**Supplementary information**

# MPL resolves genetic linkage in fitness inference from complex evolutionary histories

In the format provided by the authors and unedited

# Supplementary Material

## MPL resolves genetic linkage in fitness inference from complex evolutionary histories

Muhammad Saqib Sohail[1,*], Raymond H. Y. Louie[1,2,3,4,*], Matthew R. McKay[1,5,†], and John P. Barton[6,†]

[1]Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong.
[2]Institute for Advanced Study, Hong Kong University of Science and Technology, Hong Kong.
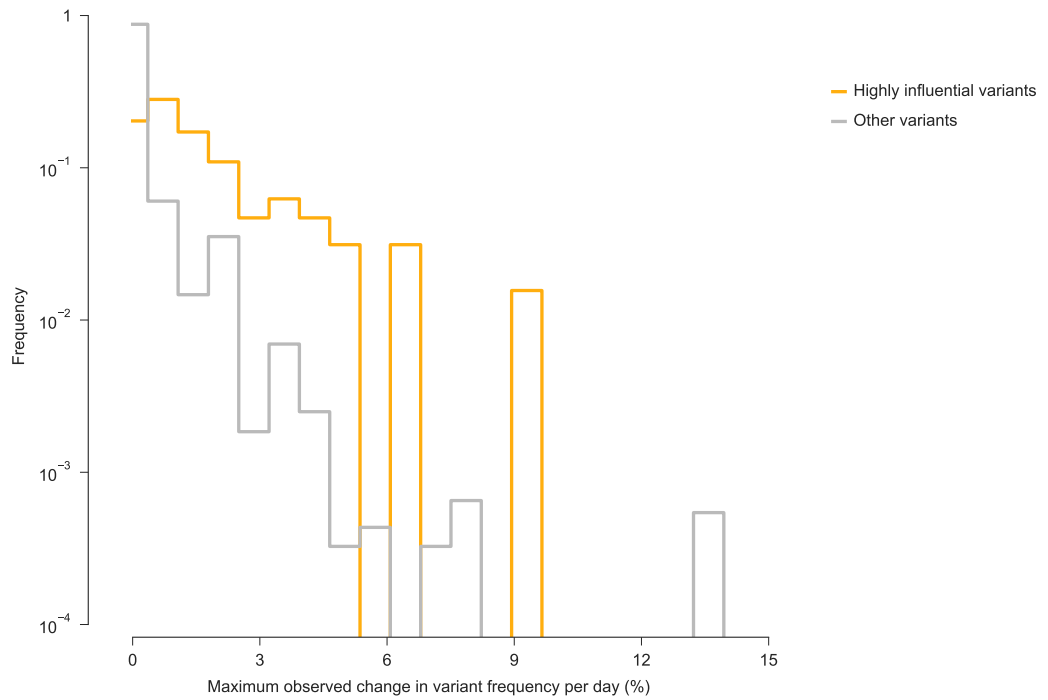[3]The Kirby Institute, University of New South Wales, Australia.
[4]School of Medical Sciences, University of New South Wales, Australia.
[5]Department of Chemical and Biological Engineering, Hong Kong University of Science and Technology, Hong Kong.
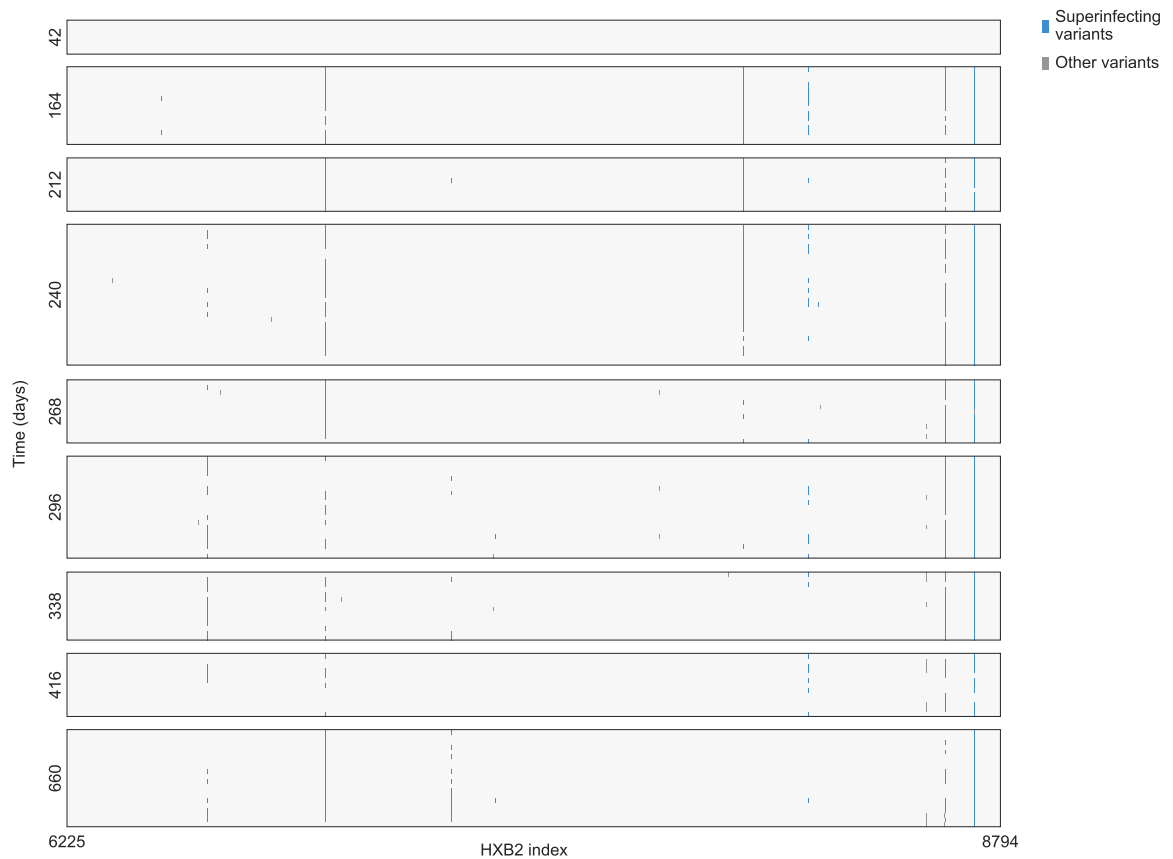[6]Department of Physics and Astronomy, University of California, Riverside, USA.
[*]These authors contributed equally to this work.
[†]Address correspondence to: m.mckay@ust.hk, john.barton@ucr.edu.

