

CLEAR: Composition of Likelihoods for Evolve And Resequencing 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

The advent of next generation sequencing technologies has made whole-genome and whole-population sampling possible, even for eukaryotes with large genomes. With this development, experimental evolution studies can be designed to observe molecular evolution “in-action” via Evolve-and-Resequencing (E&R) experiments. Among other applications, E&R studies can be used to locate the genes and variants responsible for genetic adaptation. Existing literature on time-series data analysis often assumes large population size, accurate allele frequency estimates, and wide time spans. These assumptions do not hold in many E&R studies.

In this article, we propose a method—Composition of Likelihoods for Evolve-And-Resequencing experiments (CLEAR)—to identify signatures of selection in small population E&R experiments. CLEAR takes whole-genome sequence of pool of individuals (pool-seq) as input, and properly addresses heterogeneous ascertainment bias resulting from uneven coverage. CLEAR also provides unbiased estimates of model parameters, including population size, selection strength and dominance, while being computationally efficient. Extensive simulations show that CLEAR achieves higher power in detecting and localizing selection over a wide range of parameters, and is robust to variation of coverage. We applied CLEAR statistic to multiple E&R experiments, including, data from a study of *D. melanogaster* adaptation to alternating temperatures and a study of outcrossing yeast populations, and identified multiple regions under selection with genome-wide significance.

1 Introduction

Natural selection is a key force in evolution, and a mechanism by which populations can adapt to external ‘selection’ pressure. Examples of adaptation abound in the natural world [22], including for example, classic examples like lactose tolerance in Northern Europeans [9], human adaptation to high altitudes [55, 69], but also drug resistance in pests [15], HIV [24], cancer [27, 70], malarial parasite [3, 44], and others [56]. In these examples, understanding the genetic basis of adaptation can provide valuable information, underscoring the importance of the problem.

Experimental evolution refers to the study of the evolutionary processes of a model organism in a controlled [7, 10, 28, 37, 38, 46, 47] or natural [5, 8, 16, 17, 41, 50, 68] environment. Recent advances in whole genome sequencing have enabled us to sequence populations at a reasonable cost, even for large genomes. Perhaps more important for experimental evolution studies, we can now evolve and resequence (E&R) multiple replicates of a population to obtain *longitudinal time-series data*, in order to investigate the dynamics of evolution at molecular level. Although constraints such as small sizes, limited timescales, and oversimplified laboratory environments may

34 limit the interpretation of E&R results, these studies are increasingly being used to test a wide
35 range of hypotheses [34] and have been shown to be more predictive than static data analysis
36 [12, 18, 52]. In particular, longitudinal E&R data is being used to estimate model parameters
37 including population size [33, 49, 60, 64, 65, 67], strength of selection [11, 29, 30, 40, 43, 57, 60],
38 allele age [40] recombination rate [60], mutation rate [6, 60], quantitative trait loci [4] and for tests
39 of neutrality hypotheses [8, 13, 23, 60].

40 While many E&R study designs are being used [6, 53], we restrict our attention to the adaptive
41 evolution due to standing variation in fixed size populations. This regime has been considered
42 earlier, typically with *D. melanogaster* as the model organism of choice, to identify adaptive genes
43 in longevity and aging [13, 51] (600 generations), courtship song [63] (100 generations), hypoxia
44 tolerance [71] (200 generations), adaptation to new laboratory environments [26, 46] (59 genera-
45 tions), egg size [32] (40 generations), C virus resistance [42] (20 generations), and dark-fly [31] (49
46 generations).

47 The task of identifying selection signatures can be addressed at different levels of specificity.
48 At the coarsest level, identification could simply refer to deciding whether some genomic region (or
49 a gene) is under selection or not. In the following, we refer to this task as *detection*. In contrast,
50 the task of *site-identification* corresponds to the process of finding the favored mutation/allele
51 at nucleotide level. Finally, *estimation of model parameters*, such as strength of selection and
52 dominance at the site, can provide a comprehensive description of the selection process.

53 In the effort to analyze E&R selection experiments, many authors chose to adapt existing
54 tests that were originally used for static data, pairwise comparisons (two time-points) and single
55 replicates to perform a null scan. For instance, Zhu *et al.* [71] used the ratio of the estimated
56 population size of case and control populations to compute test statistic for each genomic region.
57 Burke *et al.* [13] applied Fisher exact test to the last observation of data on case and control
58 populations. Orozco-terWengel *et al.* [46] used the Cochran-Mantel-Haenszel (CMH) test [1] to
59 detect SNPs whose read counts change consistently across all replicates of two time-point data.
60 Turner *et al.* [63] proposed the diffStat statistic to test whether the change in allele frequencies
61 of two populations deviate from the distribution of change in allele frequencies of two drifting
62 populations. Bergland *et al.* [8] calculated F_{st} to populations throughout time to signify their
63 differentiation from ancestral (two time-point data) as well as geographically different populations.
64 Jha *et al.* [32] computed test statistic of generalized linear-mixed model directly from read counts.

65 Alternatively, *direct* methods have been developed to analyze time-series data by taking a
66 likelihood approach, and estimating population genetics parameters. Bollback *et al.* [11] proposed
67 a Hidden Markov Model (HMM) to estimate the selection coefficient s and population size by
68 using a diffusion approximation to the continuous Wright Fisher Markov process. Steinrücken and
69 Song [57] proposed a general diploid selection model which takes into account of dominance of
70 the favored allele and approximates likelihood analytically. Mathieson and McVean [43] adopted
71 HMMs to structured populations and estimated parameters using an Expectation Maximization
72 (EM) procedure on discretized allele frequency. Feder *et al.* [23] modeled increments in allele
73 frequency with a Brownian motion process, proposed the Frequency Increment Test (FIT). More
74 recently, Topa *et al.* [62] proposed a Gaussian Process (GP) for modeling single-locus time-series
75 pool-seq data. Terhorst *et al.* [60] extended GP to compute joint likelihood of multiple loci under
76 null and alternative hypotheses. Recently, Schraiber *et al.* [54] proposed a Bayesian framework to
77 estimate parameters using Monte Carlo Markov chain sampling.

78 While existing methods have been successfully applied to their corresponding application, they
79 make some assumptions which may not hold in E&R studies. First, they assume that the underlying
80 population size is large, so it is reasonable to model dynamics of allele frequencies using continuous
81 state models. A number of existing methods were originally designed to process wide time spans

82 such as ancient DNA studies. Finally, they assume that input data is in the form of unbiased allele
83 frequencies, which may not be valid for shotgun sequencing experiments.

84 Here, we consider a Hidden Markov Model (HMM), similar to Williamson *et al.* [67] and Boll-
85 back *et al.*'s [11] but under a “small-population-size” regime. Specifically, we use a discrete state
86 (frequency) model. We show that for small population sizes, discrete models can compute likeli-
87 hood exactly, which improves statistical performance, especially for short time-span experiments.
88 Additionally, we add another level of sampling-noise to the traditional HMM model, allowing for
89 heterogeneous ascertainment bias due to uneven coverage among variants. We show that for a wide
90 range of parameters, CLEAR provides higher power for detecting selection, estimates model pa-
91 rameters consistently, and localizes favored allele more accurately compared to the state-of-the-art
92 methods, while being computationally efficient.

93 2 Materials and Methods

94 Consider a panmictic diploid population with fixed size of N individuals. Let $\nu = \{\nu_t\}_{t \in \mathcal{T}}$ be
95 frequencies of the derived allele at generations $t \in \mathcal{T}$ for a given variant, where at generations
96 $\mathcal{T} = \{\tau_i : 0 \leq \tau_0 < \tau_1, \dots < \tau_T\}$ samples of n individuals are chosen for pooled sequencing. The
97 experiment is replicated R times. We denote allele frequencies of the R replicates by the set $\{\nu\}_R$.
98 To identify the genes and variants that are responding to selection pressure, we use the following
99 procedure:

- 100 (i) **Estimating population size.** The procedure starts by estimating the effective population
101 size, \hat{N} , under the assumption that much of the genome is evolving neutrally.
- 102 (ii) **Estimating selection parameters.** For each polymorphic site, selection and dominance
103 parameters s, h are estimated so as to maximize the likelihood of the time series data, given
104 \hat{N} .
- 105 (iii) **Computing likelihood statistics.** For each variant, a log-odds ratio of the likelihood
106 of selection model ($s > 0$) to the likelihood of neutral evolution/drift model is computed.
107 Likelihood ratios in a genomic region are combined to compute the CLEAR statistic for the
108 region.
- 109 (iv) **Hypothesis testing.** An empirical null distribution of the CLEAR statistic is calculated using
110 genome-wide drift simulations, and used to compute p -values and thresholds for a specified
111 FDR. We perform single locus hypothesis testing within selected regions to identify significant
112 variants and report genes that intersect with the selected variants.

113 These steps are described in detail below.

114 2.1 Estimating Population Size

115 Methods for estimating population sizes from temporal neutral evolution data have been devel-
116 oped [2, 11, 33, 60, 67]. Here, we aim to extend these models to explicitly model the sampling noise
117 that arise in pool-seq data. Specifically, we model the variation in sequence coverage over different
118 locations, and the noise due to sequencing only a subset of the individuals in the population. In
119 addition, many existing methods [11, 23, 60, 62] are designed for large populations, and model
120 frequency as a continuous quantity. We show that smooth approximations may be inadequate for
121 small populations, low starting frequencies and sparse sampling (in time) that are typical in ex-
122 perimental evolution (see Results, Fig 3A-C, and Fig 2). To this end, we model the Wright-Fisher

123 Markov process for generating pool-seq data (Fig S1) via a discrete HMM (Fig 1-B). We start by
 124 computing a likelihood function for the population size given neutral pool-seq data.

125 **Likelihood for Neutral Model.** We model the allele frequency counts $2N\nu_t$ as being sampled
 126 from a Binomial distribution. Specifically,

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

127 where π is the global distribution of allele frequencies in the base population. Here we simply
 128 assume π is the site frequency spectrum of fixed sized neutral population Fig S2. Note that π
 129 may depend on the demographic history of the founder lines.

130 To estimate frequency after τ transitions, it is enough to specify the $2N \times 2N$ transition matrix
 131 $P^{(\tau)}$, where $P^{(\tau)}[i, j]$ denotes probability of change in allele frequency from $i/2N$ to $j/2N$ in τ
 132 generations:

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

$$P^{(\tau)} = P^{(\tau-1)}P^{(1)} \quad (2)$$

133 Furthermore, in an E&R experiment, $n \leq N$ individuals are randomly selected for sequencing. The
 134 sampled allele frequencies, $\{y_t\}_{t \in \mathcal{T}}$, are also Binomially distributed

$$2ny_t \sim \text{Binomial}(2n, \nu_t) \quad (3)$$

135 We introduce the $2N \times 2n$ sampling matrix Y , where $Y[i, j]$ stores the probability that the sample
 136 allele frequency is $j/2n$ given that the true allele frequency is $i/2N$.

