweep. The signal is slightly stronger under standing variation scenario. The theoretical and empirical simulations are in close correspondence.

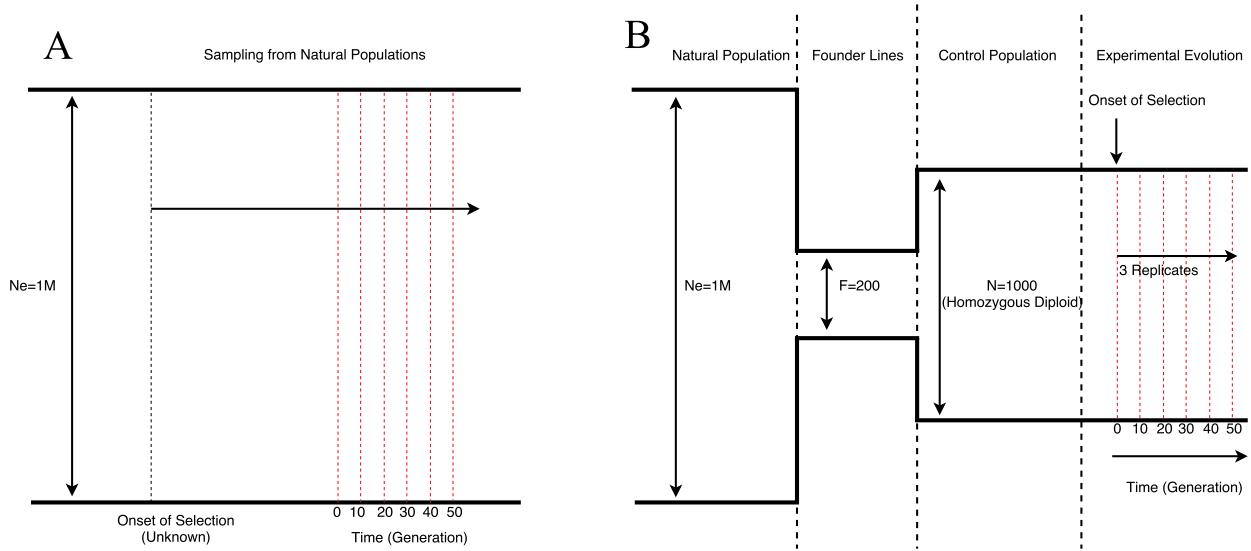


Fig 3: **Two settings for collecting genomic time series data.**

Different settings in which dynamic data is collected are depicted with typical parameters for *D. melanogaster*. In both settings, 6 samples (vertical red dashed lines) are taken every 10 generation. When sampling from naturally evolving populations (A), the time of onset of selection is unknown, and population size is larger. For (controlled) experimental evolution (B), founder lines are first sampled from a natural population to create a homogeneous population. Then, multiple replicates of this population are evolved and sampled over time.

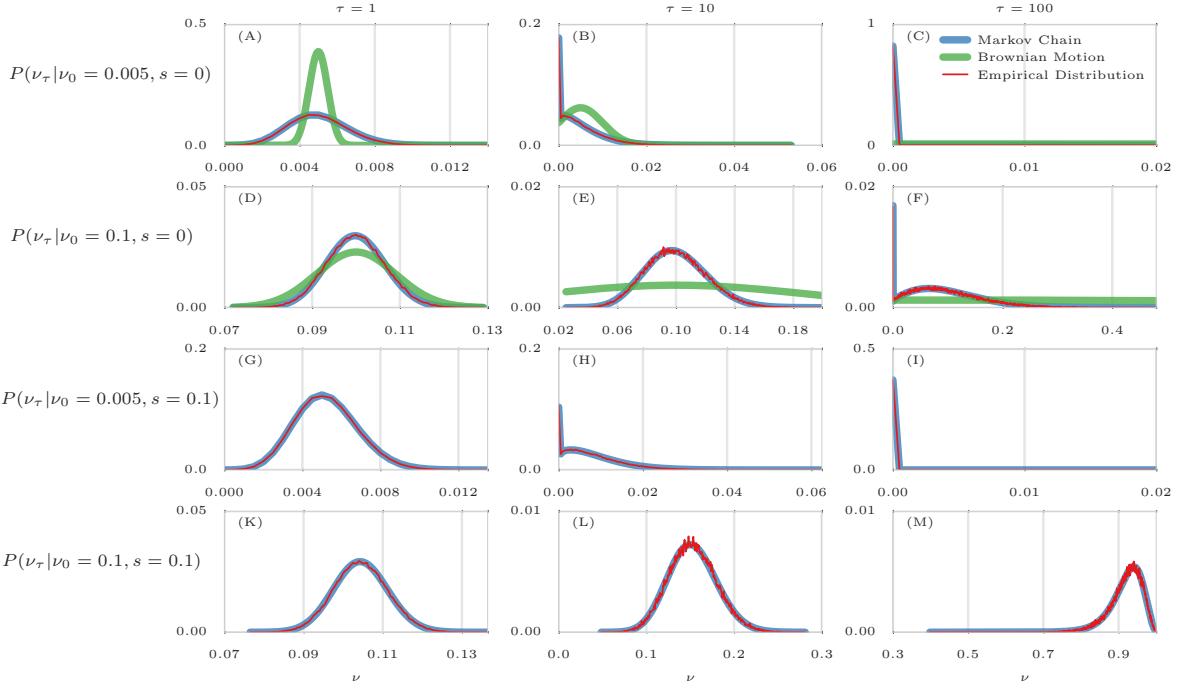


Fig 4: Comparison of empirical distributions of allele frequencies (red) versus predictions from Brownian Motion (green), and Markov chain (blue).

Comparison of empirical and theoretical distributions under neutral evolution (panels A-F) and selection (panels G-M) with different starting frequencies $\nu_0 \in \{0.005, 0.1\}$ and sampling times of $\mathcal{T} = \{0, \tau\}$, where $\tau \in \{1, 10, 100\}$. For each panel, the empirical distribution was computed over 100,000 simulations. Brownian motion (Gaussian approximation) provides poor approximations when initial frequency is far from 0.5 (A) or sampling is sparse (B,C,E,F). In addition, Brownian motion can only provide approximations under neutral evolution. In contrast, Markov chain consistently provide a good approximation in all cases.

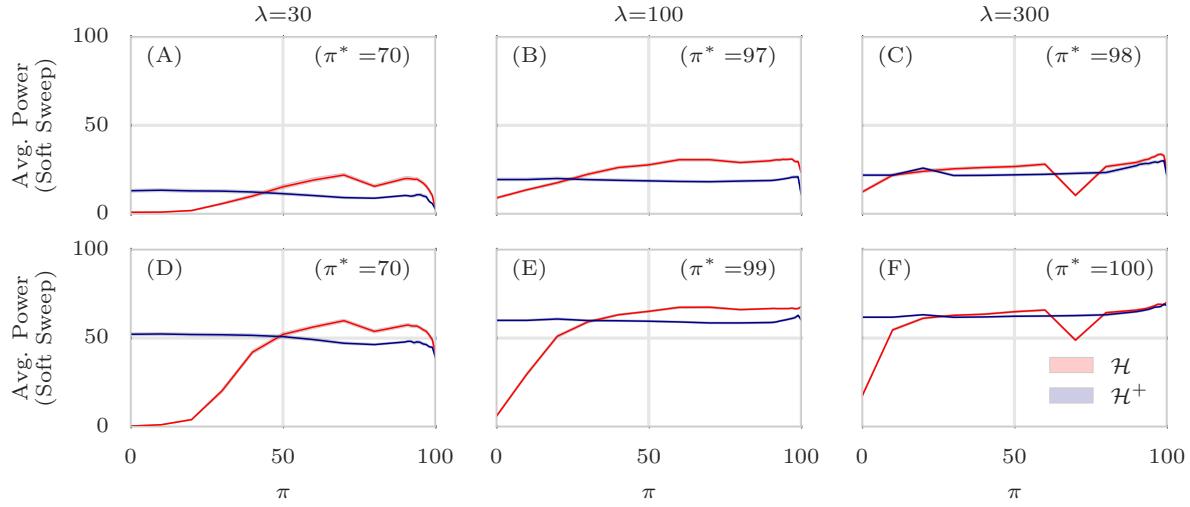


Fig 5: **Average power of CLEAR statistics as function of percentile cutoffs.**

Detection power is averaged for $s \in \{0.025, 0.05, 0.075, 0.1\}$ for CLEAR statistics, \mathcal{H}_π and \mathcal{H}_π^+ . We denote $\pi^* = \arg \max_\pi \{\text{Power}(\mathcal{H}), \text{Power}(\mathcal{H}^+)\}$. Average power was computed using 8000 simulations for each choice of π . The appropriate choice of π can be used to improve performance for different coverage values. In all simulations, 3 replicates are evolved and sampled at generations $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

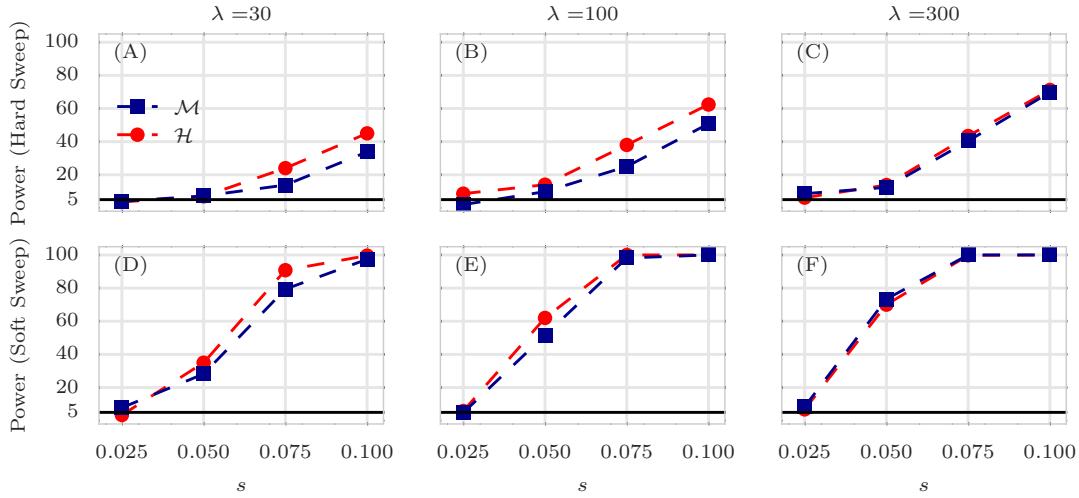


Fig 6: **Comparison of power of Markov chain and HMM.**

Detection power for Markov chain (\mathcal{M}) and HMM (\mathcal{H}) under hard (A-C) and soft sweep (D-F) scenarios, for different coverage λ and selection strength s . The y-axis measures power – sensitivity with false positive rate $\text{FPR} \leq 0.05$ – for 2000 simulations of 50Kbp regions. The horizontal line reflects the power of a random classifier. In all simulations, 3 replicates are evolved and sampled at generations $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