Supplementary Figure 1: **Highly influential variants are far more likely than other variants to change rapidly in frequency.** For all genetic variants across all individuals and genomic regions considered in this study, we computed the maximum change in frequency per day between successive sequence samples. For most variants (those for which $\sum_j |\Delta \hat{s}_{ij}| \leq 0.4$), the maximum change in frequency per day is less than 1%. Highly influential variants, which have large effects on inferred selection coefficients at other sites ($\sum_j |\Delta \hat{s}_{ij}| > 0.4$), are much more likely than other variants to change rapidly in frequency.

Supplementary Figure 2: **Extensive recombination between CAP256 primary and superinfecting strains.** Individual CAP256 was infected by a distinct superinfecting strain of HIV-1 15 weeks after the primary infection [1]. Each row represents a sequence, with nucleotide variants different from the primary infecting strain highlighted. Soon after superinfection, by 164 days after initial infection, recombinants between the primary and superinfecting strains dominate the viral population. Recombination continues throughout infection, introducing substantial variation into the VRC26 epitope region by 240 days after initial infection.

| Patient ID | Abbreviated name | Sampling times (days) | Number of sequences 3$'$ | Number of sequences 5$'$ | Subtype |
|---|---|---|---|---|---|
| 700010607 | CH607 | 0 | 29 | 30 | B |
| | | 9 | 5 | 8 | |
| | | 14 | 16 | 13 | |
| | | 21 | 23 | 25 | |
| 700010470 | CH470 | 0 | 12 | 25 | B |
| | | 13 | 14 | 11 | |
| | | 41 | 28 | 29 | |
| | | 69 | 23 | 10 | |
| | | 174 | 25 | 9 | |
| | | 420 | 11 | 10 | |
| | | 454 | | 10 | |
| 700010077 | CH77 | 0 | 1 | 1 | B |
| | | 14 | 15 | 17 | |
| | | 32 | 12 | 6 | |
| | | 102 | 5 | | |
| | | 159 | 11 | 8 | |
| 700010058 | CH58 | 0 | 1 | 1 | B |
| | | 8 | 6 | 6 | |
| | | 45 | 9 | 9 | |
| | | 85 | 9 | 9 | |
| | | 154 | | 7 | |
| | | 239 | | 8 | |
| | | 252 | | 8 | |
| | | 350 | | 4 | |
| 700010040 | CH40 | 0 | 8 | 12 | B |
| | | 16 | 9 | 9 | |
| | | 45 | 14 | 12 | |
| | | 111 | 9 | 8 | |
| | | 181 | 7 | 11 | |
| | | 283 | 11 | 4 | |
| | | 412 | 12 | 11 | |
| | | 552 | 12 | 7 | |
| 706010164 | CH164 | 0 | 1 | 1 | C |
| | | 14 | 37 | 27 | |
| | | 28 | 10 | 17 | |
| | | 70 | 23 | 15 | |
| | | 183 | 11 | 14 | |
| | | 434 | 20 | 24 | |
| 705010198 | CH198 | 0 | 11 | 14 | C |
| | | 11 | 10 | 17 | |
| | | 60 | 27 | 16 | |
| 705010185 | CH185 | 0 | 28 | 11 | C |
| | | 25 | 11 | 16 | |
| | | 67 | 27 | 22 | |
| | | 180 | 15 | | |
| | | 416 | 16 | | |
| 705010162 | CH162 | 0 | 14 | 16 | C |
| | | 21 | 17 | 12 | |
| | | 77 | 7 | 4 | |
| | | 179 | 9 | 8 | |
| | | 438 | 22 | 20 | |
| 704010042 | CH42 | 0 | 9 | 11 | C |
| | | 21 | 17 | 18 | |
| | | 60 | 6 | 7 | |
| | | 172 | 12 | 12 | |
| | | 424 | 24 | 21 | |
| | | 676 | 25 | 16 | |
| 703010256 | CH256 | 0 | 11 | 11 | C |