137 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
 138 coverage, and the read count of the derived allele, respectively. Let $\{\lambda_t\}_{t \in \mathcal{T}}$ be the sequencing
 139 coverage at different generations. Then, the observed data are sampled according to

$$d_t \sim \text{Poisson}(\lambda_t), \quad c_t \sim \text{Binomial}(d_t, y_t) \quad (4)$$

140 The emission probability for a observed tuple $x_t = \langle d_t, c_t \rangle$ is

$$\mathbf{e}_i(x_t) = \binom{d_t}{c_t} \left(\frac{i}{2n}\right)^{c_t} \left(1 - \frac{i}{2n}\right)^{d_t - c_t}. \quad (5)$$

For $1 \leq t \leq T, 1 \leq j \leq 2N$, let $\alpha_{t,j}$ denote the probability of emitting x_1, x_2, \dots, x_t and reaching
 state j at τ_t . Then, α_t can be computed using the forward-procedure [19]:

$$\alpha_t^T = \alpha_{t-1}^T P^{(\delta_t)} \text{diag}(Y \mathbf{e}(x_t)) \quad (6)$$

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

$$\mathcal{L}(N|\{\mathbf{x}\}_R, n) = \prod_{r=1}^R \mathcal{L}(N|\mathbf{x}^{(r)}, n) = \Pr(\{\mathbf{x}\}_R|N, n) = \prod_{r=1}^R \sum_i \alpha_{T,i}^{(r)} \quad (7)$$

143 where $\mathbf{x} = \{x_t\}_{t \in \mathcal{T}}$. The graphical model and the generative process for which data is being
 144 generated is depicted in Fig 1-B and Fig S1, respectively.

145 Finally, the last step is to compute an estimate \hat{N} that maximizes the likelihood of all M
 146 variants in whole genome. Let $\mathbf{x}_i^{(r)}$ denote the time-series data of the i -th variant in replicate r .
 147 Then,

$$\hat{N} = \arg \max_N \prod_{i=1}^M \prod_{r=1}^R \mathcal{L}(N|\mathbf{x}_i^{(r)}) \quad (8)$$

148 2.2 Estimating Selection Parameters

149 **Likelihood for Selection Model.** Assume that the site is evolving under selection constraints
 150 $s \in \mathbb{R}$, $h \in \mathbb{R}_+$, where s and h denote selection strength and dominance parameters, respectively.
 151 By definition, the relative fitness values of genotypes 0|0, 0|1 and 1|1 are given by $w_{00} = 1$,
 152 $w_{01} = 1 + hs$ and $w_{11} = 1 + s$. Then, ν_{t+} , the frequency at time $\tau_t + 1$ (one generation ahead), can
 153 be estimated using:

$$\begin{aligned} \hat{\nu}_{t+} = \mathbb{E}[\nu_{t+} | s, h, \nu_t] &= \frac{w_{11}\nu_t^2 + w_{01}\nu_t(1 - \nu_t)}{w_{11}\nu_t^2 + 2w_{01}\nu_t(1 - \nu_t) + w_{00}(1 - \nu_t)^2} \\ &= \nu_t + \frac{s(h + (1 - 2h)\nu_t)\nu_t(1 - \nu_t)}{1 + s\nu_t(2h + (1 - 2h)\nu_t)}. \end{aligned} \quad (9)$$

154 The machinery for computing likelihood of the selection parameters is identical to that of population
 155 size, except for transition matrices. Hence, here we only describe the definition transition matrix
 156 $Q_{s,h}$ of the selection model. Let $Q_{s,h}^{(\tau)}[i, j]$ denote the probability of transition from $i/2N$ to $j/2N$
 157 in τ generations, then (See [20], Pg. 24, Eqn. 1.58-1.59):

$$Q_{s,h}^{(1)}[i, j] = \Pr\left(\nu_{t+} = \frac{j}{2N} \mid \nu_t = \frac{i}{2N}; s, h, N\right) = \binom{2N}{j} \hat{\nu}_{t+}^j (1 - \hat{\nu}_{t+})^{2N-j} \quad (10)$$

$$Q_{s,h}^{(\tau)} = Q_{s,h}^{(\tau-1)} Q_{s,h}^{(1)} \quad (11)$$

158 The maximum likelihood estimates are given by

$$\hat{s}, \hat{h} = \arg \max_{s,h} \prod_{r=1}^R \mathcal{L}(s, h | \mathbf{x}^{(r)}, \hat{N}) \quad (12)$$

159 Using grid search, we first estimate N (Eq. 8), and subsequently, we estimate parameters s, h
 160 (Eq. 12, Fig S3). By broadcasting and vectorizing the grid search operations across all variants, the
 161 genome scan on millions of polymorphisms can be done in significantly smaller time than iterating
 162 a numerical optimization routine for each variant (see Results and Fig 4).

163 2.3 Empirical Likelihood Ratio Statistics

164 The likelihood ratio statistic for testing directional selection, to be computed for each variant, is
 165 given by

$$H = -2 \log \left(\frac{\mathcal{L}(\bar{s}, 0.5 | \{\mathbf{x}\}_R, \hat{N})}{\mathcal{L}(0, 0.5 | \{\mathbf{x}\}_R, \hat{N})} \right), \quad (13)$$

166 where $\bar{s} = \arg \max_s \prod_{r=1}^R \mathcal{L}(s, 0.5 | \mathbf{x}^{(r)}, \hat{N})$. Similarly we can define a test statistic for testing if
 167 selection is dominant by

$$D = -2 \log \left(\frac{\mathcal{L}(\hat{s}, \hat{h} | \{\mathbf{x}\}_R, \hat{N})}{\mathcal{L}(\bar{s}, 0.5 | \{\mathbf{x}\}_R, \hat{N})} \right). \quad (14)$$

168 While extending the single-locus WF model to a multiple linked-loci can improve the power of
 169 the model [60], it is computationally and statistically expensive to compute exact likelihood. In
 170 addition, computing linked-loci joint likelihood requires haplotype resolved data, which pool-seq

171 does not provide. Here, similar to Nielsen *et al* [45], we calculate *composite likelihood ratio* score
172 for a genomic region.

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

173 where L is a collection of segregating sites and H_{ℓ} is the likelihood ratio score based for each
174 variant ℓ in L . The optimal value of the hyper-parameter L depends upon a number of factors,
175 including initial frequency of the favored allele, recombination rates, linkage of the favored allele
176 to neighboring variants, population size, coverage, and time since the onset of selection (duration
177 of the experiment). In [S1 Text](#), we provide a heuristic to compute a reasonable value of L , based
178 on experimental data.

179 We work with a normalized value of \mathcal{H} , given by

$$\mathcal{H}_i^* = \frac{\mathcal{H}_i - \mu_{\mathcal{C}}}{\sigma_{\mathcal{C}}}, \quad \forall i \in \mathcal{C}, \quad (16)$$

180 where $\mu_{\mathcal{C}}$ and $\sigma_{\mathcal{C}}$ are the mean and standard deviation of \mathcal{H} values in a large region \mathcal{C} . We found
181 different chromosomes to have different distribution of \mathcal{H}_i values, and therefore decided to use single
182 chromosomes as \mathcal{C} .

183 2.4 Hypothesis Testing

184 **Single-Locus tests.** Under neutrality, Log-likelihood ratios can be approximated by χ^2 distri-
185 bution [66], and p -values can be computed directly. However, Feder *et al.* [23] showed that when
186 the number of independent samples (replicates) is small, χ^2 is a crude approximation to the true
187 null distribution and results in more false positive. Following their suggestion, we first compute the
188 empirical null distribution using simulations with the estimated population size (See [Fig S1](#)). The
189 empirical null distribution of statistic H is used to compute p -values as the fraction of null values
190 that exceed the test score. Finally, we use Storey and Tibshirani’s method [59] to control for False
191 Discovery Rate in multiple testing.

192 **Composite likelihood tests.** Similar to single-locus tests, we compute the null distribution of the
193 \mathcal{H}^* statistic using whole-genome simulations with the estimated population size, and subsequently
194 compute FDR. The simulations for generating the null distribution of \mathcal{H}^* are described next.

195 2.5 Simulations

196 We use the same simulation procedure for two purposes. First, we use them to test the power
197 of CLEAR against other methods in small genomic windows. Second, we use the simulations to
198 generate the distribution of null values for the statistic to compute empirical p -values. We mainly
199 chose parameters that are relevant to *D. melanogaster* experimental evolution [35]. See also [Fig 1-A](#)
200 for illustration.

201 **I. Creating initial founder line haplotypes.** Using `msms` [21], we created neutral popu-
202 lations for F founding haplotypes with command `$. /msms <F> 1 -t <2μWNe> -r <2rNeW>`
203 `<W>`, where $F = 200$ is number of founder lines, $N_o = 10^6$ is effective founder population size,
204 $r = 2 \times 10^{-8}$ is recombination rate, $\mu = 2 \times 10^{-9}$ is mutation rate. The window size W is
205 used to compute $\theta = 2\mu N_o W$ and $\rho = 2N_o r W$. We chose $W = 50\text{Kbp}$ for simulating indi-
206 vidual windows for performance evaluations, and $W = 20\text{Mbp}$ for simulating *D. melanogaster*
207 chromosomes for p -value computations.

208 **II. Creating initial diploid population.** An initial set of $F = 200$ haplotypes was created
209 from step I, and duplicated to create F homozygous diploid individuals to simulate generation
210 of inbred lines. N diploid individuals were generated by sampling with replacement from the
211 F individuals.

212 **III. Forward Simulation.** We used forward simulations for evolving populations under selection.
213 We also consider selection regimes which the favored allele is chosen from standing variation
214 (not *de novo* mutations). Given initial diploid population, position of the site under selection,
215 selection strength s , number of replicates $R = 3$, recombination rate $r = 2 \times 10^{-8}$ and
216 sampling times $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$, `simuPop` [48] was used to perform forward simulation
217 and compute allele frequencies for all of the R replicates. For hard sweep (respectively, soft
218 sweep) simulations we randomly chose a site with initial frequency of $\nu_0 = 0.005$ (respectively,
219 $\nu_0 = 0.1$) to be the favored allele. For generating the null distribution with drift for p -value
220 computations, we used this procedure with $s = 0$.

221 **IV. Sequencing Simulation.** Given allele frequency trajectories we sampled depth of each site
222 in each replicate identically and independently from $\text{Poisson}(\lambda)$, where $\lambda \in \{30, 100, 300\}$ is
223 the coverage for the experiment. Once depth d is drawn for the site with frequency ν , the
224 number of reads c carrying the derived allele are sampled according to $\text{Binomial}(d, \nu)$. For
225 experiments with finite depth the tuple $\langle c, d \rangle$ is the input data for each site.

226 3 Results

227 **Modeling Allele Frequency Trajectories in Small Populations.** We first tested the goodness
228 of fit of the discrete versus continuous models in modeling allele frequency trajectories, under
229 general E&R parameters. For this purpose, we conducted 100K simulations with two time samples
230 $\mathcal{T} = \{0, \tau\}$ where $\tau \in \{1, 10, 100\}$ is the parameter controlling the density of sampling in time. In
231 addition, we repeated simulations for different values of starting frequency $\nu_0 \in \{0.005, 0.1\}$ (i.e.,
232 hard and soft sweep) and selection strength $s \in \{0, 0.1\}$ (i.e., neutral and selection). Then, given
233 initial frequency ν_0 , we computed the expected distribution of the frequency of the next sample
234 ν_τ under two models to make a comparison. Fig 2A-F shows that Brownian motion (continuous
235 model) is inadequate when ν_0 is far from 0.5, or when sampling times are sparse ($\tau > 1$). If the
236 favored allele arises from standing variation in a neutral population, it is unlikely to have frequency
237 close to 0.5, and the starting frequencies are usually much smaller (see Fig S2). Moreover, in typical
238 *D. melanogaster* experiments for example, sampling is sparse. Often, the experiment is designed
239 so that $10 \leq \tau \leq 100$ [26, 35, 46, 71].