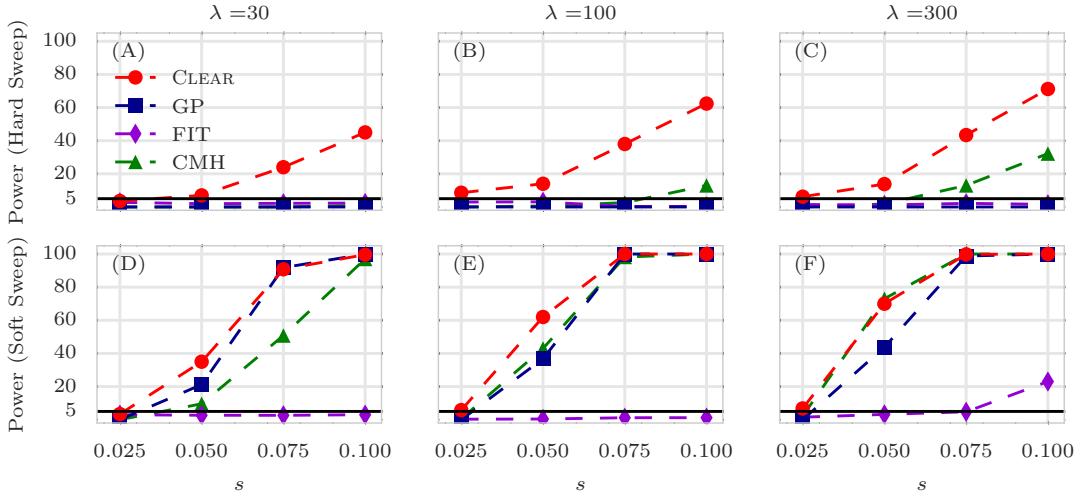


Fig 7: **Power calculations for detection of selection.**

Detection power for $\text{CLEAR}(\mathcal{H})$, Frequency Increment Test (FIT), Gaussian Process (GP), and CMH under hard (A-C) and soft sweep (D-F) scenarios. λ , s denote the mean coverage and selection coefficient, respectively. The y -axis measures power – sensitivity with false positive rate $\text{FPR} \leq 0.05$ – for 2,000 simulations of 50Kbp regions. The horizontal line reflects the power of a random classifier. In all simulations, 3 replicates are evolved and sampled at generations $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

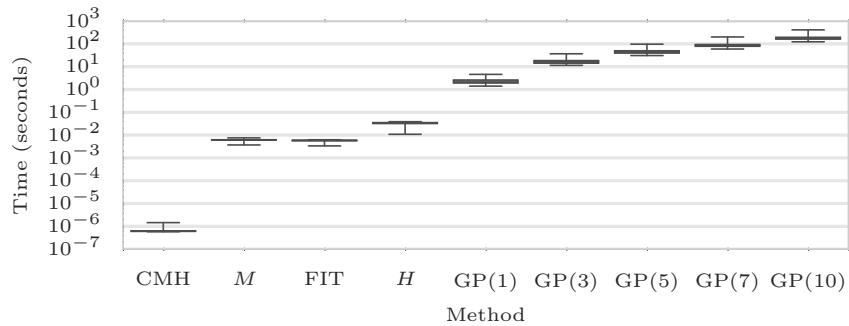


Fig 8: **Running time.**

Box plots of running time per variant (CPU-secs.) of $\text{CLEAR}(\mathcal{M}, \mathcal{H})$, CMH, FIT, and GP with single, 3, 5, 7, and 10 loci over 1000 simulations conducted on a workstation with 4th Generation Intel Core i7 processor. The average running time for each method is shown on the x-axis. In all simulations, 3 replicates are evolved and sampled at generations $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

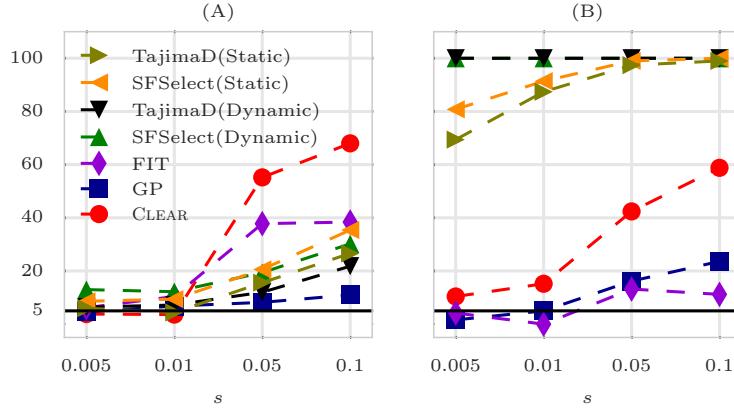


Fig 9: **Power of SFS based statistics.**

Power of detecting selection for Frequency Increment Test (FIT), Gaussian Process (GP), CLEAR(\mathcal{H}) on hard sweep natural experimental evolution with $N_e = 10^4$ and coverage $\lambda = \infty$. The measurements are conducted for a range of selection coefficients, s . Each point represents the mean of 200 simulations. For each simulation, sampling starts at a randomly chosen time, and subsequently 5 replicate samples are acquired every 10 generations. (A) Start of sampling is chosen randomly throughout the sweep $\tau_1 \sim U[1, t_{\nu=1}(s, N_e)]$, where $t_{\nu=x}(s, N_e)$ denotes the expected time to reach carrier frequency x in a hard sweep and $U[a, b]$ is discrete uniform distribution. (B) The start of sampling is chosen near fixation of the favored allele, i.e. $\tau_1 \sim U[t_{\nu=0.9}(s, N_e), t_{\nu=1}(s, N_e)]$.

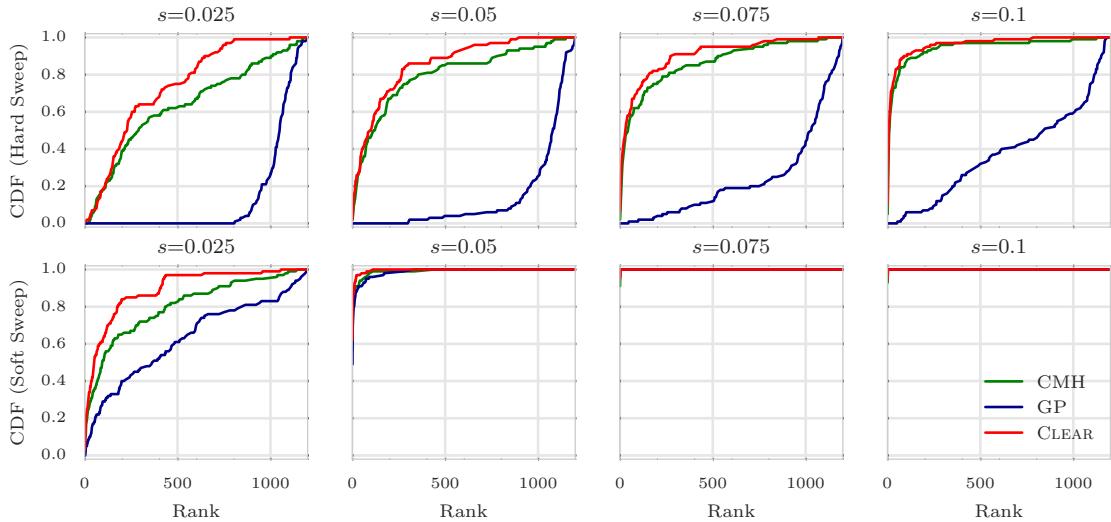


Fig 10: **Ranking performance for 100 \times coverage.**

Cumulative Distribution Function (CDF) of the distribution of the rank of the favored allele in 1000 simulations for CLEAR (H), Gaussian Process (GP), CMH, and Frequency Increment Test (FIT), for different values of selection coefficient s and initial carrier frequency. Note that the individual variant CLEAR score (H) is used to rank variants. The Area Under Curve (AUC) is computed as an overall quantitative measure to compare the performance of methods for each configuration. In all simulations, 3 replicates are evolved and sampled at generations $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

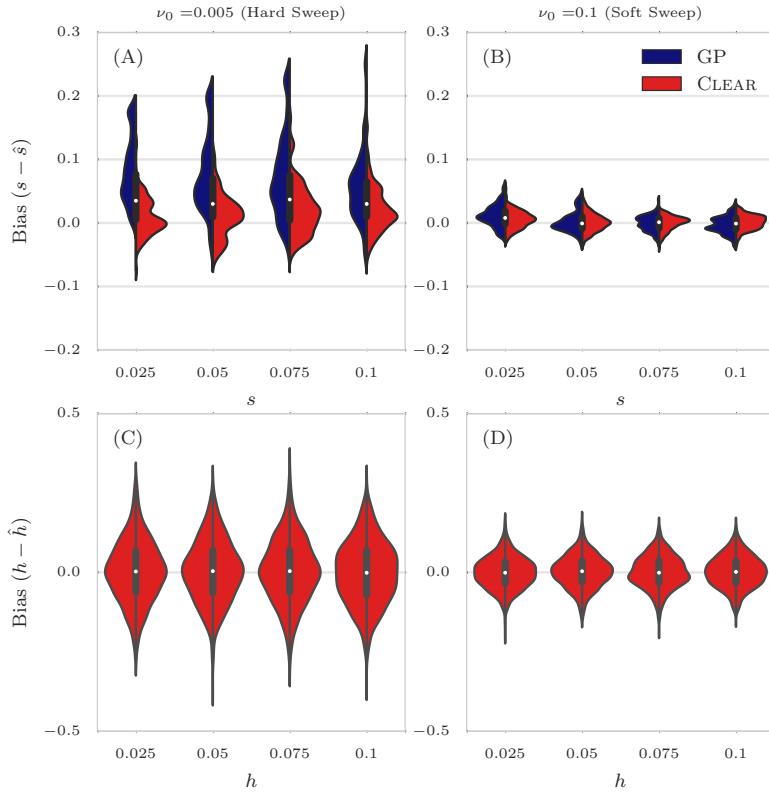


Fig 11: **Distribution of bias for 100 \times coverage.**

The distribution of bias ($s - \hat{s}$) in estimating selection coefficient over 1000 simulations using Gaussian Process (GP) and CLEAR (H) is shown for a range of choices for the selection coefficient s and starting carrier frequency ν_0 , when coverage $\lambda = 100$ (Panels A,B). GP and CLEAR have similar variance in estimates of s for soft sweep, while CLEAR provides lower variance in hard sweep. Also see S3 Table. Panels C,D show the variance in the estimation of h . In all simulations, 3 replicates are evolved and sampled at generations $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

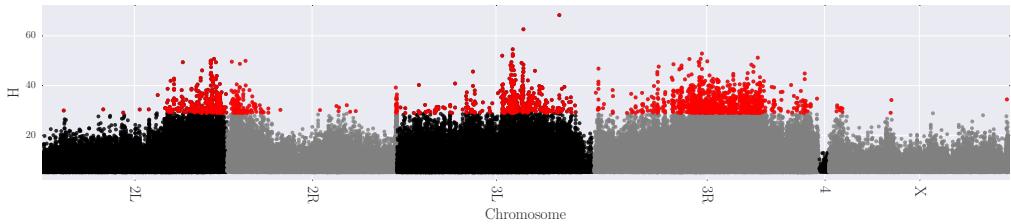


Fig 12: **CLEAR scan of the data from a study of *D. melanogaster* adaptation to alternating temperatures.**

Manhattan plot of the CLEAR (\mathcal{H}^+) statistic (A) and the number of SNPs (B) in 30Kbp sliding windows with steps of 10Kbp, excluding heterochromatic regions. Regions that exceed the local FDR threshold of \mathcal{H}^+ are shown in red dots.