| | | 28 | 12 | 10 | |
| | | 63 | 27 | 30 | |
| | | 172 | 10 | 13 | |
| | | 426 | 15 | 16 | |
| | | 684 | 24 | 30 | |
| 703010131 | CH131 | 0 | 33 | 8 | C |
| | | 21 | 10 | 7 | |
| | | 28 | 8 | 5 | |
| | | 34 | 12 | 7 | |
| | | 63 | 12 | 12 | |
| | | 91 | 9 | 6 | |
| | | 175 | 12 | 12 | |
| | | 273 | 9 | 10 | |
| | | 333 | 9 | 9 | |
| 703010159 | CH159 | 0 | 23 | 16 | C |
| | | 8 | 12 | 11 | |
| | | 12 | 13 | 10 | |
| | | 22 | 6 | 5 | |
| | | 29 | 7 | 15 | |
| | | 56 | 13 | 8 | |
| | | 85 | 10 | 10 | |
| | | 302 | 14 | 18 | |
| CAP256 | | 42 | 7 | | C |
| | | 164 | 16 | | |
| | | 212 | 11 | | |
| | | 240 | 29 | | |
| | | 268 | 13 | | |
| | | 296 | 21 | | |
| | | 338 | 14 | | |
| | | 416 | 13 | | |
| | | 660 | 20 | | |

Supplementary Table 1: **Summary of the HIV-1 data.** Sampling times and number of sequences per time point for each individual and sequencing region after data processing. Most patients have similar sampling profiles. For example, trajectory lengths are several hundred days in length for all but 3 individuals, CH198, CH77 and CH607, where the trajectory lengths are $< 200$ days. Typically, around 7 to 15 sequences are available per time point, and most of the times between successive samples $\Delta t$ are $< 50$ days, with the great majority $< 100$ days.

# Supplementary Text

## Contents

# 1   Path integral approximation and MPL estimator of selection coefficients

This section presents detailed derivations of the path integral approximation and MPL estimator, as described in Methods. To summarize, we first describe the WF model in Section 1.1, which accounts for evolutionary features including mutation, selection, drift, recombination, and incomplete temporal sampling. An assumption in this model, which will be subsequently relaxed, is that there are only two alleles at each locus, and mutation probabilities are symmetric. After introducing the model, in Section 1.2 we present conditional moment expansions that are used in Section 1.3 to derive (S21), the path integral expression for the probability of a path of *genotype* frequencies $(\boldsymbol{z}(t_1), \boldsymbol{z}(t_2), \dots, \boldsymbol{z}(t_K))$, conditioned on $\boldsymbol{z}(t_0)$. In Section 1.4, this genotype-level analysis is used to derive (S26), the path integral expression for the probability of a path of *allele* frequencies $(\boldsymbol{x}(t_1), \boldsymbol{x}(t_2), \dots, \boldsymbol{x}(t_K))$, conditioned on $\boldsymbol{x}(t_0)$. The MPL estimate of the selection coefficients obtained from the path integral is derived in Section 1.5. We then show in Section 1.6 that the same MPL estimate can also be derived from the genotype-level path integral, demonstrating no loss of optimality from derivations based on an allele-level analysis. Finally, extensions to account for multiple alleles per locus and asymmetric mutation probabilities are developed in Section 1.7.