240 In contrast to the Brownian motion approximation, discrete Markov chain predictions (Eq. 11)
241 are highly consistent with empirical data for a wide range of simulation parameters (Fig 2A-M).
242 Moreover, the discrete markov chain can be modified to model the case when the the allele is under
243 selection.

244 **Detection Power.** We compared the performance of CLEAR against other methods for detect-
245 ing selection. For each method we calculated detection power as the percentage of true-positives
246 identified with false-positive rate ≤ 0.05 . For each configuration (specified with values for selection
247 coefficient s , starting allele frequency ν_0 and coverage λ), power of each method is evaluated over
248 2000 distinct simulations, half of which modeled neutral evolution and the rest modeled positive
249 selection.

250 We compared the power of CLEAR with Gaussian process (GP) [60], FIT [23], and CMH [1]
251 statistics. FIT and GP convert read counts to allele frequencies prior to computing the test statistic.

252 CLEAR shows the highest power in all cases and the power stays relatively high even for low coverage
253 (Fig 3 and Table S1). In particular, the difference in performance of CLEAR with other methods
254 is pronounced when starting frequency is low. The advantage of CLEAR stems from the fact that
255 favored allele with low starting frequency might be missed by low coverage sequencing. In this
256 case, incorporating the signal from linked sites becomes increasingly important. We note that
257 methods using only two time points, such as CMH, do relatively well for high selection values and
258 high coverage. However, the use of time-series data can increase detection power in low coverage
259 experiments or when starting frequency is low. Moreover, time-series data provide means for
260 estimating selection parameters s, h (see below). Finally, as CLEAR is robust to change of coverage,
261 our results (Fig 3B,C) suggest that taking many samples with lower coverage is preferable to sparse
262 sampling with higher coverage.

263 **Site-identification.** In general, localizing the favored variant, using pool-seq data is a nontrivial
264 task due to extensive linkage disequilibrium [61]. To measure performance, we sorted variants by
265 their H scores and computed rank of the favored allele for each method. For each setting of ν_0
266 and s , we conducted 1000 simulations and computed the rank of the favored mutation in each
267 simulation. The cumulative distribution of the rank of the favored allele in 1000 simulation for
268 each setting (Fig 5) shows that CLEAR outperforms other statistics.

269 An interesting observation is revisiting the contrast between site-identification and detection [39,
270 61]. When selection strength is high, detection is easier (Fig 3A-F), but site-identification is harder,
271 due to the high LD between flanking variants and the favored allele (Fig 5A-F). Moreover, site-
272 identification becomes more difficult whenever the initial frequency of the favored allele is low, i.e.,
273 at the onset of selection, LD between favored allele and its nearby variants is high. For example,
274 when coverage $\lambda = 100$ and selection coefficient $s = 0.1$, the detection power is 75% for hard sweep,
275 but 100% for soft sweep (Fig 3B-E). In contrast, the favored site was ranked as the top in 14% of
276 hard sweep cases, compared to and 95% of soft sweep simulations.

277 **Estimating Parameters.** CLEAR estimates effective population size \hat{N} and selection parameters,
278 \hat{s} and \hat{h} , as a byproduct of the hypothesis testing. We computed bias of selection fitness ($s - \hat{s}$)
279 and dominance ($h - \hat{h}$) for of CLEAR and GP for 1000 simulations in each setting. The distribution
280 of the error (bias) for $100\times$ coverage is presented in Fig 6 for different configurations. Fig S4 and
281 Fig S5 provide the distribution of estimation errors for $30\times$, and $300\times$ coverage, respectively. For
282 hard sweep, CLEAR provides estimates of s with lower variance of bias (Fig 6A). In soft sweep, GP
283 and CLEAR both provide unbiased estimates of s with low variance (Fig 6B). Fig 6 C-D shows that
284 CLEAR provides unbiased estimates of h as well when $h \in \{0, 0.5, 1, 2\}$ and $s = 0.1$. We also tested
285 if CLEAR provide unbiased estimates of N , by estimating population size on 1000 simulations when
286 $N \in \{200, 600, 1000\}$. As shown in Fig 7A-C, maximum likelihood is attained at true value of the
287 parameter.

288 **Running Time.** As CLEAR does not compute exact likelihood of a region (i.e., does not explicitly
289 model linkage between sites), the complexity of scanning a genome is linear in number of polymor-
290 phisms. Calculating score of each variant requires and $\mathcal{O}(TRN^3)$ computation for \mathcal{H} . However,
291 most of the operations are can be vectorized for all replicates to make the effective running time
292 for each variant. We conducted 1000 simulations and measured running times for computing site
293 statistics H , FIT, CMH and GP with different number of linked-loci. Our analysis reveals (Fig 4)
294 that CLEAR is orders of magnitude faster than GP, and comparable to FIT. While slower than
295 CMH on the time per variant, the actual running times are comparable after vectorization and
296 broadcasting over variants (see below).

297 These times can have a practical consequence. For instance, to run GP in the single locus
298 mode on the entire pool-seq data of the *D. melanogaster* genome from a small sample ($\approx 1.6\text{M}$
299 variant sites), it would take 1444 CPU-hours (≈ 1 CPU-month). In contrast, after vectorizing and
300 broadcasting operations for all variants operations using `numba` package, CLEAR took 75 minutes
301 to perform an scan, including precomputation, while the fastest method, CMH, took 17 minutes.

302 3.1 Analysis of a *D. melanogaster* Adaptation to Alternating Temperatures

303 We applied CLEAR to the data from a study of *D. melanogaster* adaptation to alternating temper-
304 atures [26, 46], where 3 replicate samples were chosen from a population of *D. melanogaster* for
305 59 generations under alternating 12-hour cycles of hot stressful (28°C) and non-stressful (18°C)
306 temperatures and sequenced. In this dataset, sequencing coverage is different across replicates and
307 generations (see S2 Fig of [60]) which makes variant depths highly heterogeneous (Fig S8).

308 We first filtered out heterochromatic, centromeric and telomeric regions [25], and those variants
309 that have collective coverage of more than 1500 in all 13 populations: three replicates at the
310 base population, two replicates at generation 15, one replicate at generation 23, one replicate at
311 generation 27, three replicates at generation 37 and three replicates at generation 59. After filtering,
312 we ended up with 1,605,714 variants.

313 Next, we estimated genome-wide population size $\hat{N} = 250$ (Fig 7-E) which is consistent with
314 previous studies [33, 46]. The likelihood curves of CLEAR are sharper around the optimum compared
315 to that of Bollback et. al [11]’s method (see Supplementary Fig. 1 in [46]). Also, chromosomes 3L
316 and 3R appear to have smaller population size Fig 7-D, $\hat{N} = 200, 150$, respectively. Others have
317 made similar observations on this data. In particular, Jónás *et al.* [33] shown that the chromosome-
318 wise population size varies even more when it is computed for each replicate separately (see Table
319 1 in [33]). For instance, \hat{N} is 131 for chromosome 3R replicate 1, while it is 328 for chromosome X
320 replicate 2.

321 While it would be ideal to compute CLEAR statistic for each replicate and chromosome sepa-
322 rately, computing empirical p -values and significant regions become computationally intensive as
323 empirical null distribution of each replicate and each chromosome needs to be computed. Hence, we
324 use a single genome-wide estimate $\hat{N} = 250$ in all analyses, but we normalize statistic \mathcal{H}^* separately
325 for each chromosome.

326 We use a heuristic calculation (See S1 Text) to choose the sliding window size L as the dis-
327 tance where the LD between the favored mutation and a site $L/2\text{bp}$ away remains strong. For *D.*
328 *melanogaster* parameters, we obtained $L = 30\text{kbp}$. We computed the normalized test statistic \mathcal{H}^*
329 on sliding windows of size of 30Kbp and step size of 5Kbp over the genome (See Fig 8-A).

330 Empirical null distribution of \mathcal{H}^* was estimated by creating 100 whole genome simulations
331 (400K statistic values) as described in Section 2.5. Then, p -value of the test statistic in each
332 region in the experimental data was calculated as the fraction of the null statistic values that are
333 greater than or equal to the test statistic(see Fig S9). After correcting for multiple testing, we
334 identified 5 contiguous intervals (Fig 8) satisfying $\text{FDR} \leq 0.05$, and covering 2,829 polymorphic
335 sites. We further performed single-locus hypothesis testing on the 2,829 sites to identify 174
336 individual variants with $\text{FDR} \leq 0.01$ (Fig 8-B).

337 The final set of 174 variants fall within 32 genes (Table S3) including many Serine inhibitory
338 proteases (serpins), and other genes involved in endocytosis. Recycling of synaptic vesicles is seen
339 to be blocked at high temperature in temperature sensitive *Drosophila* mutants [36]. This is also
340 supported by GO enrichment analysis, where a single GO term ‘inhibition of proteolysis’ is found
341 to enriched (corrected p -value:0.0041). To test for dominant selection, we computed D statistic on
342 simulated neutral and experimental data, and computed p -values accordingly. After correcting for

343 multiple testing, 96 variants were discovered with $FDR \leq 0.01$ (Fig S10).

344 3.2 Analysis of Outcrossing Yeast Populations

345 We also applied CLEAR to 12 replicate samples of outcrossing yeast populations [14], where sam-
346 ples are taken at generations $\mathcal{T} = \{0, 180, 360, 540\}$. We observed a significant variation in the
347 genome-wide site frequency spectrum of certain populations over different time points for some
348 replicates (Fig S11). The variation does not have an easily identifiable cause. Therefore, we focused
349 analysis on seven replicates $r \in \{3, 7, 8, 9, 10, 11, 12\}$ with genome-wide site-frequency spectrum over
350 the time range (Fig S12).

351 We estimated population size to be $\hat{N} = 2000$ haplotypes, and computed \hat{s} , \hat{h} and H statistic
352 accordingly. To compute p -values, we created 1M single-locus neutral simulations according to
353 experimental data's initial frequency and coverage. By setting FDR cutoff to 0.05, only 18 and
354 16 variants show significant signal for directional and dominant selection, respectively (Fig S10).
355 Selected variants for directional selection are clustered in two regions, which match 2 of the 5 regions
356 (regions C and E in Fig. 2-a in [14]) identified by Burke *et al.* in their preliminary analysis.

357 4 Discussion

358 We developed a computational tool, CLEAR, that can detect regions and variants under selection
359 E&R experiments. Using extensive simulations, we show that CLEAR outperforms existing methods
360 in detecting selection, locating the favored allele, and estimating model parameters. Also, while be-
361 ing computationally efficient, CLEAR provide means for estimating populations size and hypothesis
362 testing.

363 Many factors such as small population size, finite coverage, linkage disequilibrium, finite sam-
364 pling for sequencing, duration of the experiment and the small number of replicates can limit the
365 power of tools for analyzing E&R. Here, by an discrete modeling, CLEAR estimates population size,
366 and provides unbiased estimates of s, h . It adjusts for heterogeneous coverage of pool-seq data, and
367 exploits presence of linkage within a region to compute composite likelihood ratio statistic.