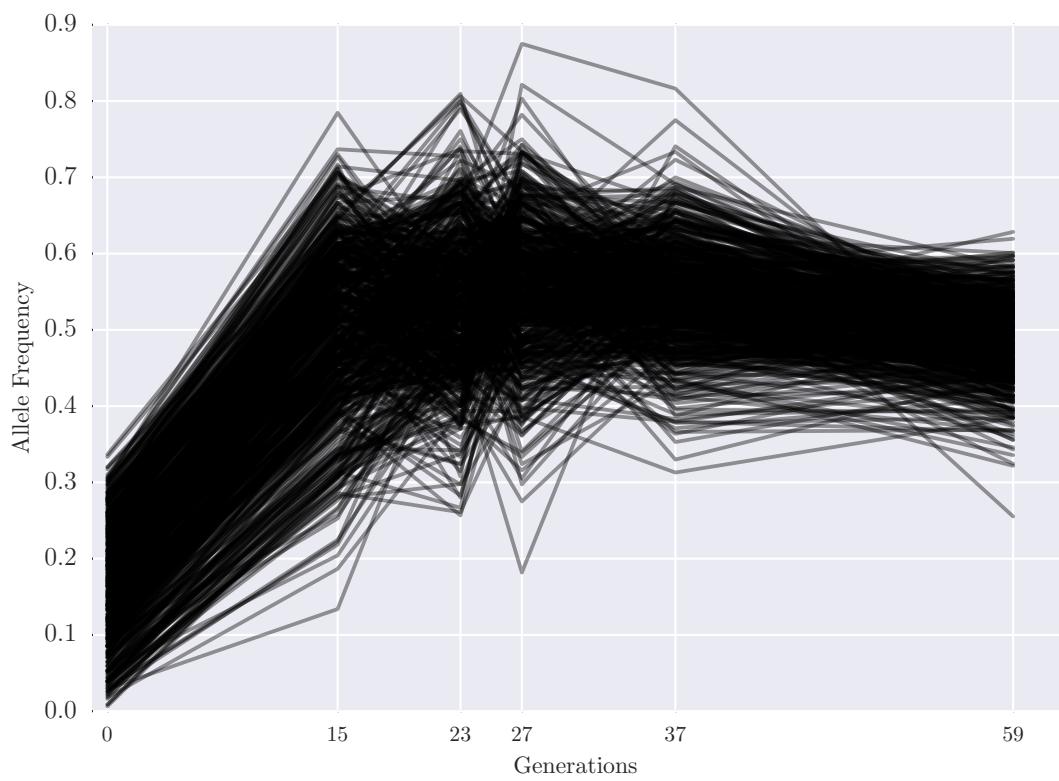


Fig 13: **dominant variants.**

465 **S1 Text Dynamics of Site Frequency Spectrum-based Statistics**
 466 **and Linkage Disequilibrium under Selection**

467 **S1.1 Text An approximate logistic function for allele frequency dynamics**

468 Assume that a site is evolving under selection constraints $s, h \in \mathbb{R}$, where s and h denote selection
 469 strength and overdominance, respectively. Let ν_t denote the frequency of the site at time $\tau_t \in \mathcal{T}$.
 470 Then, ν_{t+} , the frequency at time $\tau_t + 1$ can be estimated using:

$$\hat{\nu}_{t+} = \nu_t + \frac{s(h + (1 - 2h)\nu_t)\nu_t(1 - \nu_t)}{1 + s\nu_t(2h + (1 - 2h)\nu_t)}. \quad (\text{S1})$$

471 We can show that the dynamic of the favored allele can be modeled via a logistic function, in the
 472 case of directional selection ($h = 0.5$). Taking derivatives of Eq. S1, we have

$$\frac{d\nu_t}{dt} = \frac{s\nu_t(1 - \nu_t)}{2 + 2s\nu_t} \quad (\text{S2})$$

473 To, solve the differential equation, note that for small s , $2 + 2s\nu_t \approx 2$. Substituting,

$$\nu_t = \frac{1}{1 + \frac{1-\nu_0}{\nu_0} e^{-st/2}} = \sigma(st/2 + \eta(\nu_0)) \quad (\text{S3})$$

474 where $\sigma(\cdot)$ is the logistic function and $\eta(\cdot)$ is logit function (inverse of the logistic function).

475 **SFS appendices removed.**

476 **S1.2 Text Linkage Disequilibrium**

477 Nonrandom associations, Linkage Disequilibrium (LD), between polymorphisms are established in
 478 the substitution process, broken by recombination events and reinforced by selection. Although
 479 LD can not be measured in pooled sequencing data (phased haplotype data is required), it is still
 480 worthwhile to examine the behavior of LD as a result of the interaction between recombination and
 481 natural selection. In this part we theoretically overview expected LD in short EEs.

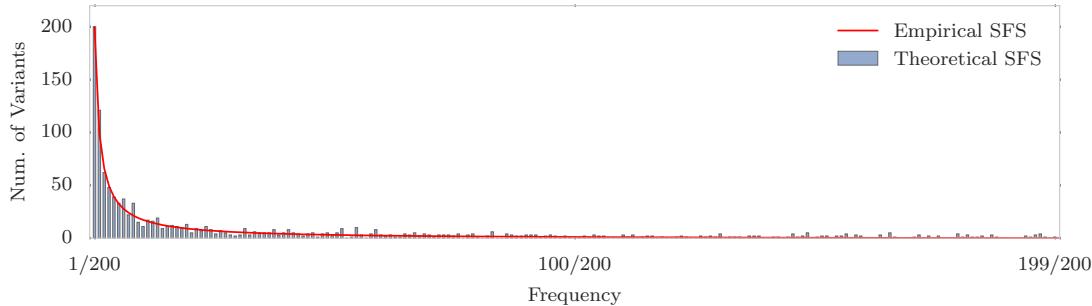
482 Let ρ_0 be the LD at time zero between the favored allele and a segregating site l base-pairs
 483 away, then under natural selection we have

$$\rho_t = \alpha_t \beta_t \rho_0 = e^{-rtl} \left(\frac{K_t}{K_0} \right) \rho_0 \quad (\text{S4})$$

484 where $K_t = 2\nu_t(1 - \nu_t)$ is the heterozygosity at the selected site, r is the recombination rate/bp/gen.
 485 The *decay factor*, $\alpha_t = e^{-rtl}$, and *growth factor*, β_t (see Eqs. 30-31 in [73]), are result of recom-
 486 bination and selection, respectively. S13 Fig presents the expected theoretical value of LD when
 487 $\rho_0 = 0.5$ between favored allele (site at position 500K) and the rest of genome, and $\nu_0 = 0.1$. For
 488 neutral evolution (top), LD decays exponentially through space and time, while in natural selection
 489 (bottom), LD increases and then decreases. Interestingly, LD increases to its maximum value, 1,
 490 for the nearby region (the plateau in S13 Fig bottom) of the favored allele.

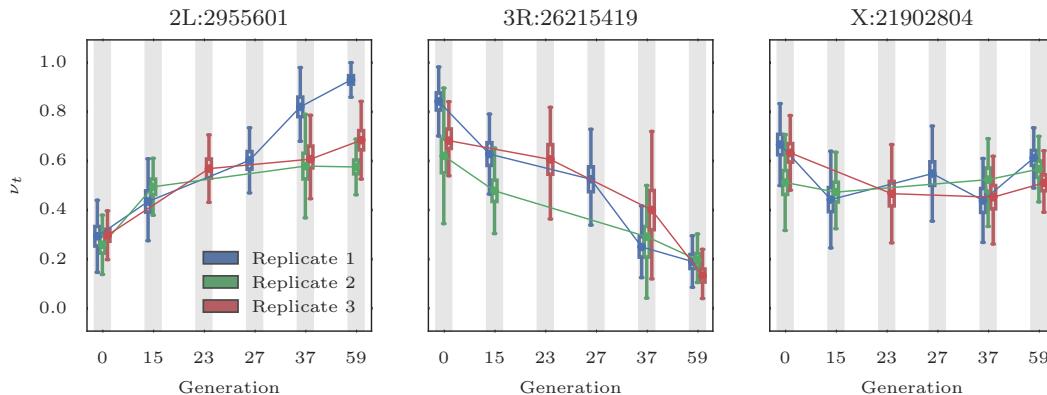
491 In principle, LD increases after the onset of selection, until $\log(\alpha_t) + \log(\beta_t) > 0$, see Eq. S4.
 492 Specifically, log of decay term is linear and, using Eq. S3, we write growth factor in term of initial
 493 frequency ν_0 and selection strength s . S14 Fig, S15 Fig, S16 Fig, and S17 Fig depict interaction of
 494 decay and growth factors for weak and strong selection and soft and hard sweeps. In all the case,
 495 LD of the favored allele with a segregating site 50Kbp away, increases in the first 50 generations,
 496 which give rise to increasing number of *hitchhikers*.

497 Increase of LD in a large (100Kbp) region is particularly advantageous to the task of identifying
 498 the region under selection, if the composite statistics is used. As a result, \mathcal{H} statistic outperforms
 499 existing (single-loci) tools in identifying selection. In contrast, augmentation of LD, increases the
 500 number of candidates for the favored allele, which makes it difficult to localize the favored allele.



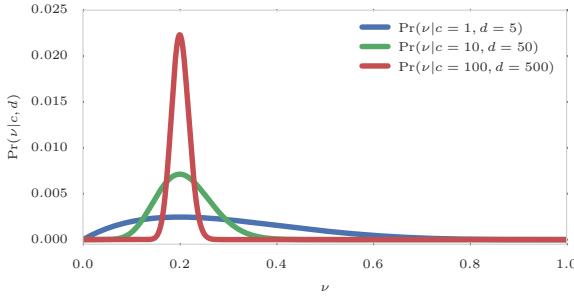
S1 Fig: **Site Frequency Spectrum.**

Theoretical and Empirical SFS in a 50Kbp region for a neutral population of 200 individuals when $N_e = 10^6$ and $\mu = 10^{-9}$. The x -axis corresponds to site frequency, and the y -axis to the number of variants with a specific frequency. In a neural population, majority of the variations stand in low frequency.



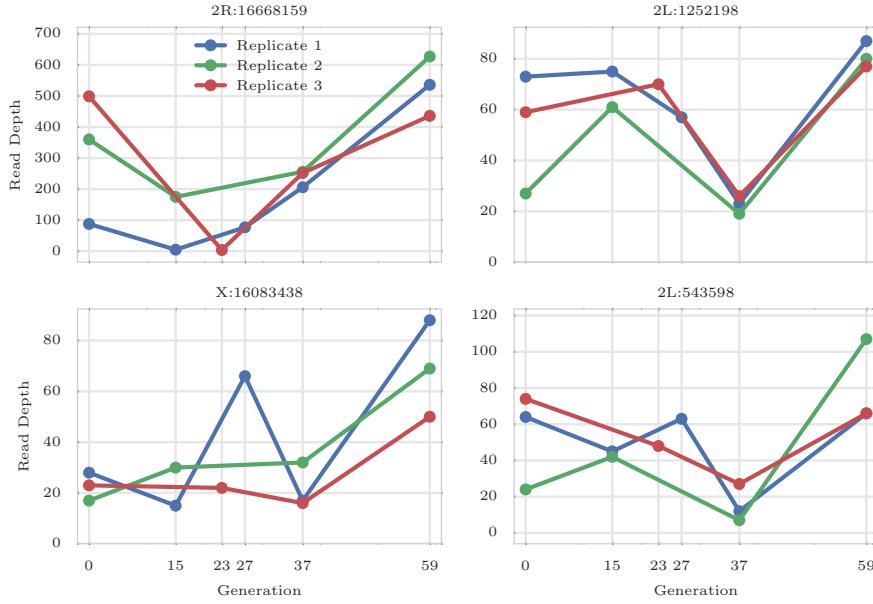
S2 Fig: **Trajectory of pool-sequenced variants.**

Trajectory of three different variants that are increasing in frequency over time. Note that for read count data, the true allele frequency is not known. Here we draw the posterior distribution of the allele frequency at each time point using box plot. The median of each distribution is denoted by dots. The variance of each box is seen to be inversely related to the depth of the measurement. For instance, generation 59 is sequenced with higher coverage than generation 37. As a result, variance of observations in generation 59 is considerably smaller.