Since our technical developments build on the basic machinery of diffusion approximations, tools that are commonly applied in population genetics (see e.g., refs.[2, 3, 4, 5, 6]), a few words on the novelty of our analysis are in order. Most notably, in addition to considering a multi-locus model which accounts for the effects of mutation and selection, a novel feature of our work is that we explicitly account for recombination and incomplete temporal sampling. These effects are practically important, but their incorporation complicates the analysis. Further extension to account for multi-allele multi-locus dynamics with asymmetric mutational probabilities also presents a new contribution, and one that is necessary for analyzing the HIV data in the main text. While path integral methods are familiar in physics[7], these are relatively less known and employed in population genetics. Some exceptions include refs. [8, 9], where they were employed for purposes other than statistical inference, and were developed under less generalized model assumptions. Other related work has proposed strategies for inferring selection from observed evolutionary histories based on deterministic models of population dynamics that ignore genetic drift[10]. The analysis we report here presents a first path integral based framework for statistical inference which applies for generalized multi-locus WF models. As a further conceptual novelty, our analysis shows that, for the model setting that we consider and under the diffusion limit, the fitness effect of mutations are completely captured by the individual and pairwise mutation probabilities, and higher order mutational correlations contain no further information.

## 1.1   Wright Fisher model

We start with a brief introduction to the WF model and the generalized parameter settings that we consider. At generation $t$, denote $\boldsymbol{Z}(t) = (Z_1(t), \dots, Z_M(t))$ as the random genotype frequency vector, and thus $\boldsymbol{z}(t) = (z_1(t), \dots, z_M(t))$

corresponds to an observed realization of this random vector. (These vectors, as well as all other vectors that we define, are taken as column vectors.) The stochastic dynamics of the genotype frequencies are governed by the transition probabilities reported in Eq. (2)-(4) of Methods. At generation $t+1$, $y_a(t)$ recombination occurs, and as a consequence of this action, the mean frequency of genotype $a$, conditioned on $\boldsymbol{Z}(t) = \boldsymbol{z}(t)$, is given by

$$y_a(t) = (1-r)^{L-1} z_a(t) + \left(1 - (1-r)^{L-1}\right) \psi_a(\boldsymbol{z}(t))$$

where the factor $(1-r)^{L-1}$ represents the probability of an individual not undergoing recombination, and $\psi_a(\boldsymbol{z}(t))$ the probability of forming genotype $a$ after recombination from individuals in generation $t$. This latter quantity is given by

$$\psi_a(\boldsymbol{z}(t)) = \sum_{c=1}^{M} \sum_{d=1}^{M} R_{a,cd} z_c(t) z_d(t) \tag{S1}$$

with $R_{a,cd}$ denoting the probability that genotypes $c$ and $d$ recombine to form genotype $a$. The probability $R_{a,cd}$ is a complicated function of the number of breakpoints and the particular genotypes $a$, $c$ and $d$; however as we will show, we do not require the exact form of $R_{a,cd}$ for our derivations. After recombination, the mean genotype frequencies at generation $t+1$ are further shaped through selection and mutation. Specifically, after recombination, selection and mutation, the mean frequency of genotype $a$, conditioned on $\boldsymbol{Z}(t) = \boldsymbol{z}(t)$, admits

$$p_a(\boldsymbol{z}(t)) := \mathrm{E}\left[Z_a(t+1)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right]$$

$$= \frac{y_a(t) f_a + \sum_{b=1, b \neq a}^{M} \left(\mu_{ba} y_b(t) f_b - \mu_{ab} y_a(t) f_a\right)}{\sum_{b=1}^{M} y_b(t) f_b}.$$

Then, the WF dynamics, specified by the transition probability

$$P\left(\boldsymbol{z}(t+1)\Big|\boldsymbol{z}(t)\right) = N! \prod_{a=1}^{M} \frac{\left(p_a(\boldsymbol{z}(t))\right)^{N z_a(t+1)}}{(N z_a(t+1))!},$$

simply results from performing random multinomial sampling (of $N$ individuals) with these mean frequencies.