368 It should be noted that, even though we described CLEAR for small fixed-size populations,
369 the statistic can be adjusted for other scenarios, including changing population sizes when the
370 demography is known. For large populations, transitions can be computed on sparse data structures,
371 as for large N the transition matrices become increasingly sparse. Alternatively, frequencies can
372 be binned to reduce dimensionality.

373 The comparison of hard and soft sweep scenarios showed that initial frequency of the favored
374 allele can have a nontrivial effect on the statistical power for identifying selection. Interestingly,
375 while it is easier to detect a region undergoing strong selection, it is harder to locate the favored
376 allele in that region.

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

381 **Software and Data Availability.** The source code and running scripts for CLEAR are publicly
382 available at <https://github.com/airanmehr/clear>.

383 *D. melanogaster* data originally published [26, 46]. The dataset of the *D. melanogaster* study,
384 until generation 37, is obtained from Dryad digital repository (<http://datadryad.org>) under
385 accession DOI: 10.5061/dryad.60k68. Generation 59 of the *D. melanogaster* study is accessed

386 from European Sequence Read Archive (<http://www.ebi.ac.uk/ena/>) under the project accession
387 number: PRJEB6340. The dataset containing experimental evolution of Yeast populations [14]
388 is downloaded from <http://wfitcch.bio.uci.edu/~tdlong/PapersRawData/BurkeYeast.gz> (last
389 accessed 01/24/2017). UCSC browser tracks for *D. melanogaster* and Yeast data analysis are found
390 in Suppl. Data 1 and 2, respectively.

391 **Acknowledgments**

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

394 **Conflict of interest**

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

399 **Figures**

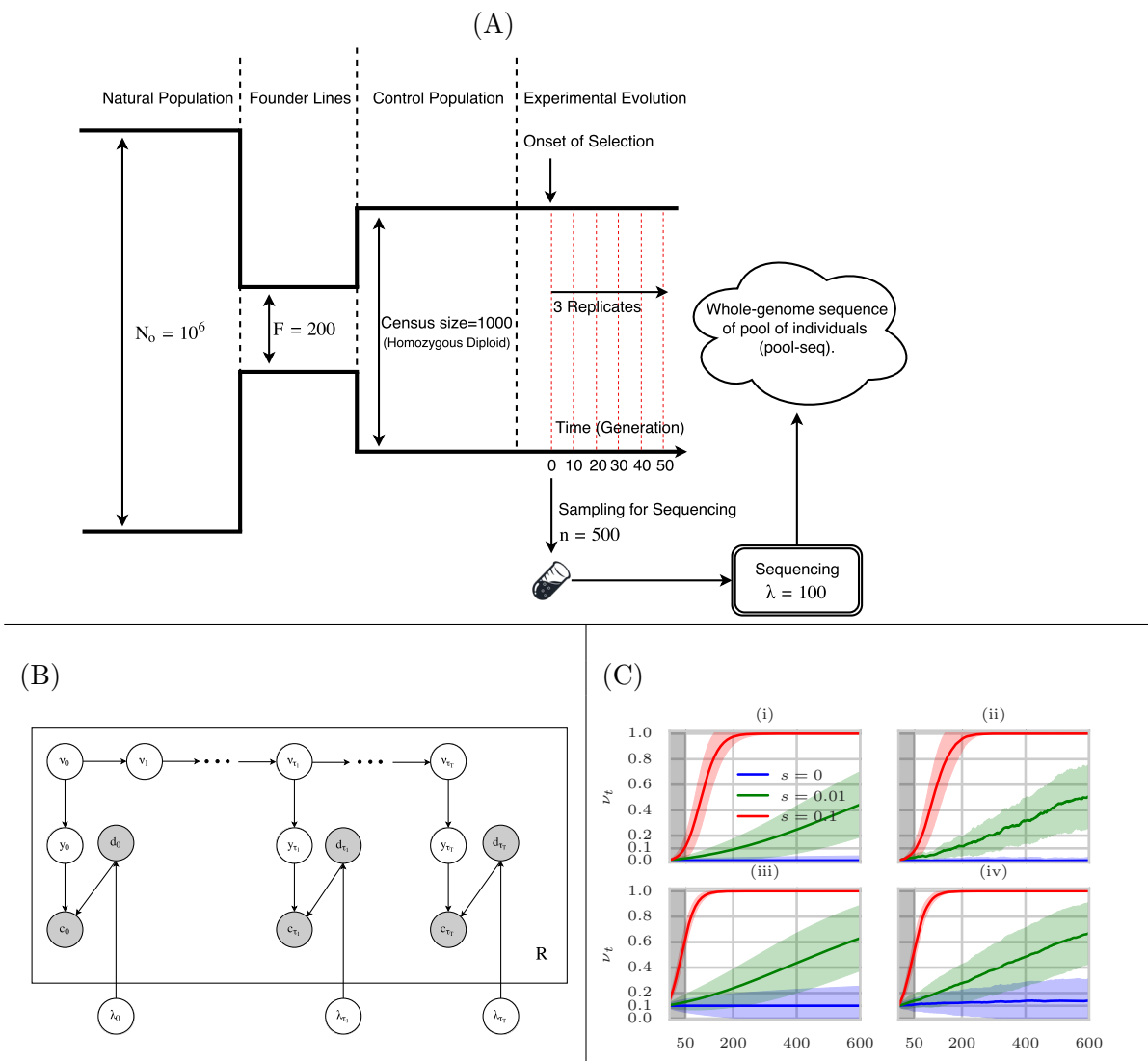


Fig 1: Evolve and Resequence Selection Experiments on *D. melanogaster*. (A) Typical configuration in which time-series data is collected for *D. melanogaster*. A small set of founder lines ($F = 200$) is selected from a large population ($N_o = 10^6$), and used to create a sub-population of isofemale lines. Multiple replicates of the population are evolved and resequenced to collect time-series genomic data. For sequencing, n individuals are randomly sampled and sequenced with coverage λ . (B) Graphical model showing dependence of the random variables in the single-locus model used to compute CLEAR statistics. Observed variables, c (derived allele read count) and d (total read count) are shaded. The variables ν, y, λ denote allele frequency, sampled allele frequency, and mean sequencing coverage, respectively. (C) Mean and 95% confidence interval of the theoretical (i,iii) and empirical (ii,iv) trajectories of the favored allele for hard (i,ii) and soft (iii,iv) sweep scenarios and $N = 1000$. The first 50 generations are shaded in gray to represent the sampling span of sampling in short-term experiments, illustrating the difficulty in predicting selection at early stages of selective sweep.

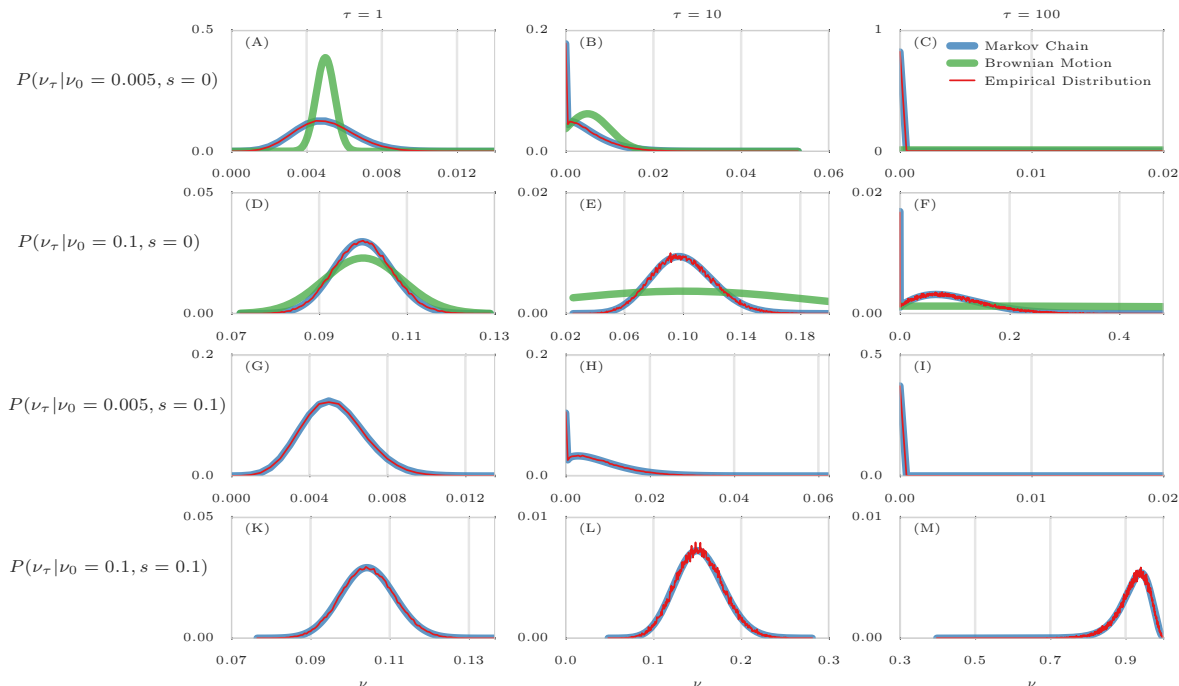


Fig 2: Comparison of empirical distributions of allele frequencies (red) versus predictions from Brownian Motion (green), and Markov chain (blue).

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 provides a good approximation in all cases.

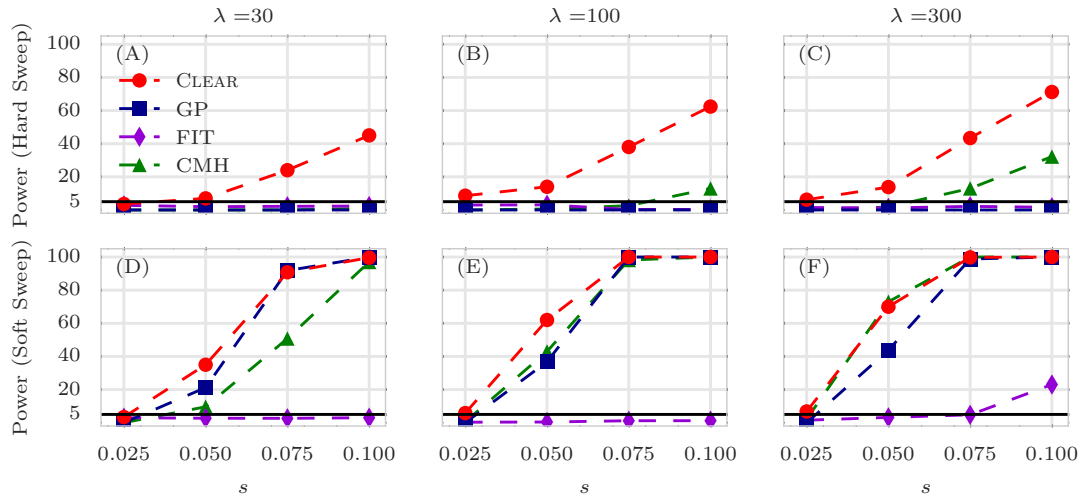


Fig 3: **Power calculations for detection of selection.**

Detection power for 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 $FPR \leq 0.05$ – for 2,000 simulations with $N = 1,000$, $L = 50\text{Kbp}$. The horizontal line reflects the power of a random classifier. In all simulations, 3 replicates are evolved and sampled at generations $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