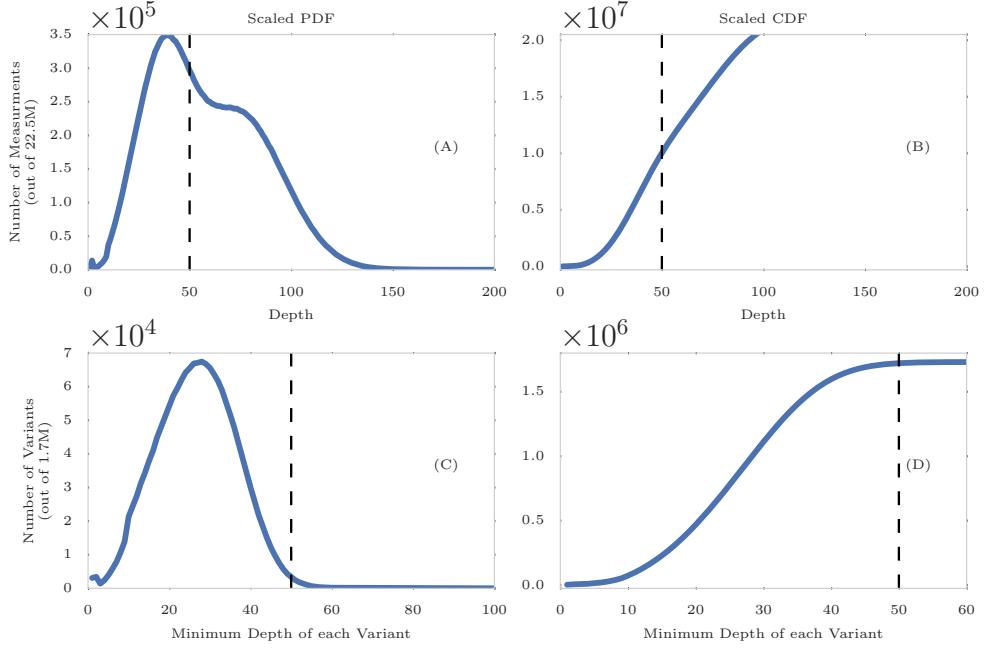
S3 Fig: Posterior distribution of allele frequency.

Distribution of hidden allele frequency for different values of depth $d = \{5, 50, 500\}$. In all cases, the true frequency is 0.2. The estimated frequency values are binomially distributed, with different variances, around the true value in all cases with.



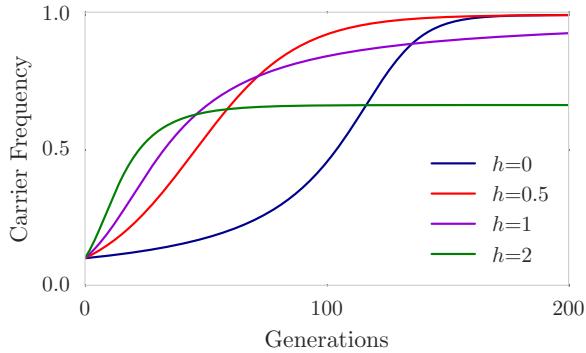
S4 Fig: Coverage heterogeneity in time series data.

Each panel shows the read depth for 3 replicates of the data from a study of *D. melanogaster* adaptation to alternating temperatures data (see section 3.1). Heterogeneity in depth of coverage is seen between replicates, and also at different time points, in all 4 variants. None of these sites pass the hard filtering with minimum depth of 30.



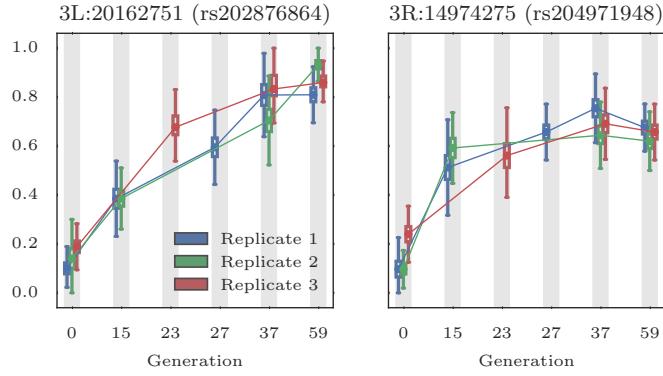
S5 Fig: Distribution of depth in the real data.

Scaled PDF (A) and CDF (B) of the read depths of all ($\approx 22.1M$) measurements, i.e., all replicates and time points of the all ($\approx 1.7M$) variants. Scaled PDF (C) and CDF (D) of the minimum depth of sites. While more than half most ($\approx 12.5M$) of the measurements have depth of 50 or greater (dashed line in (A),(B)), only a small fraction ($\approx 11K$) of variants (dashed line in (C),(D)) pass the filter of having minimum depth of 50.

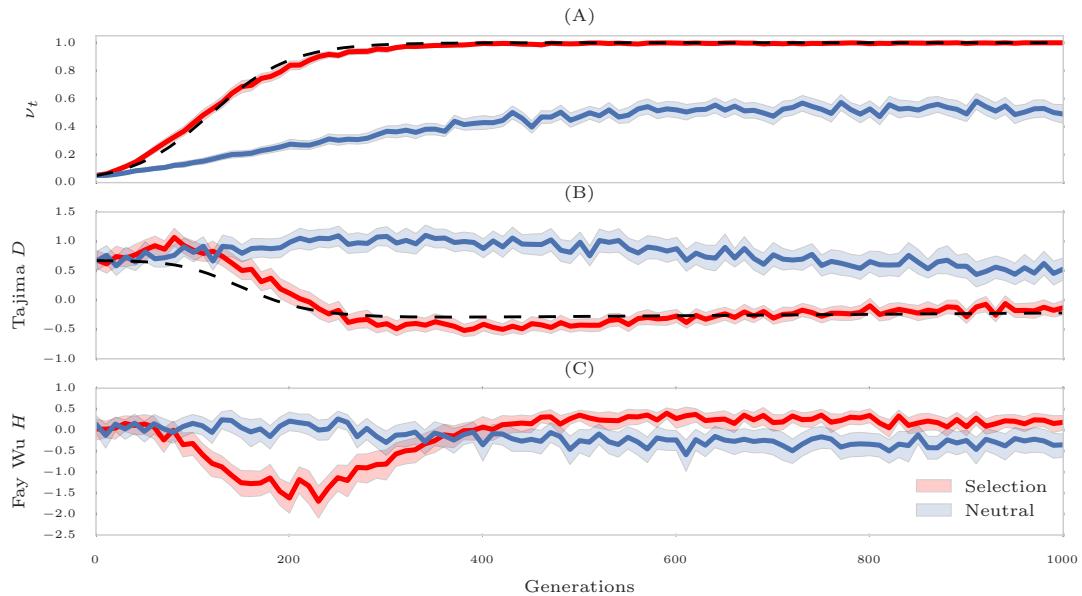


S6 Fig: Natural selection in infinite population.

Trajectory of the favored allele in an infinite population with $s = 0.1$ for $h = 0$ (recessive favored allele), $h = 0.5$ (directional selection), $h = 1$ (dominant favored allele) and $h = 2$ (overdominant favored allele).

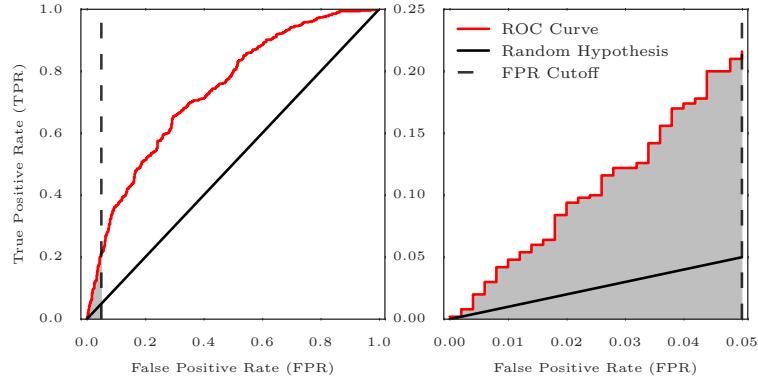


S7 Fig: Variant showing strong signal of directional (Left) and over-dominant, aka balancing selection (Right).



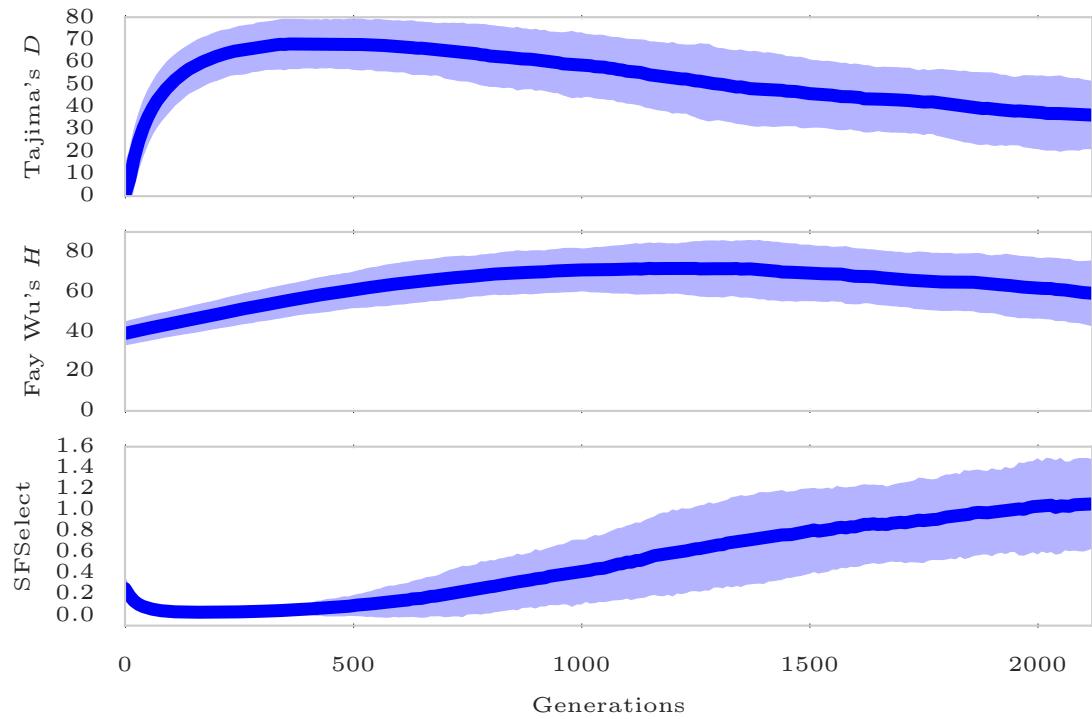
S8 Fig: Dynamic SFS-based statistics.

Mean and 95% CI of 100 simulations for neutral (blue trajectories) selection with $s = 0.1$ (red trajectories). In all case, statistic computed for a 50Kbp window and $N_e = 10^4$, $\mu = 10^{-9}$. The dashed line shows the parametric models derived in ??.