Further, we recall the assumption that $\mu_{ab} = \mu^{d_{ab}}$, with $d_{ab}$ the Hamming distance between genotypes $a$ and $b$, and the additive model for genotype fitness, such that the selection coefficient of genotype $a$ admits

$$h_a = f_a - 1 = \sum_{j=1}^{L} g_j^a s_j. \tag{S2}$$

We may then write

$$p_a(\boldsymbol{z}(t)) = \frac{(1+h_a) y_a(t) + \sum\limits_{b=1}^{M} \mu^{d_{ab}} \left((1+h_b) y_b(t) - (1+h_a) y_a(t)\right)}{\sum\limits_{b=1}^{M} (1+h_b) y_b(t)}. \tag{S3}$$

We work under the assumption that the population size $N$ is large, and that as $N \to \infty$,

$$s_i = \frac{\bar{s}_i}{N} + O\left(\frac{1}{N^2}\right), \quad \mu = \frac{\bar{\mu}}{N} + O\left(\frac{1}{N^2}\right), \quad r = \frac{\bar{r}}{N} + O\left(\frac{1}{N^2}\right), \tag{S4}$$

and consequently

$$h_a = \frac{\bar{h}_a}{N} + O\left(\frac{1}{N^2}\right), \tag{S5}$$

where $\bar{r}$, $\bar{h}_a$, $\bar{s}_i$ and $\bar{\mu}$ are constants that are independent of $N$.

## 1.2 Conditional moment expansions

Here we present large-$N$ expansions for a generalized multi-locus model, considering the joint effects of selection, mutation, temporal sampling and recombination. Specifically, we will derive expansions for

$$p_a(\boldsymbol{z}(t), \Delta t) := \mathrm{E}\left[Z_a(t+\Delta t)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right], \quad \mathrm{Var}\left(Z_a(t+\Delta t)\Big|\boldsymbol{Z}(t)\right), \quad \mathrm{Covar}\left(Z_a(t+\Delta t), Z_b(t+\Delta t)\Big|\boldsymbol{Z}(t)\right). \tag{S6}$$

Before presenting these results, given below in (S10), (S14), and (S15) respectively, we require some preliminary expansions that apply for $\Delta t = 1$.

6

### 1.2.1  Preliminary expansions for $\Delta t = 1$

We start with an expansion for $p_a(\boldsymbol{z}(t)) \equiv p_a(\boldsymbol{z}(t), 1)$. From (S4), note that

$$y_a(t) = z_a(t) - r(L-1)\big(z_a(t) - \psi_a(\boldsymbol{z}(t))\big) + O\left(\frac{1}{N^2}\right)$$

$$= z_a(t) + O\left(\frac{1}{N}\right) .$$

Applying this in (S3), together with (S4) and (S5), it leads after some basic algebra to

$$p_a(\boldsymbol{z}(t)) = z_a(t)\left(1 + h_a - \sum_{b=1}^{M} h_b z_b(t)\right) + \mu\left(-L z_a(t) + \sum_{b=1, d_{ab}=1}^{M} z_b(t)\right) - r(L-1)\big(z_a(t) - \psi_a(\boldsymbol{z}(t))\big) + O\left(\frac{1}{N^2}\right)$$

$$= z_a(t) + O\left(\frac{1}{N}\right) . \tag{S7}$$

The following expansions also hold:

$$\mathrm{E}\left[Z_a^2(t+1)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right] = z_a^2(t) + O\left(\frac{1}{N}\right) \tag{S8}$$

$$\mathrm{E}\left[Z_a(t+1)Z_b(t+1)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right] = z_a(t)z_b(t) + O\left(\frac{1}{N}\right)$$

$$\mathrm{Var}\left[Z_a(t+1)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right] = \frac{z_a(t)(1 - z_a(t))}{N} + O\left(\frac{1}{N^2}\right)$$

$$\mathrm{Var}\left[Z_a(t+1)Z_b(t+1)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right] = O\left(\frac{1}{N}\right) .$$

These are obtained by computing the conditional moment generating function,

$$M(\boldsymbol{w}, \boldsymbol{z}(t)) = \mathrm{E}\left[\exp\left(\boldsymbol{w}^T \boldsymbol{Z}(t+1)\right)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right]$$

$$= \left(\sum_{a=1}^{M} p_a(\boldsymbol{z}(t)) \exp\left(\frac{w_a}{N}\right)\right)^N ,$$

where $\boldsymbol{w} = (w_1, \ldots, w_M)$, evaluating relevant derivatives at zero in the standard way, and taking $N$ large using (S7).