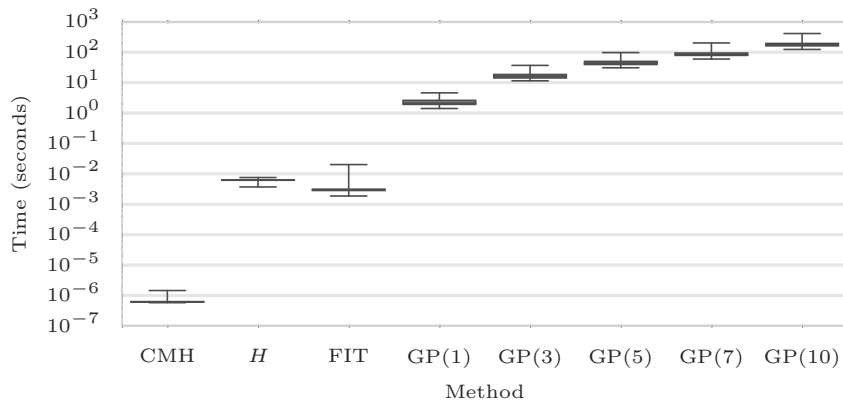


Fig 4: **Running time.**

Box plots of running time per variant (CPU-secs.) of CLEAR(\mathcal{H}), CMH, FIT, and GP with single, 3, 5, 7, and 10 loci over 1000 simulations conducted on a workstation with 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 $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

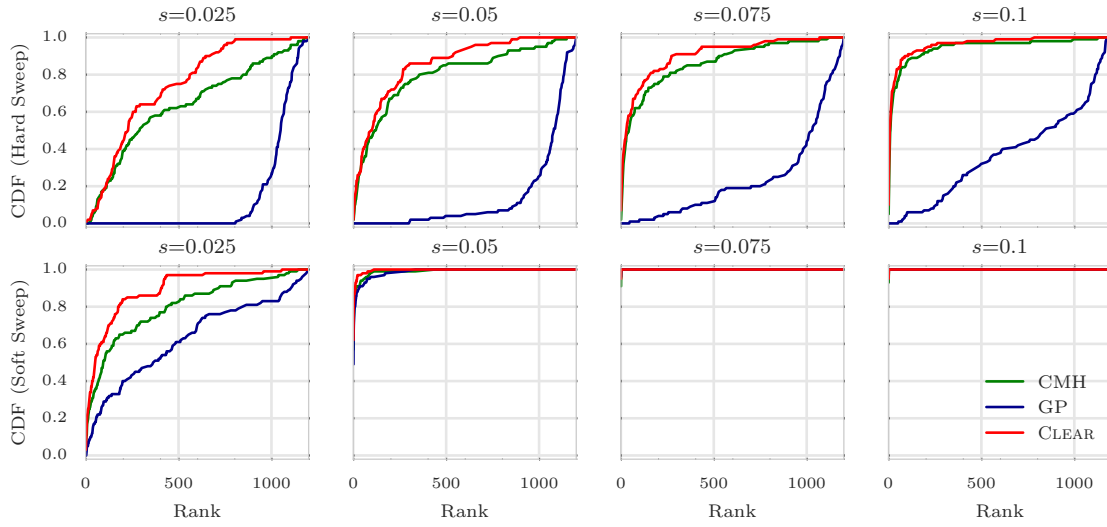


Fig 5: Ranking performance for 100x coverage.

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 $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

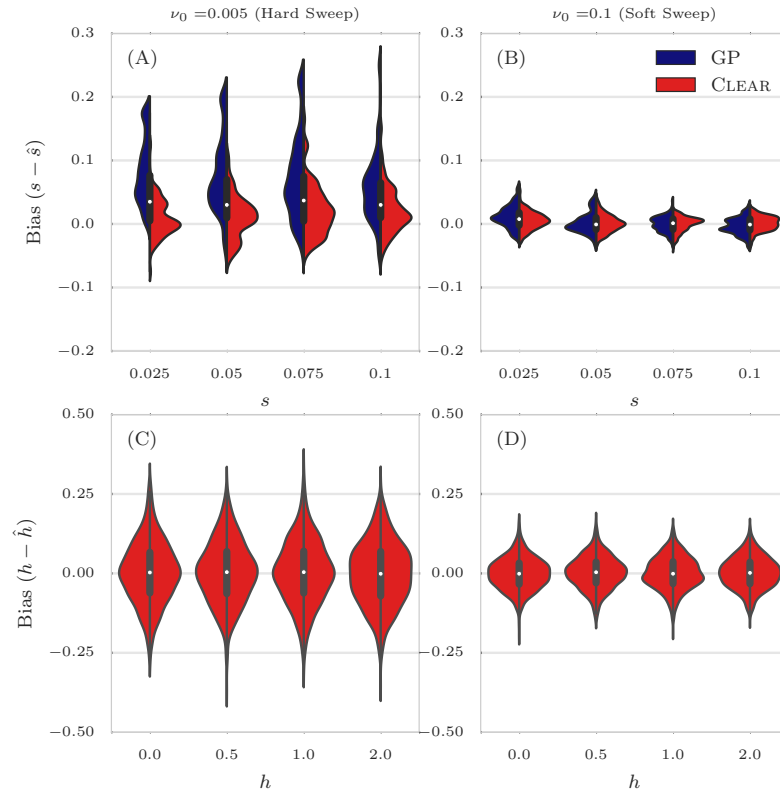


Fig 6: Distribution of bias for $100\times$ coverage.

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 [Table S2](#). Panels C,D show the variance in the estimation of h . In all simulations, 3 replicates are evolved and sampled at generations $\mathcal{T} = \{0, 10, 20, 30, 40, 50\}$.

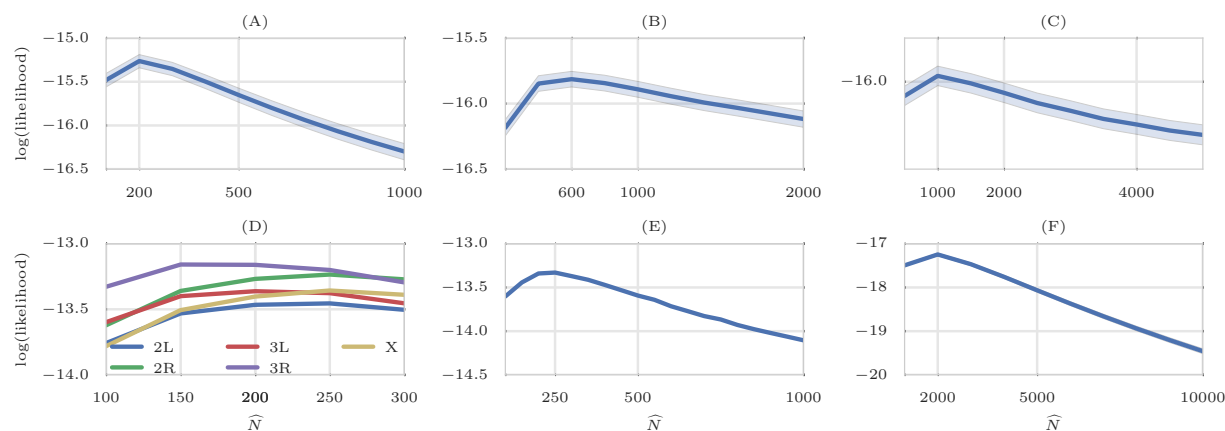


Fig 7: **Maximum likelihood Estimates of N .** Mean and 95% confidence interval of likelihoods of N on simulated data with $N = 200$ (A), $N = 600$ (B), and $N = 1000$ individuals, over 1000 simulations. Chromosome-wise (D) and genome-wide (E) likelihood of population size for data from a study of *D. melanogaster* adaptation to alternating temperatures. Likelihood of the Chromosome 3R is attained at 150, while genome-wide maximum likelihood estimate for population size is 250. (F) Likelihood of the population size with respect to all the variants in the yeast dataset. Despite large census population size ($10^6 - 10^7$ [14]), this dataset exhibits much smaller effective population size ($\hat{N} = 2000$).

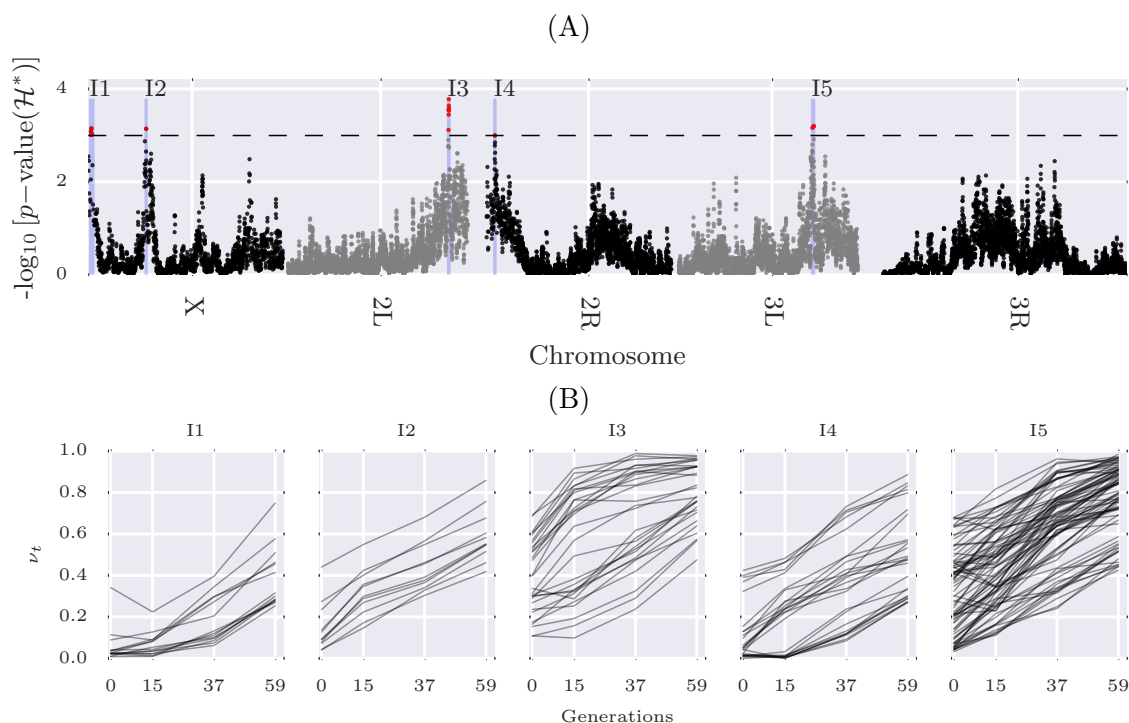


Fig 8: **Scan of CLEAR statistic on data from a study of *D. melanogaster* adaptation to alternating temperatures.** (A) Manhattan plot of scan for \mathcal{H}^* statistic over the genome. The dashed line represents cutoff for genome-wide FDR ≤ 0.05 , and identifies 5 contiguous intervals, I1-I5, which are shaded in blue. (B) Trajectories of the selected variants within intervals I1-I5.