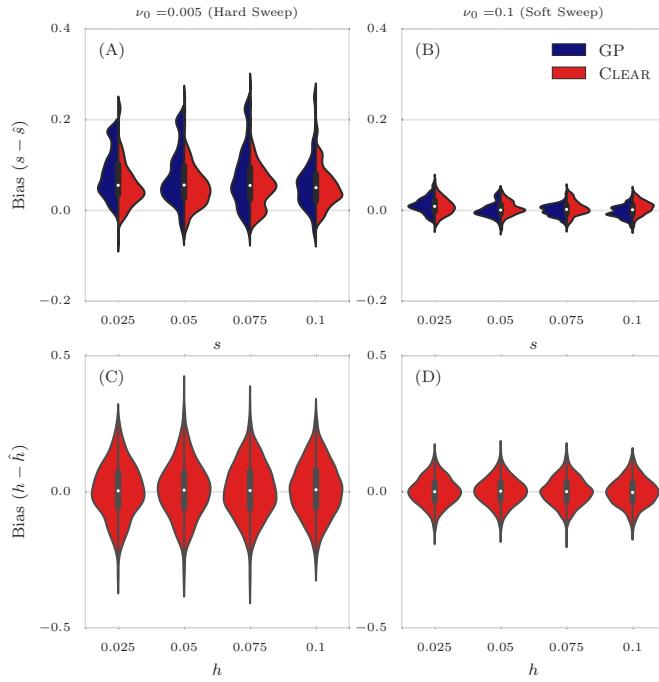
S9 Fig: Schematic for computing power of detecting selection.

Receiver Operating Characteristics (ROC) curve (left) for classification of 2000 simulations (1000 selection and 1000 neutral). The Area Under the Curve (AUC) represents overall performance. The diagonal black line represents performance of a random hypothesis which achieves Area under the curve (AUC) of 0.5. To avoid computing AUC for the regions where FPR is unacceptably high, we restrict ROC curve to the region where $FPR \leq 0.05$ (right). In this case, we define power to be the (scaled) AUC of the restricted region.



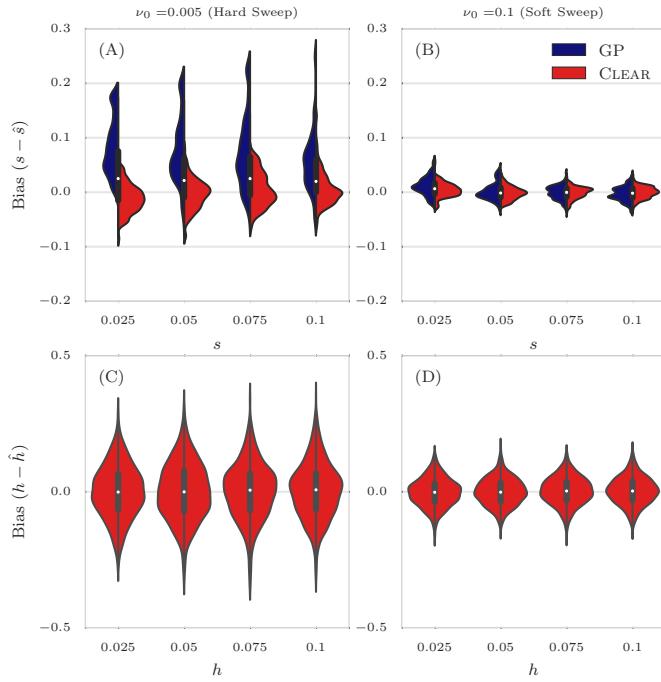
S10 Fig: **Effect of bottleneck in a neutral experimental evolution experiment with restricted number of founder lines.**

Dynamic of SFS based statistics under neutral evolution when $F = 200$ founders were selected from a larger population size ($N_e = 10^6$). The statistics for Tajima's D , Fay Wu's H and SFSelect were computed for 1000 neutral simulations and the mean and 95% confidence interval plotted. Under neutral evolution, all the statistics are expected to vary around a fixed mean through time. However, under selective constraint, D and H take negative values, while SFSelect take positive values. In experimental evolution, bottleneck effect will suppress the signal of selection, especially in early generations.



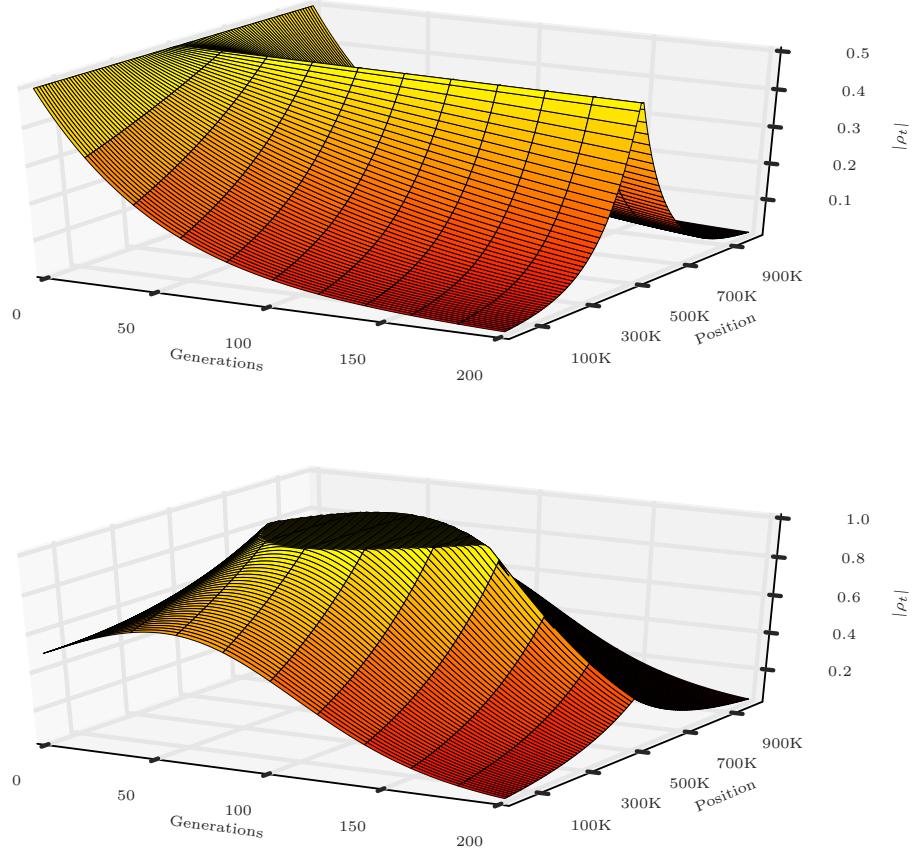
S11 Fig: Distribution of bias for $30\times$ coverage.

The distribution of bias ($s - \hat{s}$) in estimating selection coefficient over 1000 simulations using Gaussian Process (GP) and CLEAR (H) is shown for a range of choices for the selection coefficient s and starting carrier frequency ν_0 , when coverage $\lambda = 30$ (Panels A,B). GP and CLEAR have similar variance in estimates of s for soft sweep, while CLEAR provides lower variance in hard sweep. Also see [S3 Table](#). Panels C,D show the variance in the estimation of h .



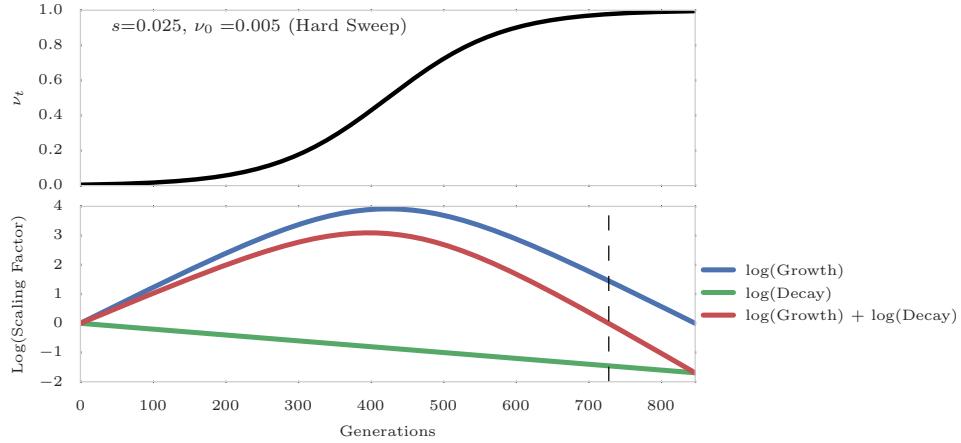
S12 Fig: **Distribution of bias for infinite coverage.**

The distribution of bias ($s - \hat{s}$) in estimating selection coefficient over 1000 simulations using Gaussian Process (GP) and CLEAR (H) is shown for a range of choices for the selection coefficient s and starting carrier frequency ν_0 , when coverage $\lambda = \infty$ (Panels A,B). GP and CLEAR have similar variance in estimates of s for soft sweep, while CLEAR provides lower variance in hard sweep. Also see [S3 Table](#). Panels C,D show the variance in the estimation of h .



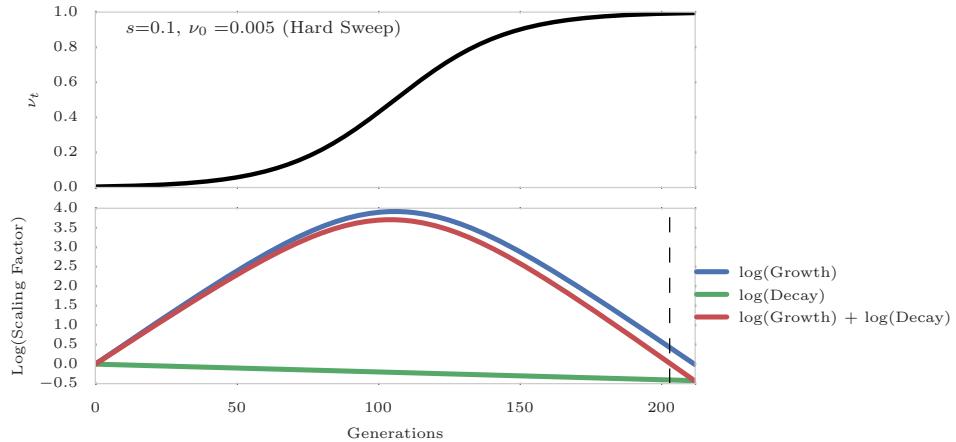
S13 Fig: **Expected dynamic of LD under selection and neutral evolution.**

Dynamic of LD (ρ_t) of a 1Mbp genome to the favored allele (at position 500K) is drawn as function of position and time for neutral (top) and selection(bottom) regimes. For sake of illustration, we assumed that at generations 0, LD of all variants with the favored allele is 0.5, initial frequency of the favored allele is 0.1, recombination rate is $r = 2 \times 10^{-8}$ (top). The selection strength is 0 and 0.05 for neutral and selection regimes, respectively. As expected LD decay exponentially through space and time. However, selection causes LD to increase then decrease.



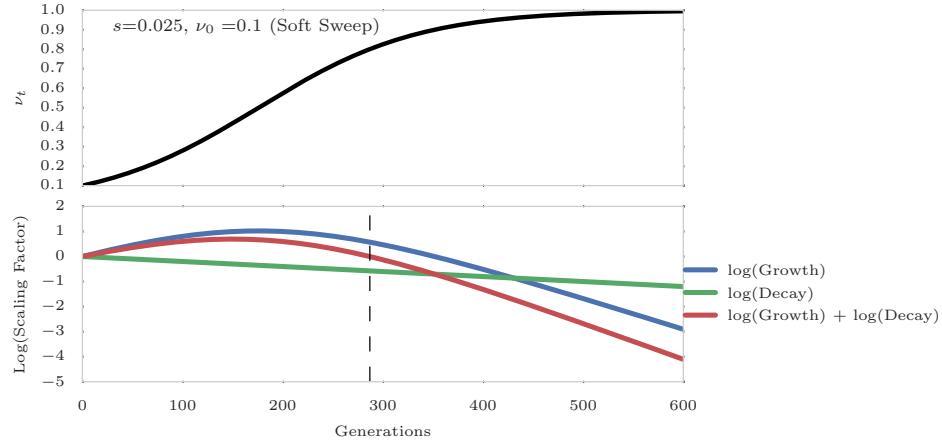
S14 Fig: Interaction between growth and decay factors of LD in hard sweep with weak selection.