### 1.2.2  Main derivations for arbitrary $\Delta t$

We first present an expansion for $p_a(\boldsymbol{z}(t), \Delta t)$ in (S6). This is obtained by applying the law of total expectation and the Markov property of the WF process, giving

$$p_a(\boldsymbol{z}(t), \Delta t) = \mathrm{E}\left[\mathrm{E}\left[Z_a(t + \Delta t)\Big|\boldsymbol{Z}(t + \Delta t - 1)\right]\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right]$$

where the innermost conditional expectation is recognized as $p_a(\boldsymbol{z}(t + \Delta t - 1))$. Hence, using the first line of (S7), we get

$$p_a(\boldsymbol{z}(t), \Delta t) = p_a(\boldsymbol{z}(t), \Delta t - 1)(1 + h_a) + \mu\left(-L p_a(\boldsymbol{z}(t), \Delta t - 1) + \sum_{b=1, d_{ab}=1}^{M} p_b(\boldsymbol{z}(t), \Delta t - 1)\right) \tag{S9}$$

$$- \sum_{b=1}^{M} h_b \mathrm{E}\left[Z_a(t + \Delta t - 1)Z_b(t + \Delta t - 1)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right]$$

$$- r(L-1)\left(p_a(\boldsymbol{z}(t), \Delta t - 1) - \mathrm{E}\left[\psi_a(\boldsymbol{Z}(t + \Delta t - 1))\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right]\right) + O\left(\frac{1}{N^2}\right) .$$

For the remaining expectations, repeated application of the law of total expectation, together with (S8) and (S1), leads to:

$$\mathrm{E}\left[Z_a(t + \Delta t - 1)Z_b(t + \Delta t - 1)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right] = z_a(t)z_b(t) + O\left(\frac{1}{N}\right)$$

$$\mathrm{E}\left[\psi_a(\boldsymbol{Z}(t + \Delta t - 1))\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right] = \psi_a(\mathbf{z}(t)) + O\left(\frac{1}{N}\right) .$$

Hence, plugging these into (S9), we have expressed $p_a(\boldsymbol{z}(t), \Delta t)$ in terms of $p_a(\boldsymbol{z}(t), \Delta t - 1)$. This recursion may be iterated $\Delta t$ times, yielding the desired expansion:

$$p_a(\boldsymbol{z}(t), \Delta t) \tag{S10}$$

$$= z_a(t) + \Delta t \left( z_a(t) \left( h_a - \sum_{b=1}^{M} h_b z_b(t) \right) + \mu \left( -L z_a(t) + \sum_{b=1, d_{ab}=1}^{M} z_b(t) \right) - r(L-1)\big(z_a(t) - \psi_a(\boldsymbol{z}(t))\big) \right) + O\left(\frac{1}{N^2}\right).$$

Now consider the variance and covariance terms in (S6). For the variance term, using the law of total variance and the Markov property of the WF process, we have

$$\mathrm{Var}\left( Z_a(t + \Delta t) \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right) = \mathrm{E}\left[ \mathrm{Var}\left( Z_a(t + \Delta t) \bigg| \boldsymbol{Z}(t + \Delta t - 1) \right) \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right]$$

$$+ \mathrm{Var}\left( \mathrm{E}\left[ Z_a(t + \Delta t) \bigg| \boldsymbol{Z}(t + \Delta t - 1) \right] \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right). \tag{S11}$$

The first term admits

$$\mathrm{E}\left[ \mathrm{Var}\left( Z_a(t + \Delta t) \bigg| \boldsymbol{Z}(t + \Delta t - 1) \right) \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right] = \mathrm{E}\left[ \frac{\boldsymbol{Z}(t + \Delta t - 1)(1 - \boldsymbol{Z}(t + \Delta t - 1))}{N} \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right] + O\left(\frac{1}{N^2}\right)$$

$$= \frac{z_a(t)(1 - z_a(t))}{N} + O\left(\frac{1}{N^2}\right) \tag{S12}$$

where the first line follows from (S8), the second line from repeated application of the law of total expectation.

The second term admits

$$\mathrm{Var}\left( \mathrm{E}\left[ Z_a(t + \Delta t) \bigg| \boldsymbol{Z}(t + \Delta t - 1) \right] \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right) = \mathrm{Var}\left( Z_a(t + \Delta t - 1) \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right) + O\left(\frac{1}{N^2}\right) \tag{S13}$$

which follows from the first line in (S7), recalling again (S4), and using the fact that $\mathrm{Var}\left(\sum_{i=1}^{n} a_i X_i\right) = \sum_{i=1}^{n} a_i^2 \mathrm{Var}(X_i) + 2 \sum_{1 \le i \le j \le n} a_i a_j \mathrm{Cov}(X_i, X_j)$ for arbitrary random variables $\{X_i\}_{i=1}^{n}$.