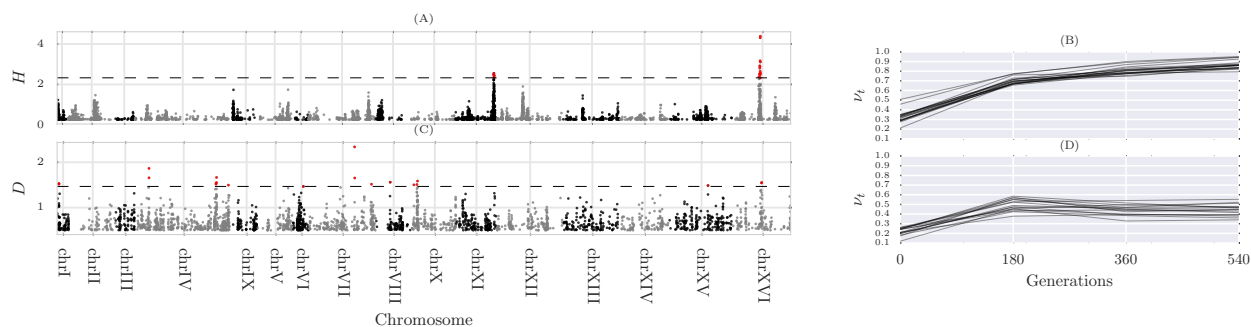


Fig 9: Single locus analysis of the yeast outcrossing populations. Manhattan plot of scan for testing directional selection (A) and dominant selection (C). The dashed line represents cutoff for genome-wide $FDR \leq 0.05$. Trajectories of the selected variants are depicted in panels (B) and (D).

400 S1 Text Choosing Window Size

401 In genome-wide scans for detecting selection, we apply the CLEAR statistic on sliding windows of
402 length L bp. The single locus statistic values within the window are averaged to get the composite
403 statistic. While the statistic is robust to variation in window-size, choosing a very large window
404 where LD has decayed will weaken the composite signal, and choosing a small window will decrease
405 the power of composite likelihoods. Here, we use a systematic calculation to choose L as the
406 distance where the LD between the favored mutation and a site $L/2$ bp away remains strong.

407 Consider a segregating site l bp away from the favored allele in a selective sweep. Let ρ_τ be
408 the LD between the favored allele and the site, τ generations after the onset of selection. Then, we
409 have (see Eqs. 30-31 in [58]):

$$\rho_\tau = \alpha_\tau \beta_\tau \rho_0 = e^{-r\tau l} \left(\frac{K^{(\tau)}}{K^{(0)}} \right) \rho_0, \quad (\text{S1})$$

410 where $K^{(\tau)} = 2\nu_\tau(1 - \nu_\tau)$ is the heterozygosity at the selected site, r is the recombination rate
411 (crossovers/bp/gen). The ‘decay factor’, $\alpha_\tau = e^{-r\tau l}$, and ‘growth factor’, β_τ , are due to recombina-
412 tion and selection, respectively. Under regular parameter settings, linkage to the favored allele is
413 expected to increase after onset of selection and then decreases due to crossover events (See Fig S13-
414 A). While ρ_0 is unknown in pool-seq E&R experiments, we compute the value of l so that

$$\alpha_\tau \beta_\tau = 1. \quad (\text{S2})$$

415 In E&R scenarios, we let τ be the time of the last sampling. For given s , we aim to compute the
416 smallest window size L over all possible starting frequencies. Specifically,

$$L = 2 \min_{\nu_0} \left\{ \frac{1}{r\tau} \log \left(\frac{\hat{\nu}_\tau(1 - \hat{\nu}_\tau)}{\nu_0(1 - \nu_0)} \right) \right\}, \quad (\text{S3})$$

417 where the term $\hat{\nu}_\tau$ depends on initial frequency ν_0 and selection strength s (Eq. 9).

418 We used *D. melanogaster* dataset parameters, $N = 250$, $r = 2 \times 10^{-8}$ and $\tau = 59$ to compute
419 the optimal window size for different values of Ns , ranging from weak selection to strong selection:
420 $Ns \in \{20, 100, 200, 500\}$, or $s \in \{0.08, 0.4, 0.8, 2\}$. We set $L = 30$ Kbp (See Fig S13-B) to provide
421 good resolution for detecting weak selection.

Generative Process 1: The Generative Process for Neutral Wright-Fisher Time-series Pool-seq Data.

Input: $N, n, R, \{\lambda_t\}_{t \in \mathcal{T}}, \mathcal{T} = \{\tau_0, \dots, \tau_T\}$