Expected dynamics of LD under weak selection ($s = 0.025$) with hard sweep ($\nu_0 = 0.005$). In addition to recombination, initial frequency of the favored allele and selection strength determine the dynamics of LD. The vertical dashed line denotes the time in which LD start to decrease. LD between the favored allele and the rest of genome increase for ≈ 700 generations after onset of selection, implying that localizing adaptive allele in short term experimental evolution is a difficult task, especially when the frequency of the favored allele is low.



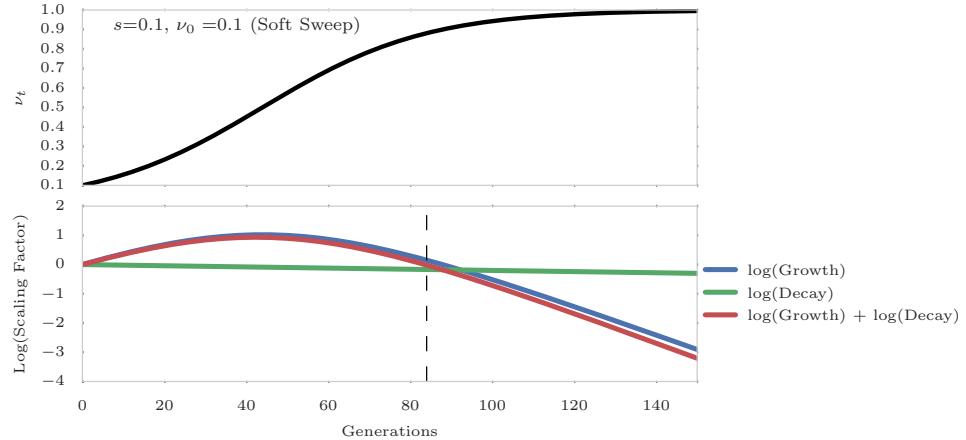
S15 Fig: Interaction between growth and decay factors of LD in hard sweep with strong selection.

Expected dynamics of LD under strong selection ($s = 0.1$) with hard sweep ($\nu_0 = 0.005$). In addition to recombination, initial frequency of the favored allele and selection strength determine the dynamics of LD. The vertical dashed line denotes the time in which LD start to decrease. LD between the favored allele and the rest of genome increase for ≈ 200 generations after onset of selection, implying that localizing adaptive allele in short term experimental evolution is a difficult task, especially when the frequency of the favored allele is low.



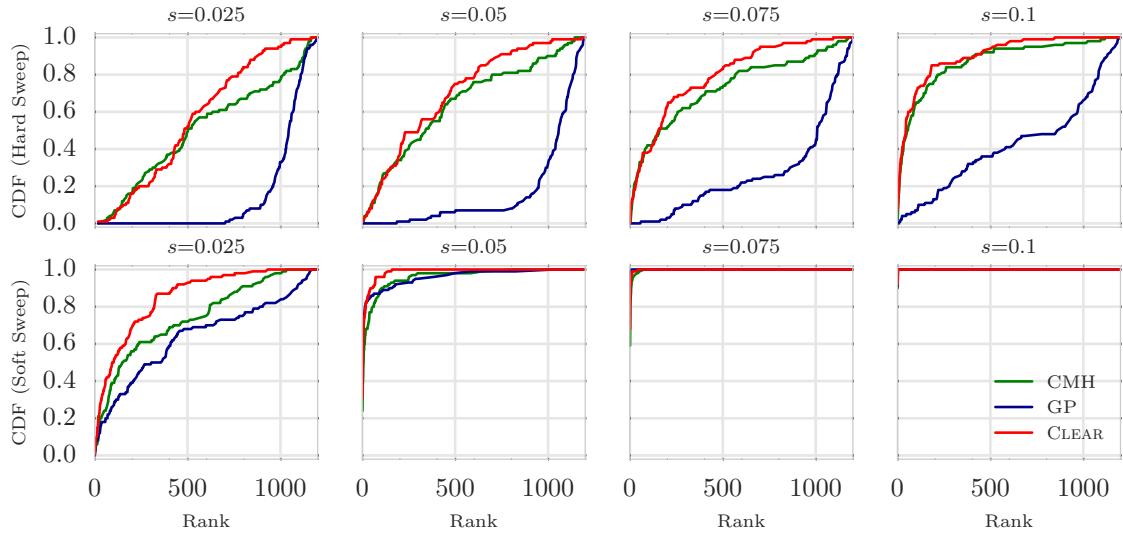
S16 Fig: Interaction between growth and decay factors of LD in soft sweep with weak selection.

Expected dynamics of LD under weak selection ($s = 0.025$) for soft sweep ($\nu_0 = 0.1$). In addition to recombination, initial frequency of the favored allele and selection strength determine the dynamics of LD. The vertical dashed line denotes the time in which LD start to decrease. LD between the favored allele and the rest of genome increase for ≈ 300 generations after onset of selection, implying that localizing adaptive allele in short term experimental evolution is a difficult task.



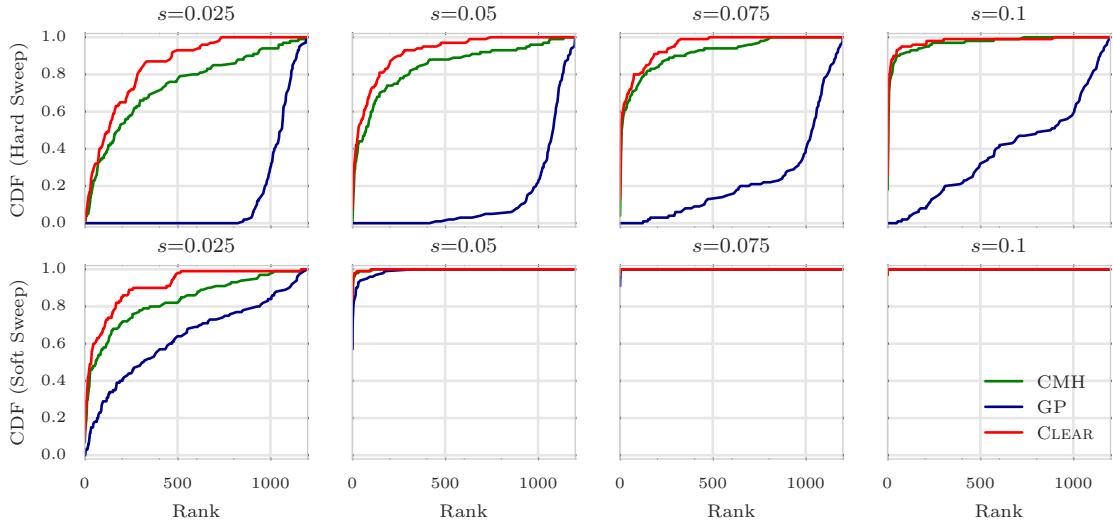
S17 Fig: Interaction between growth and decay factors of LD in soft sweep with strong selection.

Expected dynamics of LD under strong selection ($s = 0.1$) for soft sweep ($\nu_0 = 0.1$). In addition to recombination, initial frequency of the favored allele and selection strength determine the dynamics of LD. The vertical dashed line denotes the time in which LD start to decrease. LD between the favored allele and the rest of genome increase for ≈ 80 generations after onset of selection, implying that localizing adaptive allele in short term experimental evolution is a difficult task.



S18 Fig: **Ranking performance for $30\times$ coverage.**

Cumulative Distribution Function (CDF) of the distribution of the rank of the favored allele in 1000 simulations for CLEAR (H score), Gaussian Process (GP), and Cochran Mantel Haenszel (CMH), for different values of selection coefficient s and initial carrier frequency.



S19 Fig: **Ranking performance for $300\times$ coverage.**

Cumulative Distribution Function (CDF) of the distribution of the rank of the favored allele in 1000 simulations for CLEAR (H score), Gaussian Process (GP), and Cochran Mantel Haenszel (CMH), for different values of selection coefficient s and initial carrier frequency.

S1 Table: Average of power for detecting selection.

Hard Sweep			Soft Sweep		
λ	Method	Avg Power	λ	Method	Avg Power
300	CLEAR	34	300	CLEAR	69
300	CMH	12	300	CMH	69
300	FIT	2	300	GP	61
300	GP	0	300	FIT	8
100	CLEAR	31	100	CLEAR	67
100	CMH	4	100	CMH	60
100	FIT	2	100	GP	59
100	GP	0	100	FIT	1
30	CLEAR	20	30	CLEAR	57
30	FIT	2	30	GP	53
30	CMH	0	30	CMH	39
30	GP	0	30	FIT	3

Average power is computed for 8000 simulations with $s \in \{0.025, 0.05, 0.075, 0.1\}$. Frequency Increment Test (FIT), Gaussian Process (GP), CLEAR (H statistic) and Cochran Mantel Haenszel (CMH) are compared for different initial carrier frequency ν_0 . For all sequencing coverages, CLEAR outperform other methods. When coverage is not high ($\lambda \in \{30, 100\}$) and initial frequency is low (hard sweep), CLEAR significantly perform better than others.

S2 Table: Average running time per variant in seconds for different methods.

Method	Avg. Time per Locus
CMH	0.0
M	0.006
FIT	0.006
H	0.033
GP(1)	2.551
GP(3)	19.177
GP(5)	50.291
GP(7)	95.602
GP(10)	202.017

S3 Table: Mean and standard deviation of the distribution of bias ($s - \hat{s}$) of 8000 simulations with coverage $\lambda = 100 \times$ and $s \in \{0.025, 0.05, 0.075, 0.1\}$.

Method	ν_0	Mean	STD
GP	0.005	0.073	0.061
CLEAR	0.005	0.016	0.035
GP	0.1	0.002	0.016
CLEAR	0.1	0.002	0.013

S4 Table: GO enrichment analysis of data from a study of *D. melanogaster* adaptation to alternating temperatures using Gowinda.

GO ID	GO Term	-log(<i>p</i> -value)
GO:0001558	regulation of cell growth	4.1
GO:0001700	embryonic development via the syncytial blastoderm	4.1
GO:0003341	cilium movement	4.1
GO:0006030	chitin metabolic process	3.8
GO:0006355	regulation of transcription, DNA-templated	4.1
GO:0006367	transcription initiation from RNA polymerase II promoter	4.1
GO:0006508	proteolysis	4.1
GO:0006719	juvenile hormone catabolic process	4.1
GO:0006839	mitochondrial transport	4.1
GO:0007018	microtubule-based movement	4.1
GO:0007269	neurotransmitter secretion	3.6
GO:0007291	sperm individualization	4.1
GO:0007298	border follicle cell migration	4.1
GO:0007475	apposition of dorsal and ventral imaginal disc-derived wing surfaces	4.1
GO:0007552	metamorphosis	3.8
GO:0007602	phototransduction	4.1
GO:0008104	protein localization	3.1
GO:0008340	determination of adult lifespan	4.1
GO:0008362	chitin-based embryonic cuticle biosynthetic process	4.1
GO:0009312	oligosaccharide biosynthetic process	3.0
GO:0009408	response to heat	4.1
GO:0015991	ATP hydrolysis coupled proton transport	4.1
GO:0016079	synaptic vesicle exocytosis	4.1
GO:0016485	protein processing	4.1
GO:0031146	SCF-dependent proteasomal ubiquitin-dependent protein catabolic process	4.1
GO:0035556	intracellular signal transduction	4.1
GO:0042742	defense response to bacterium	3.8
GO:0043066	negative regulation of apoptotic process	3.1
GO:0045494	photoreceptor cell maintenance	4.1
GO:0045664	regulation of neuron differentiation	4.1
GO:0045861	negative regulation of proteolysis	4.1
GO:0048675	axon extension	4.1
GO:0055114	oxidation-reduction process	3.1
GO:0061024	membrane organization	4.1

S5 Table: Enriched genes of analysis of data from a study of *D. melanogaster* adaptation to alternating temperatures associated with GO terms of “cold acclimation” and “response to heat”.