Substituting (S12) and (S13) into (S11), and iterating, we arrive at the desired expansion

$$\mathrm{Var}\left( Z_a(t + \Delta t) \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right) = \Delta t \frac{z_a(t)(1 - z_a(t))}{N} + O\left(\frac{1}{N^2}\right). \tag{S14}$$

For the covariance term in (S6), a similar procedure is applied. For $a \ne b$, this produces

$$\mathrm{Covar}\left( Z_a(t + \Delta t), Z_b(t + \Delta t) \bigg| \boldsymbol{Z}(t) = \boldsymbol{z}(t) \right) = -\Delta t \frac{z_a(t) z_b(t)}{N} + O\left(\frac{1}{N^2}\right). \tag{S15}$$

## 1.3   Genotype-level path integral

Here, making use of the moment expansions just derived, we present a diffusion approximation and path integral expression for the stochastic genotype frequency dynamics under our generalized WF model setting. We start by recalling some standard concepts from diffusion process theory, which provides a continuous-time continuous-frequency approximation to the discrete-time discrete-frequency WF process. In our setting, this approximation is valid under the large-$N$ parameter scalings (S4) and (S5), and corresponds to the continuous process

$$\check{\boldsymbol{Z}}(\tau) = \big( \check{Z}_1(\tau), \ldots, \check{Z}_M(\tau) \big) \coloneqq \boldsymbol{Z}(\lfloor N\tau \rfloor), \quad \tau \ge 0 \tag{S16}$$

taken in the limit $N \to \infty$, where $\lfloor \cdot \rfloor$ denotes the floor function. The variable $\tau$ is a continuous time variable with units of $N$ generations, with one generation in discrete time (i.e., from $t$ to $t + 1$) thus taking

$$\delta\tau = \frac{1}{N} \tag{S17}$$

continuous time units. The diffusion process is described by the probability density function $\phi$, the solution to

$$\frac{\partial \phi}{\partial \tau} = \left[ -\sum_{a=1}^{M} \frac{\partial}{\partial \check{z}_a} \bar{d}_a(\check{\boldsymbol{z}}(\tau)) + \sum_{a=1}^{M} \sum_{b=1}^{M} \frac{\partial}{\partial \check{z}_a} \frac{\partial}{\partial \check{z}_b} \bar{C}_{ab}(\check{\boldsymbol{z}}(\tau)) \right] \phi. \tag{S18}$$

8

This is fully determined by the drift vector $\bar{d}(\check{\boldsymbol{z}}(\tau))$, which describes the rate of expected changes in genotype frequencies at time $\tau$ (see Eq. 4.99 of Risken[7]), and the diffusion matrix $\bar{C}(\check{\boldsymbol{z}}(\tau))$, which describes the scaled covariance of the genotype frequency changes (see equation 4.100 of Risken[7]).

Under our general model setting, where we assume genotype frequencies are observed at generations $t$ and $t + \Delta t$ with no observations in-between, the drift vector has $a$th entry

$$
\begin{aligned}
\bar{d}_a(\check{\boldsymbol{z}}(\tau)) &= \lim_{\delta\tau\Delta t \to 0} \frac{1}{\delta\tau\Delta t} \mathrm{E}\left[\left(\check{Z}_a\left(\tau + \delta\tau\Delta t\right) - \check{Z}_a(\tau)\right)\bigg|\check{\boldsymbol{Z}}(\tau) = (\check{z}_1(\tau), \ldots, \check{z}_M(\tau))\right] \\
&= \lim_{N\to\infty} N \left(\check{z}_a(\tau)\left(h_a - \sum_{b=1}^{M} h_b\check{z}_b(\tau)\right) + \mu\left(-L\check{z}_a(\tau) + \sum_{b=1,d_{ab}=1}^{M}\check{z}_b(\tau)\right) - r(L-1)\left(\check{z}_a(\tau) - \psi_a\left(\check{\boldsymbol{z}}(\tau)\right)\right)\right) \\
&= \check{z}_a(\tau)\left(\bar{h}_a - \sum_{b=1}^{M} \bar{h}_b\check{z}_b(\tau)\right) + \bar{\mu}\left(-L\check{z}_a(\tau) + \sum_{b=1,d_{ab}=1}^{M}\check{z}_b(\tau)\right) - \bar{r}(L-1)\left(\check{z}_a(\tau) - \psi_a\left(\check{\boldsymbol{z}}(\tau)\right)\right),
\end{aligned}
$$

where the second line follows from (S10), (S16) and (S17), while the diffusion matrix has $(a,b)$th entry