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  $\tau_T$  do
     $2N\nu_t \sim \text{Binomial}(2N, \nu_{t-1});$ 
    if  $t \in \mathcal{T}$  then
       $d_t^{(r)} \sim \text{Poisson}(\lambda_{\tau_t});$ 
       $2ny_t \sim \text{Binomial}(2n, \nu_t);$ 
       $c_t^{(r)} \sim \text{Binomial}(d_t^{(r)}, y_t);$ 
    end
  end
end

```

Fig S1: The Generative Process for Neutral Wright-Fisher Time-series Pool-seq Data.

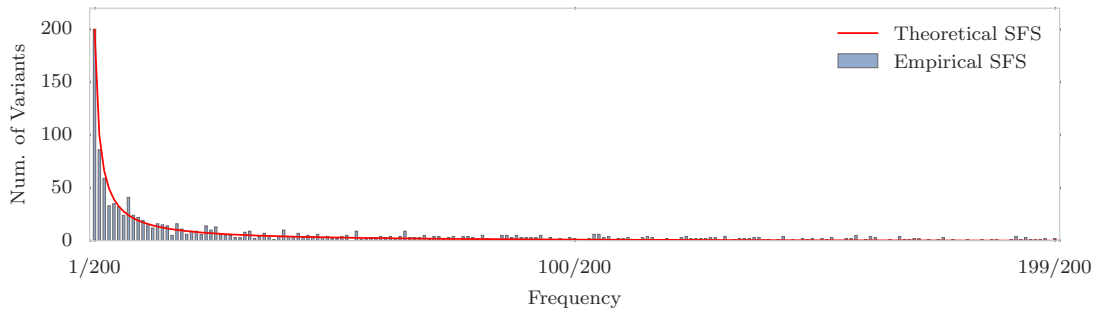


Fig S2: Site Frequency Spectrum.

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 neutral population, majority of the variations stand in low frequency.

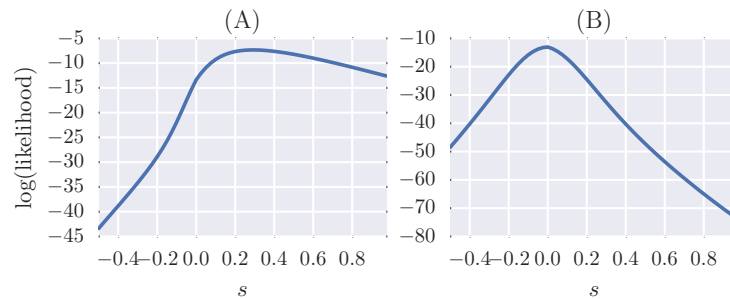


Fig S3: Likelihoods of the parameter s .

Likelihood of the parameter s in *D. melanogaster* data for a variant with $\hat{s} = 0.2$ (A) and $\hat{s} = 0$ (B).

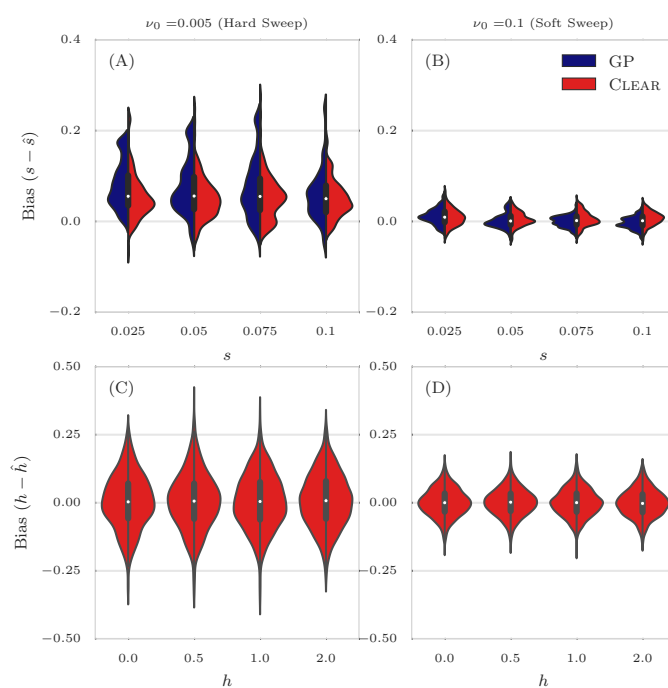


Fig S4: **Distribution of bias for $30\times$ coverage.**

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 Table S2. Panels C,D show the variance in the estimation of h .

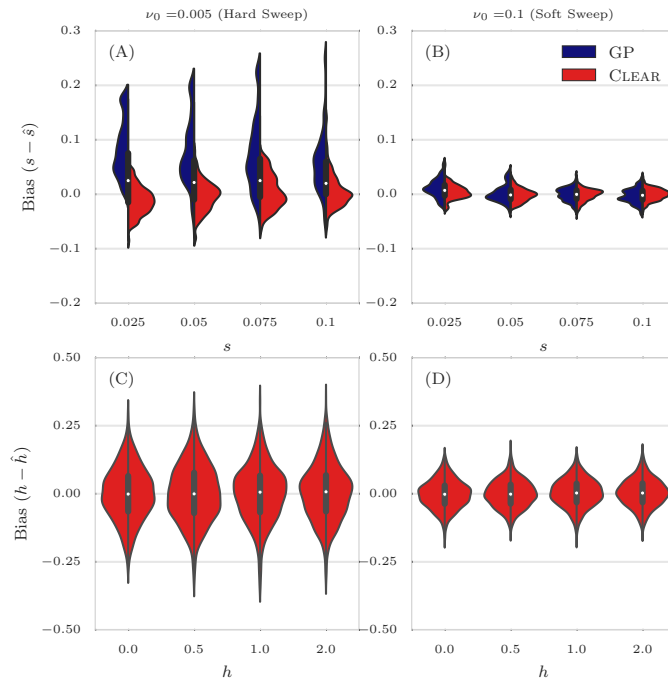


Fig S5: **Distribution of bias for 300 \times coverage.**

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 $\lambda = \infty$ (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 Table S2. Panels C,D show the variance in the estimation of h .

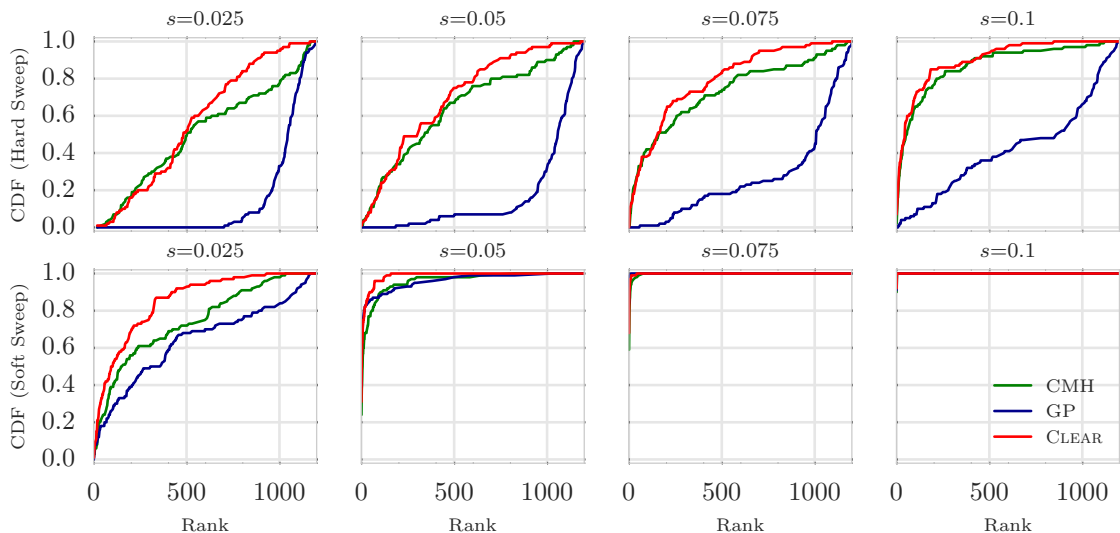


Fig S6: Ranking performance for $30\times$ coverage.

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.

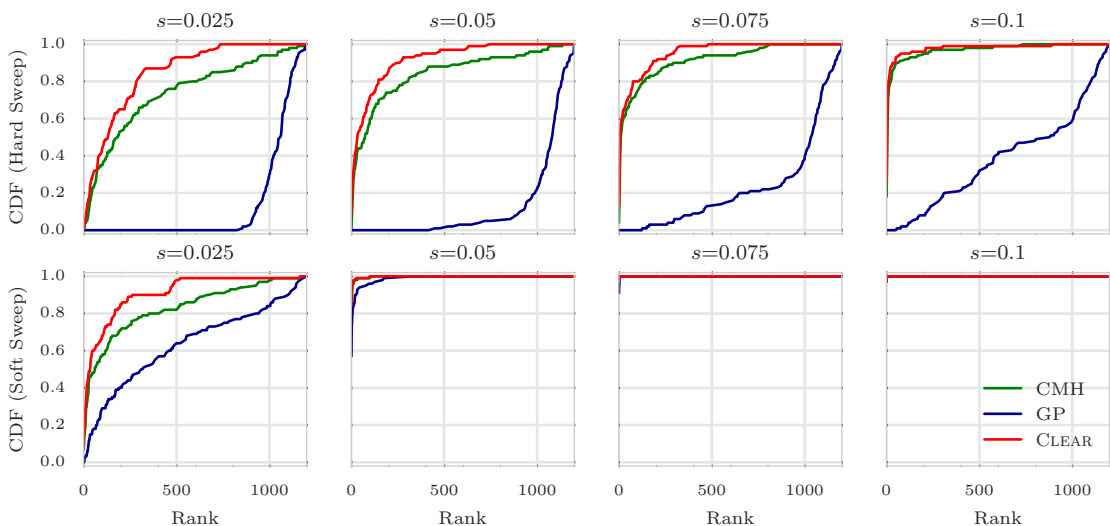


Fig S7: Ranking performance for $300\times$ coverage.

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.

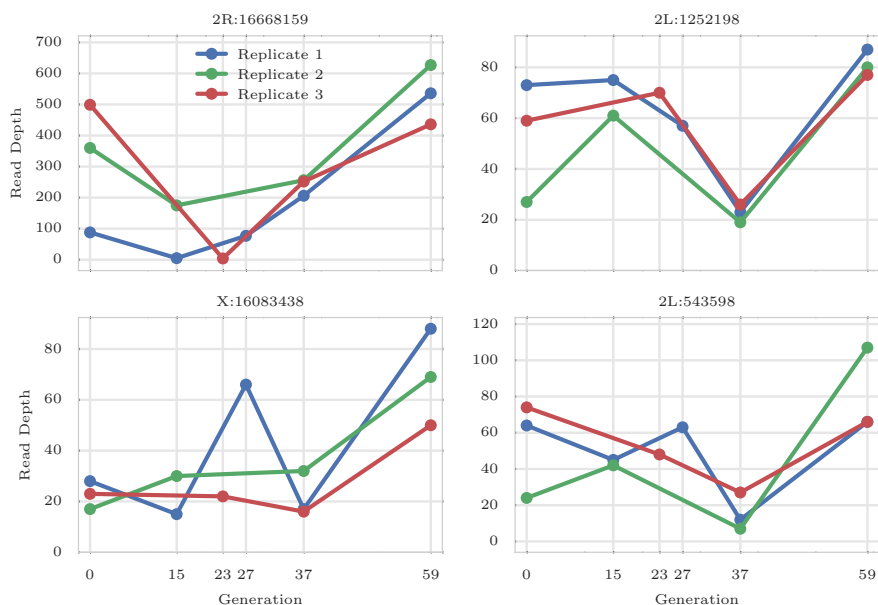


Fig S8: **Coverage heterogeneity in time series data.**

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 the hard filtering with minimum depth of 30.

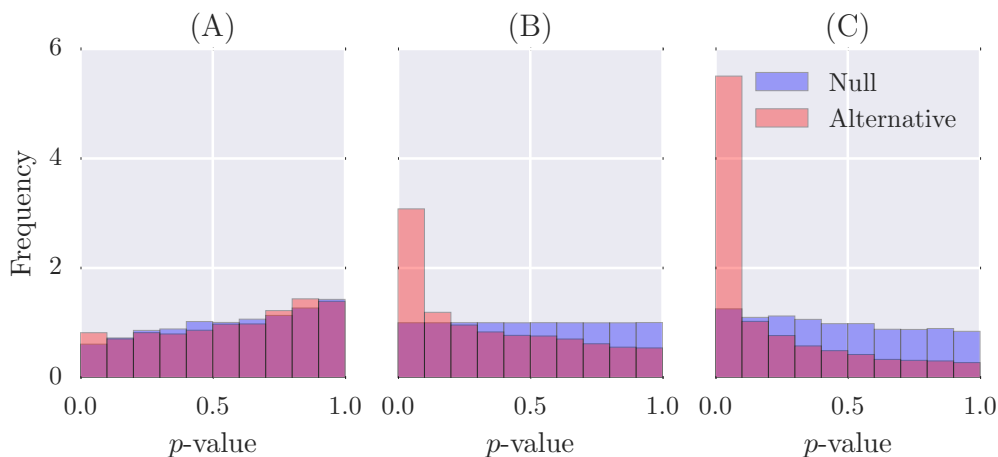


Fig S9: **Distribution of p -values.** Distribution of p -values of CLEAR in null simulations and experimental data when $N = 250$. Panel (A),(C) shows the effect of under estimations ($\hat{N} = 100$) and over-estimation ($\hat{N} = 500$) of population size in computing p -values, and panel (B) shows the distribution of p -values when unbiased estimate is used to create simulations. .

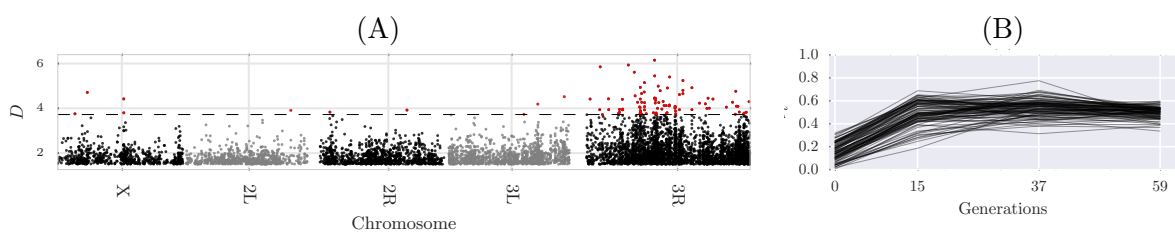


Fig S10: Single locus analysis of the data from a study of *D. melanogaster* adaptation to alternating temperatures.

Manhattan plot of scan for testing dominant selection (A). Significant variants with $FDR \leq 0.01$ are denoted in red, and their trajectories are depicted in panel (B).

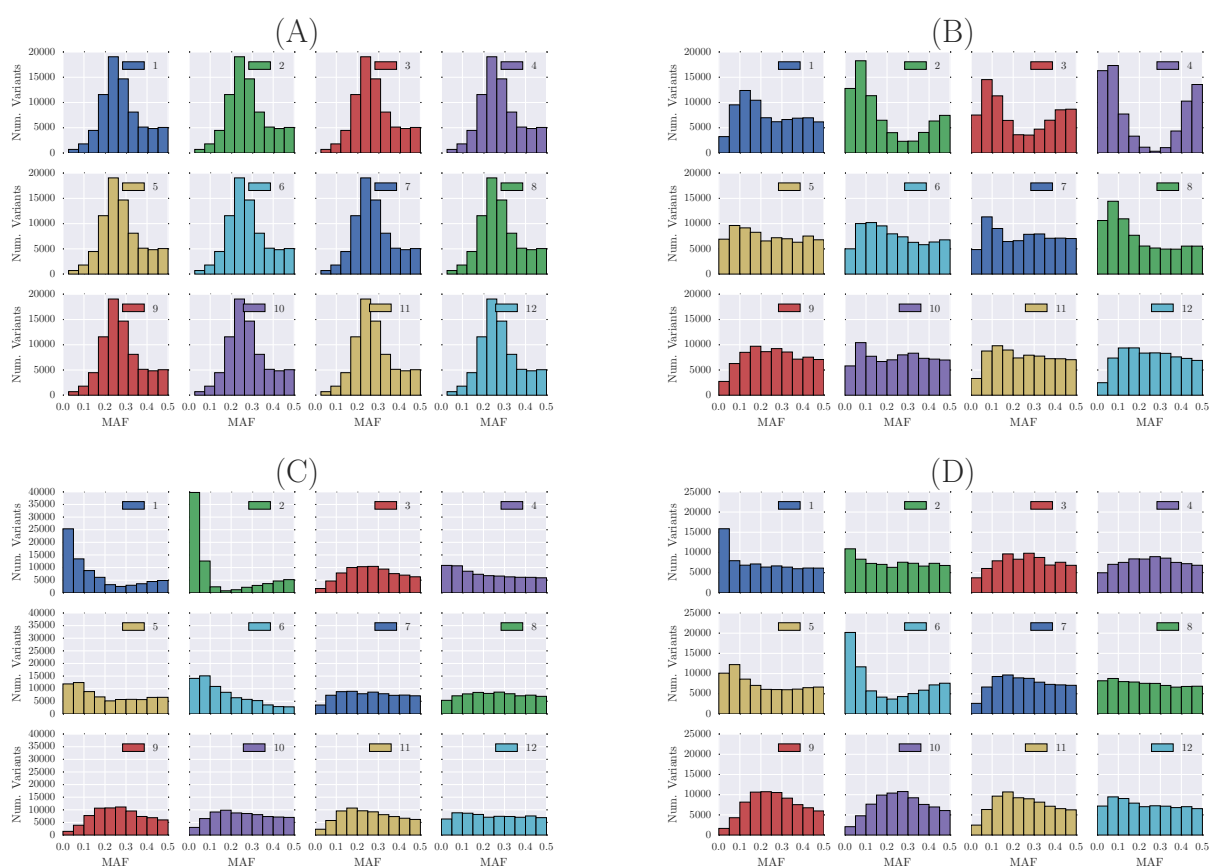


Fig S11: Site frequency spectrum of the Yeast dataset. Whole-genome site frequency spectrum of the Yeast dataset at generations 0 (A), 180 (B), 360 (C) and 540 (D). Some replicates, e.g. replicate 2, undergoing severe demographic events.

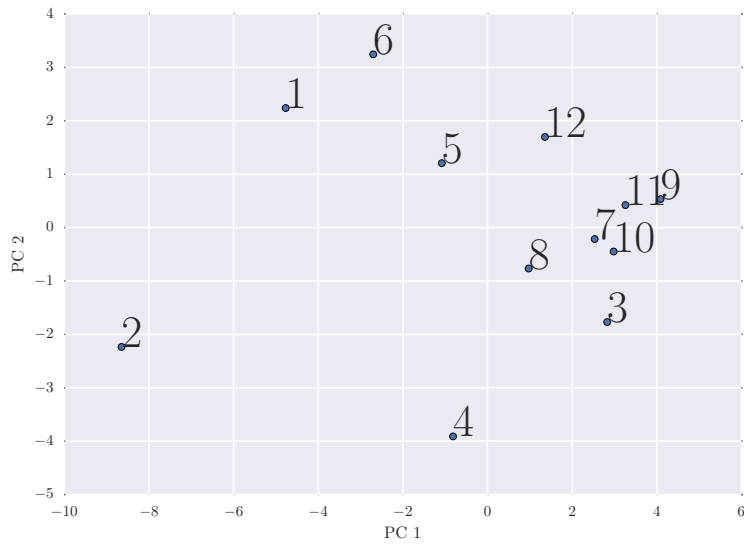


Fig S12: **Population similarity.** Principle component analysis of the 12 replicates throughout the experiment, showing that some populations exhibiting distinct frequency spectra.

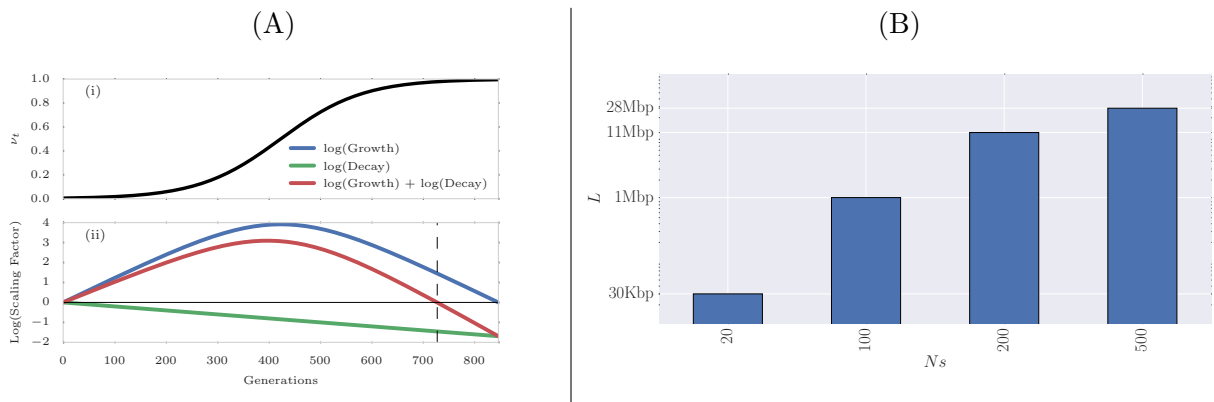


Fig S13: **Choosing window size for CLEAR statistic.** (A) Expected dynamics of LD between favored allele ($s = 0.025$) and a variant 50Kbp away, with initial frequency $\nu_0 = 0.01$. (A-i) depicts the dynamic of the favored allele during the selective sweep. (A-ii) illustrates interaction of the growth and decay factors introduced in Eq. S1, with the red line describing overall effect of selection and recombination on LD. The vertical dashed line points to the time when the LD value starts to decrease below original LD. (B) Alternatively, we can fix time, and find the window-size at which LD decays below the original LD (Eq. S3). The plot shows the window size as a function of Ns (20,100,200,500), after fixing other model parameters to match *D. melanogaster* E&R experiments ($N = 250$, $r = 2 \times 10^8$, $\tau = 59$).

Table S1: **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 (\mathcal{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.

Table S2: **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

Table S3: **Overlapping genes with the 174 candidate variants.**

Interval	Position	FBgn	Gene Name	GO Function
I1	X:1.567-1.824M	FBgn0023531	CG32809	NA
		FBgn0023130	a6	embryonic development via the syncytial blastoderm
		FBgn0025378	CG3795	serine-type endopeptidase activity
		FBgn0025391	Scgdelta	heart contraction, mesoderm development
		FBgn0261548	CG42666	NA
		FBgn0026086	Adar	RNA editing
		FBgn0026090	CG14812	negative regulation of cysteine-type endopeptidase activity involved in apoptotic process
I2	X:7.175-7.241M	FBgn0023522	CG11596	NA
		FBgn0029941	CG1677	NA
		FBgn0029944	Dok	stress activated protein kinase signaling
I3	2L:16.878-16.993M	FBgn0029946	CG15034	NA
		FBgn0052832	CG32832	mitochondrial pyruvate transport
		FBgn0032618	CG31743	sulfotransferase activity
		FBgn0085342	CG34313	NA
		FBgn0040985	CG6115	NA
		FBgn0261671	tweek	synaptic vesicle endocytosis
		FBgn0026150	ApepP	metalloaminopeptidase activity
I4	2R:2.725-2.810M	FBgn0262355	CR43053	NA
		FBgn0053179	beat-IIIb	NA
		FBgn0040674	CG9445	NA
		FBgn0265935	coro	adult somatic muscle development
		FBgn0033110	CG9447	NA
		FBgn0033113	Spn42Dc	Inhibitory Serpins
I5	3L:14.362-14.514M	FBgn0028988	Spn42Dd	Inhibitory Serpins
		FBgn0033115	Spn42De	Inhibitory Serpins
		FBgn0050158	CG30158	small GTPase mediated signal transduction
		FBgn0036421	CG13481	ubiquitin-protein transferase activity
		FBgn0262580	CG43120	NA
I5	3L:14.362-14.514M	FBgn0036422	CG3868	NA
		FBgn0087007	bbg	PDZ domain
		FBgn0036426	CG9592	NA
		FBgn0036427	CG4613	serine-type endopeptidase activity

References

- [1] Alan Agresti and Maria Kateri. *Categorical data analysis*. Springer, 2011.
- [2] 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.
- [3] Frédéric Arieu, 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.
- [4] 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.
- [5] 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.