FlyBase ID	GO Term	Gene Name
FBgn0001224	cold acclimation	Hsp23
FBgn0001225	cold acclimation	Hsp26
FBgn0001233	cold acclimation	Hsp83
FBgn0034758	cold acclimation	CG13510
FBgn0001224	response to heat	Hsp23
FBgn0001225	response to heat	Hsp26
FBgn0001233	response to heat	Hsp83
FBgn0001223	response to heat	Hsp22
FBgn0001226	response to heat	Hsp27
FBgn0001227	response to heat	Hsp67Ba
FBgn0001228	response to heat	Hsp67Bb
FBgn0001229	response to heat	Hsp67Bc
FBgn0003301	response to heat	rut
FBgn0004575	response to heat	Syn
FBgn0010303	response to heat	hep
FBgn0019949	response to heat	Cdk9
FBgn0023517	response to heat	Pgam5
FBgn0025455	response to heat	CycT
FBgn0026086	response to heat	Adar
FBgn0035982	response to heat	CG4461

S6 Table: General statistics of analysis of data from a study of *D. melanogaster* adaptation to alternating temperatures.

Statistic	Value
Num. of Variants	1,608,032
Num. of Candidate Intervals	89
Total Num. of Genes	17,293
Num. of Variant Genes	12,834
Num. of Genes within Candidate Intervals	968
Total Num. of GO	6,983
Num. of GO with 3 or More Genes	3,447
Num. of Candidate Variants for Gowinda	2,886

References

- [1] Guillaume Achaz. Frequency spectrum neutrality tests: one for all and all for one. *Genetics*, 183(1):249–258, 2009.
- [2] Alan Agresti and Maria Kateri. *Categorical data analysis*. Springer, 2011.
- [3] Joshua M Akey. Constructing genomic maps of positive selection in humans: Where do we go from here? *Genome research*, 19(5):711–722, 2009.
- [4] Eric C Anderson, Ellen G Williamson, and Elizabeth A Thompson. Monte Carlo evaluation of the likelihood for Ne from temporally spaced samples. *Genetics*, 156(4):2109–2118, 2000.
- [5] Frédéric Ariey, Benoit Witkowski, Chanaki Amaratunga, Johann Beghain, Anne-Claire Langlois, Nimol Khim, Saorin Kim, Valentine Duru, Christiane Bouchier, Laurence Ma, and Others. A molecular marker of artemisinin-resistant Plasmodium falciparum malaria. *Nature*, 505(7481):50–55, 2014.
- [6] James G Baldwin-Brown, Anthony D Long, and Kevin R Thornton. The power to detect quantitative trait loci using resequenced, experimentally evolved populations of diploid, sexual organisms. *Molecular biology and evolution*, page msu048, 2014.
- [7] Rowan D H Barrett, Sean M Rogers, and Dolph Schluter. Natural selection on a major armor gene in threespine stickleback. *Science*, 322(5899):255–257, 2008.
- [8] Jeffrey E Barrick and Richard E Lenski. Genome dynamics during experimental evolution. *Nature Reviews Genetics*, 14(12):827–839, 2013.
- [9] Jeffrey E Barrick, Dong Su Yu, Sung Ho Yoon, Haeyoung Jeong, Tae Kwang Oh, Dominique Schneider, Richard E Lenski, and Jihyun F Kim. Genome evolution and adaptation in a long-term experiment with Escherichia coli. *Nature*, 461(7268):1243–1247, 2009.
- [10] Alan O Bergland, Emily L Behrman, Katherine R O'Brien, Paul S Schmidt, and Dmitri A Petrov. Genomic evidence of rapid and stable adaptive oscillations over seasonal time scales in *Drosophila*. *PLoS Genet*, 10(11):e1004775, 2014.
- [11] Todd Bersaglieri, Pardis C Sabeti, Nick Patterson, Trisha Vanderploeg, Steve F Schaffner, Jared A Drake, Matthew Rhodes, David E Reich, and Joel N Hirschhorn. Genetic signatures of strong recent positive selection at the lactase gene. *The American Journal of Human Genetics*, 74(6):1111–1120, 2004.
- [12] Jonathan P Bollback and John P Huelsenbeck. Clonal interference is alleviated by high mutation rates in large populations. *Molecular biology and evolution*, 24(6):1397–1406, 2007.
- [13] Jonathan P Bollback, Thomas L York, and Rasmus Nielsen. Estimation of 2Nes from temporal allele frequency data. *Genetics*, 179(1):497–502, 2008.
- [14] Adam R Boyko, Scott H Williamson, Amit R Indap, Jeremiah D Degenhardt, Ryan D Hernandez, Kirk E Lohmueller, Mark D Adams, Steffen Schmidt, John J Sninsky, Shamil R Sunyaev, and Others. Assessing the evolutionary impact of amino acid mutations in the human genome. *PLoS Genet*, 4(5):e1000083, 2008.

- [15] Molly K Burke, Joseph P Dunham, Parvin Shahrestani, Kevin R Thornton, Michael R Rose, and Anthony D Long. Genome-wide analysis of a long-term evolution experiment with *Drosophila*. *Nature*, 467(7315):587–590, 2010.
- [16] Molly K Burke, Gianni Liti, and Anthony D Long. Standing genetic variation drives repeatable experimental evolution in outcrossing populations of *Saccharomyces cerevisiae*. *Molecular biology and evolution*, page msu256, 2014.
- [17] Josep M Comeron, Ramesh Ratnappan, and Samuel Bailin. The many landscapes of recombination in *Drosophila melanogaster*. *PLoS Genet*, 8(10):e1002905, 2012.
- [18] P Daborn, S Boundy, J Yen, B Pittendrigh, and Others. DDT resistance in *Drosophila* correlates with Cyp6g1 over-expression and confers cross-resistance to the neonicotinoid imidacloprid. *Molecular Genetics and Genomics*, 266(4):556–563, 2001.
- [19] Rachel Daniels, Hsiao-Han Chang, Papa Diogoye Séne, Danny C Park, Daniel E Neafsey, Stephen F Schaffner, Elizabeth J Hamilton, Amanda K Lukens, Daria Van Tyne, Souleymane Mboup, and Others. Genetic surveillance detects both clonal and epidemic transmission of malaria following enhanced intervention in Senegal. *PLoS One*, 8(4):e60780, 2013.
- [20] Vincent J Denef and Jillian F Banfield. In situ evolutionary rate measurements show ecological success of recently emerged bacterial hybrids. *Science*, 336(6080):462–466, 2012.
- [21] Michael M Desai and Joshua B Plotkin. The polymorphism frequency spectrum of finitely many sites under selection. *Genetics*, 180(4):2175–2191, 2008.
- [22] Richard Durbin, Sean R Eddy, Anders Krogh, and Graeme Mitchison. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge university press, 1998.
- [23] Eyal Elyashiv, Shmuel Sattath, Tina T Hu, Alon Strutsovsky, Graham McVicker, Peter Andolfatto, Graham Coop, and Guy Sella. A Genomic Map of the Effects of Linked Selection in *Drosophila*. *PLoS Genet*, 12(8):1–24, 2016.
- [24] Warren J Ewens. *Mathematical Population Genetics 1: Theoretical Introduction*, volume 27. Springer Science & Business Media, 2012.
- [25] Gregory Ewing and Joachim Hermisson. MSMS: a coalescent simulation program including recombination, demographic structure and selection at a single locus. *Bioinformatics*, 26(16):2064–2065, 2010.
- [26] Shaohua Fan, Matthew E B Hansen, Yancy Lo, and Sarah A Tishkoff. Going global by adapting local: A review of recent human adaptation. *Science*, 354(6308):54–59, 2016.
- [27] Justin C Fay and Chung-I Wu. Hitchhiking under positive Darwinian selection. *Genetics*, 155(3):1405–1413, 2000.
- [28] Alison F Feder, Sergey Kryazhimskiy, and Joshua B Plotkin. Identifying signatures of selection in genetic time series. *Genetics*, 196(2):509–522, 2014.
- [29] Alison F Feder, Soo-Yon Rhee, Susan P Holmes, Robert W Shafer, Dmitri A Petrov, and Pleuni S Pennings. More effective drugs lead to harder selective sweeps in the evolution of drug resistance in HIV-1. *eLife*, 5, jan 2016.

- [30] Anna-Sophie Fiston-Lavier, Nadia D Singh, Mikhail Lipatov, and Dmitri A Petrov. Drosophila melanogaster recombination rate calculator. *Gene*, 463(1):18–20, 2010.
- [31] Susanne U Franssen, Viola Nolte, Ray Tobler, and Christian Schlötterer. Patterns of linkage disequilibrium and long range hitchhiking in evolving experimental Drosophila melanogaster populations. *Molecular biology and evolution*, 32(2):495–509, 2015.
- [32] Nandita R Garud, Philipp W Messer, Erkan O Buzbas, and Dmitri A Petrov. Recent selective sweeps in North American Drosophila melanogaster show signatures of soft sweeps. *PLoS Genet*, 11(2):e1005004, 2015.
- [33] John H Gillespie. *Population genetics: a concise guide*. JHU Press, 2010.
- [34] Michael M Gottesman. Mechanisms of cancer drug resistance. *Annual review of medicine*, 53(1):615–627, 2002.
- [35] Torsten Günther and Graham Coop. Robust identification of local adaptation from allele frequencies. *Genetics*, 195(1):205–220, 2013.
- [36] Matthew Hegeness, Noam Shresh, Daniel Hartl, and Roy Kishony. An equivalence principle for the incorporation of favorable mutations in asexual populations. *Science*, 311(5767):1615–1617, 2006.
- [37] Kent E Holsinger and Bruce S Weir. Genetics in geographically structured populations: defining, estimating and interpreting FST. *Nature Reviews Genetics*, 10(9):639–650, 2009.
- [38] Christopher J R Illingworth and Ville Mustonen. Distinguishing driver and passenger mutations in an evolutionary history categorized by interference. *Genetics*, 189(3):989–1000, 2011.
- [39] Christopher J R Illingworth, Leopold Parts, Stephan Schiffels, Gianni Liti, and Ville Mustonen. Quantifying selection acting on a complex trait using allele frequency time series data. *Molecular biology and evolution*, 29(4):1187–1197, 2012.
- [40] Minako Izutsu, Atsushi Toyoda, Asao Fujiyama, Kiyokazu Agata, and Naoyuki Fuse. Dynamics of Dark-Fly Genome Under Environmental Selections. *G3: Genes—Genomes—Genetics*, pages g3—115, 2015.
- [41] Aashish R Jha, Cecelia M Miles, Nodia R Lippert, Christopher D Brown, Kevin P White, and Martin Kreitman. Whole-genome resequencing of experimental populations reveals polygenic basis of egg-size variation in Drosophila melanogaster. *Molecular biology and evolution*, 32(10):2616–2632, 2015.
- [42] Agnes Jónás, Thomas Taus, Carolin Kosiol, Christian Schlötterer, and Andreas Futschik. Estimating the Effective Population Size from Temporal Allele Frequency Changes in Experimental Evolution. *Genetics*, aug 2016.
- [43] Tadeusz J Kawecki, Richard E Lenski, Dieter Ebert, Brian Hollis, Isabelle Olivieri, and Michael C Whitlock. Experimental evolution. *Trends in ecology & evolution*, 27(10):547–560, 2012.
- [44] Robert Kofler and Christian Schlötterer. Gowinda: unbiased analysis of gene set enrichment for genome-wide association studies. *Bioinformatics*, 28(15):2084–2085, 2012.

- [45] Robert Kofler and Christian Schlötterer. A guide for the design of evolve and resequencing studies. *Molecular biology and evolution*, page mst221, 2013.
- [46] Gregory I Lang, David Botstein, and Michael M Desai. Genetic variation and the fate of beneficial mutations in asexual populations. *Genetics*, 188(3):647–661, 2011.
- [47] Gregory I Lang, Daniel P Rice, Mark J Hickman, Erica Sodergren, George M Weinstock, David Botstein, and Michael M Desai. Pervasive genetic hitchhiking and clonal interference in forty evolving yeast populations. *Nature*, 500(7464):571–574, 2013.
- [48] Quan Long, Fernando A Rabanal, Dazhe Meng, Christian D Huber, Ashley Farlow, Alexander Platzer, Qingrun Zhang, Bjarni J Vilhjálmsson, Arthur Korte, Viktoria Nizhynska, and Others. Massive genomic variation and strong selection in *Arabidopsis thaliana* lines from Sweden. *Nature genetics*, 45(8):884–890, 2013.
- [49] Michael Lynch, Darius Bost, Sade Wilson, Takahiro Maruki, and Scott Harrison. Population-genetic inference from pooled-sequencing data. *Genome biology and evolution*, 6(5):1210–1218, 2014.
- [50] Anna-Sapfo Malaspinas, Orestis Malaspinas, Steven N Evans, and Montgomery Slatkin. Estimating allele age and selection coefficient from time-serial data. *Genetics*, 192(2):599–607, 2012.
- [51] Frank Maldarelli, Mary Kearney, Sarah Palmer, Robert Stephens, JoAnn Mican, Michael A Polis, Richard T Davey, Joseph Kovacs, Wei Shao, Diane Rock-Kress, and Others. HIV populations are large and accumulate high genetic diversity in a nonlinear fashion. *Journal of virology*, 87(18):10313–10323, 2013.
- [52] Nelson E Martins, Vítor G Faria, Viola Nolte, Christian Schlötterer, Luis Teixeira, Élio Sucena, and Sara Magalhães. Host adaptation to viruses relies on few genes with different cross-resistance properties. *Proceedings of the National Academy of Sciences*, 111(16):5938–5943, 2014.
- [53] Iain Mathieson and Gil McVean. Estimating selection coefficients in spatially structured populations from time series data of allele frequencies. *Genetics*, 193(3):973–984, 2013.
- [54] Philipp W Messer and Dmitri A Petrov. Population genomics of rapid adaptation by soft selective sweeps. *Trends in ecology & evolution*, 28(11):659–669, 2013.
- [55] Shalini Nair, Denae Nash, Daniel Sudimack, Anchalee Jaidee, Marion Barends, Anne-Catrin Uhlemann, Sanjeev Krishna, François Nosten, and Tim J C Anderson. Recurrent gene amplification and soft selective sweeps during evolution of multidrug resistance in malaria parasites. *Molecular Biology and Evolution*, 24(2):562–573, 2007.
- [56] Rasmus Nielsen and James Signorovitch. Correcting for ascertainment biases when analyzing SNP data: applications to the estimation of linkage disequilibrium. *Theoretical population biology*, 63(3):245–255, 2003.
- [57] Rasmus Nielsen, Scott Williamson, Yuseob Kim, Melissa J Hubisz, Andrew G Clark, and Carlos Bustamante. Genomic scans for selective sweeps using SNP data. *Genome research*, 15(11):1566–1575, 2005.

- [58] Pablo Orozco-ter Wengel, Martin Kapun, Viola Nolte, Robert Kofler, Thomas Flatt, and Christian Schlötterer. Adaptation of *Drosophila* to a novel laboratory environment reveals temporally heterogeneous trajectories of selected alleles. *Molecular ecology*, 21(20):4931–4941, 2012.
- [59] Tugce Oz, Aysegul Guvenek, Sadik Yildiz, Enes Karaboga, Yusuf Talha Tamer, Nirva Mumcuyan, Vedat Burak Ozan, Gizem Hazal Senturk, Murat Cokol, Pamela Yeh, and Others. Strength of selection pressure is an important parameter contributing to the complexity of antibiotic resistance evolution. *Molecular biology and evolution*, page msu191, 2014.
- [60] Bo Peng and Marek Kimmel. simuPOP: a forward-time population genetics simulation environment. *Bioinformatics*, 21(18):3686–3687, 2005.
- [61] Edward Pollak. A new method for estimating the effective population size from allele frequency changes. *Genetics*, 104(3):531–548, 1983.
- [62] Susan E Ptak and Molly Przeworski. Evidence for population growth in humans is confounded by fine-scale population structure. *Trends in Genetics*, 18(11):559–563, 2002.
- [63] Sebastian E Ramos-Onsins and Julio Rozas. Statistical properties of new neutrality tests against population growth. *Molecular biology and evolution*, 19(12):2092–2100, 2002.
- [64] Brian J Reid, Rumen Kostadinov, and Carlo C Maley. New strategies in Barrett’s esophagus: integrating clonal evolutionary theory with clinical management. *Clinical Cancer Research*, 17(11):3512–3519, 2011.
- [65] Silvia C Remolina, Peter L Chang, Jeff Leips, Sergey V Nuzhdin, and Kimberly A Hughes. Genomic basis of aging and life-history evolution in *Drosophila melanogaster*. *Evolution*, 66(11):3390–3403, 2012.
- [66] Roy Ronen, Nitin Udpa, Eran Halperin, and Vineet Bafna. Learning natural selection from the site frequency spectrum. *Genetics*, 195(1):181–193, 2013.
- [67] P C Sabeti, S F Schaffner, B Fry, J Lohmueller, P Varilly, O Shamovsky, A Palma, T S Mikkelsen, D Altshuler, and E S Lander. Positive natural selection in the human lineage. *science*, 312(5780):1614–1620, 2006.
- [68] Stanley A Sawyer and Daniel L Hartl. Population genetics of polymorphism and divergence. *Genetics*, 132(4):1161–1176, 1992.
- [69] Christian Schlötterer, R Kofler, E Versace, R Tobler, and S U Franssen. Combining experimental evolution with next-generation sequencing: a powerful tool to study adaptation from standing genetic variation. *Heredity*, 114(5):431–440, 2015.
- [70] Tatum S Simonson, Yingzhong Yang, Chad D Huff, Haixia Yun, Ga Qin, David J Witherspoon, Zhenzhong Bai, Felipe R Lorenzo, Jinchuan Xing, Lynn B Jorde, and Others. Genetic evidence for high-altitude adaptation in Tibet. *Science*, 329(5987):72–75, 2010.
- [71] Brad Spellberg, Robert Guidos, David Gilbert, John Bradley, Helen W Boucher, W Michael Scheld, John G Bartlett, John Edwards, Infectious Diseases Society of America, and Others. The epidemic of antibiotic-resistant infections: a call to action for the medical community from the Infectious Diseases Society of America. *Clinical Infectious Diseases*, 46(2):155–164, 2008.

- [72] Matthias Steinrücken, Anand Bhaskar, and Yun S Song. A novel spectral method for inferring general diploid selection from time series genetic data. *The annals of applied statistics*, 8(4):2203, 2014.
- [73] Wolfgang Stephan, Yun S Song, and Charles H Langley. The hitchhiking effect on linkage disequilibrium between linked neutral loci. *Genetics*, 172(4):2647–2663, 2006.
- [74] Fumio Tajima. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics*, 123(3):585–595, 1989.
- [75] Jonathan Terhorst, Christian Schlötterer, and Yun S Song. Multi-locus Analysis of Genomic Time Series Data from Experimental Evolution. *PLoS Genet*, 11(4):e1005069, 2015.
- [76] Ray Tobler, Susanne U Franssen, Robert Kofler, Pablo Orozco-terWengel, Viola Nolte, Joachim Hermissen, and Christian Schlötterer. Massive habitat-specific genomic response in *D. melanogaster* populations during experimental evolution in hot and cold environments. *Molecular biology and evolution*, 31(2):364–375, 2014.
- [77] Hande Topa, Ágnes Jónás, Robert Kofler, Carolin Kosiol, and Antti Honkela. Gaussian process test for high-throughput sequencing time series: application to experimental evolution. *Bioinformatics*, page btv014, 2015.
- [78] Thomas L Turner, Andrew D Stewart, Andrew T Fields, William R Rice, and Aaron M Tarone. Population-based resequencing of experimentally evolved populations reveals the genetic basis of body size variation in *Drosophila melanogaster*. *PLoS Genet*, 7(3):e1001336, 2011.
- [79] Joseph J Vitti, Sharon R Grossman, and Pardis C Sabeti. Detecting natural selection in genomic data. *Annual review of genetics*, 47:97–120, 2013.
- [80] Jinliang Wang. A pseudo-likelihood method for estimating effective population size from temporally spaced samples. *Genetical research*, 78(03):243–257, 2001.
- [81] Robin S Waples. A generalized approach for estimating effective population size from temporal changes in allele frequency. *Genetics*, 121(2):379–391, 1989.
- [82] David Williams and David Williams. *Weighing the odds: a course in probability and statistics*, volume 548. Springer, 2001.
- [83] Ellen G Williamson and Montgomery Slatkin. Using maximum likelihood to estimate population size from temporal changes in allele frequencies. *Genetics*, 152(2):755–761, 1999.
- [84] Scott H Williamson, Melissa J Hubisz, Andrew G Clark, Bret A Payseur, Carlos D Bustamante, and Rasmus Nielsen. Localizing recent adaptive evolution in the human genome. *PLoS Genet*, 3(6):e90, 2007.
- [85] Mark A Winters, Robert M Lloyd Jr, Robert W Shafer, Michael J Kozal, Michael D Miller, and Mark Holodniy. Development of elvitegravir resistance and linkage of integrase inhibitor mutations with protease and reverse transcriptase resistance mutations. *PloS one*, 7(7):e40514, 2012.
- [86] Xin Yi, Yu Liang, Emilia Huerta-Sánchez, Xin Jin, Zha Xi Ping Cuo, John E Pool, Xun Xu, Hui Jiang, Nicolas Vinckenbosch, Thorfinn Sand Korneliussen, and Others. Sequencing of 50 human exomes reveals adaptation to high altitude. *Science*, 329(5987):75–78, 2010.

- [87] Hiba Zahreddine and K L Borden. Mechanisms and insights into drug resistance in cancer. *Front Pharmacol*, 4(28.10):3389, 2013.
- [88] Dan Zhou, Nitin Udpa, Merril Gersten, DeeAnn W Visk, Ali Bashir, Jin Xue, Kelly A Frazer, James W Posakony, Shankar Subramaniam, Vineet Bafna, and Gabriel G. Haddad. Experimental selection of hypoxia-tolerant *Drosophila melanogaster*. *Proceedings of the National Academy of Sciences*, 108(6):2349–2354, 2011.