- [6] Jeffrey E Barrick and Richard E Lenski. Genome dynamics during experimental evolution. *Nature Reviews Genetics*, 14(12):827–839, 2013.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] Jonathan P Bollback, Thomas L York, and Rasmus Nielsen. Estimation of 2Nes from temporal allele frequency data. *Genetics*, 179(1):497–502, 2008.
- [12] 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.
- [13] 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.
- [14] 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.

- [15] 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.
- [16] 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.
- [17] Vincent J Deneff and Jillian F Banfield. In situ evolutionary rate measurements show ecological success of recently emerged bacterial hybrids. *Science*, 336(6080):462–466, 2012.
- [18] Michael M Desai and Joshua B Plotkin. The polymorphism frequency spectrum of finitely many sites under selection. *Genetics*, 180(4):2175–2191, 2008.
- [19] Richard Durbin, Sean R Eddy, Anders Krogh, and Graeme Mitchison. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge university press, 1998.
- [20] Warren J Ewens. *Mathematical Population Genetics 1: Theoretical Introduction*, volume 27. Springer Science & Business Media, 2012.
- [21] 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.
- [22] 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.
- [23] Alison F Feder, Sergey Kryazhimskiy, and Joshua B Plotkin. Identifying signatures of selection in genetic time series. *Genetics*, 196(2):509–522, 2014.
- [24] 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.
- [25] Anna-Sophie Fiston-Lavier, Nadia D Singh, Mikhail Lipatov, and Dmitri A Petrov. *Drosophila melanogaster* recombination rate calculator. *Gene*, 463(1):18–20, 2010.
- [26] Susanne U Fransen, 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.
- [27] Michael M Gottesman. Mechanisms of cancer drug resistance. *Annual review of medicine*, 53(1):615–627, 2002.
- [28] Matthew Hegreness, Noam Shores, Daniel Hartl, and Roy Kishony. An equivalence principle for the incorporation of favorable mutations in asexual populations. *Science*, 311(5767):1615–1617, 2006.
- [29] 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.

- [30] 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.
- [31] 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.
- [32] 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.
- [33] Ágnes 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.
- [34] 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.
- [35] Robert Kofler and Christian Schlötterer. A guide for the design of evolve and resequencing studies. *Molecular biology and evolution*, page mst221, 2013.
- [36] Toshio Kosaka and Kazuo Ikeda. Reversible Blockage of Membrane Retrieval and Endocytosis in the Garland Cell of the Temperature-sensitive. *The Journal of cell biology*, 97, 1983.
- [37] 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.
- [38] 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.
- [39] 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.
- [40] 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.
- [41] 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.
- [42] 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.

- [43] 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.
- [44] 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.
- [45] 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.
- [46] 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.
- [47] 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.
- [48] Bo Peng and Marek Kimmel. simuPOP: a forward-time population genetics simulation environment. *Bioinformatics*, 21(18):3686–3687, 2005.
- [49] Edward Pollak. A new method for estimating the effective population size from allele frequency changes. *Genetics*, 104(3):531–548, 1983.
- [50] 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.
- [51] 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.
- [52] Stanley A Sawyer and Daniel L Hartl. Population genetics of polymorphism and divergence. *Genetics*, 132(4):1161–1176, 1992.
- [53] 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.
- [54] Joshua G Schraiber, Steven N Evans, and Montgomery Slatkin. Bayesian inference of natural selection from allele frequency time series. *Genetics*, 203(1):493–511, 2016.
- [55] 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.
- [56] 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.
- [57] 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.
- [58] 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.
- [59] John D Storey and Robert Tibshirani. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences*, 100(16):9440–9445, 2003.
- [60] 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.
- [61] Ray Tobler, Susanne U Franssen, Robert Kofler, Pablo Orozco-terWengel, Viola Nolte, Joachim Hermisson, 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.
- [62] 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.
- [63] 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.
- [64] Jinliang Wang. A pseudo-likelihood method for estimating effective population size from temporally spaced samples. *Genetical research*, 78(03):243–257, 2001.
- [65] Robin S Waples. A generalized approach for estimating effective population size from temporal changes in allele frequency. *Genetics*, 121(2):379–391, 1989.
- [66] David Williams and David Williams. *Weighing the odds: a course in probability and statistics*, volume 548. Springer, 2001.
- [67] 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.
- [68] 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.
- [69] Xin Yi, Yu Liang, Emilia Huerta-Sanchez, Xin Jin, Zha Xi Ping Cuo, John E Pool, Xun Xu, Hui Jiang, Nicolas Vinckenbosch, Thorfinn Sand Korneliusen, and Others. Sequencing of 50 human exomes reveals adaptation to high altitude. *Science*, 329(5987):75–78, 2010.
- [70] Hiba Zahreddine and K L Borden. Mechanisms and insights into drug resistance in cancer. *Front Pharmacol*, 4(28.10):3389, 2013.

- [71] Dan Zhou, Nitin Udpa, Merrill 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.