$$
\begin{aligned}
\bar{C}_{ab}(\check{\boldsymbol{z}}(\tau)) &= \frac{1}{2}\lim_{\delta\tau\Delta t \to 0}\frac{1}{\delta\tau\Delta t}\mathrm{E}\left[\left(\left(\check{Z}_a\left(\tau + \delta\tau\Delta t\right) - \check{Z}_a(\tau)\right)\right)\left(\left(\check{Z}_b\left(\tau + \delta\tau\Delta t\right) - \check{Z}_b(\tau)\right)\right)\bigg|\check{\boldsymbol{Z}}(\tau) = (\check{z}_1(\tau), \ldots, \check{z}_M(\tau))\right] \\
&= \frac{1}{2}\lim_{\delta\tau\Delta t \to 0}\frac{1}{\delta\tau\Delta t}\mathrm{Covar}\left(\check{Z}_a\left(\tau + \delta\tau\Delta t\right), \check{Z}_b\left(\tau + \delta\tau\Delta t\right)\bigg|\check{\boldsymbol{Z}}(\tau) = (\check{z}_1(\tau), \ldots, \check{z}_M(\tau))\right) \\
&= \frac{1}{2}\begin{cases}\check{z}_a(\tau)(1 - \check{z}_a(\tau)) & a = b \\ -\check{z}_a(\tau)\check{z}_b(\tau) & a \neq b\end{cases},
\end{aligned}
$$

following from (S14) and (S15).

The path integral approach (see equation 4.109 of Risken[7]) approximates the transition density of a diffusion process over a small time period. Applying this approach to the process described by (S18), for small $\delta\tau\Delta t$ (equivalently large $\frac{N}{\Delta t}$), we obtain for the transition probability density over $\Delta t$ generations,

$$
\begin{aligned}
&\phi(\check{\boldsymbol{z}}(\tau + \delta\tau\Delta t)|\check{\boldsymbol{z}}(\tau), N, \mu, \boldsymbol{h}) \\
&\approx \frac{\exp\left(-\frac{1}{4\delta\tau\Delta t}\left(\check{\boldsymbol{z}}(\tau + \delta\tau\Delta t) - \check{\boldsymbol{z}}(\tau) - \bar{d}(\check{\boldsymbol{z}}(\tau))\delta\tau\Delta t\right)^{\mathrm{T}}\bar{C}(\check{\boldsymbol{z}}(\tau))^{-1}\left(\check{\boldsymbol{z}}(\tau + \delta\tau\Delta t) - \check{\boldsymbol{z}}(\tau) - \bar{d}(\check{\boldsymbol{z}}(\tau))\delta\tau\Delta t\right)\right)}{(4\pi\delta\tau\Delta t)^{M/2}\sqrt{\det(\bar{C}(\check{\boldsymbol{z}}(\tau)))}}.
\end{aligned} \tag{S19}
$$

Note that here we have shown the explicit dependence on $N$, $\mu$ and $\boldsymbol{h}$. From this result, and recalling (S17), the transition probability from time $t_k$ to $t_{k+1}$ of the original discrete-time discrete-frequency WF process can (for large $\frac{N}{\Delta t}$) be approximated by

$$
\begin{aligned}
&P\left(\boldsymbol{z}(t_{k+1})|\boldsymbol{z}(t_k), N, \mu, \boldsymbol{h}\right) \\
&\approx \phi(\boldsymbol{z}(t_{k+1})|\boldsymbol{z}(t_k), N, \mu, \boldsymbol{h})\mathrm{d}\boldsymbol{z}(t_{k+1}) \\
&= \left(\frac{N}{2\pi\Delta t_k}\right)^{M/2}\frac{\mathrm{d}\boldsymbol{z}(t_{k+1})}{\sqrt{\det C(\boldsymbol{z}(t_k))}} \\
&\quad \times \exp\left(-\frac{N}{2\Delta t_k}\sum_{a=1}^{M}\sum_{b=1}^{M}\left[z_a(t_{k+1}) - z_a(t_k) - d_a(\boldsymbol{z}(t_k))\Delta t_k\right]\left(C^{-1}(\boldsymbol{z}(t_k))\right)_{ab}\left[z_b(t_{k+1}) - z_b(t_k) - d_b(\boldsymbol{z}(t_k))\Delta t_k\right]\right),
\end{aligned}
$$

where $\Delta t_k = t_{k+1} - t_k$, the $\mathrm{d}\boldsymbol{z}(t_{k+1}) = \prod_{a=1}^{M}\mathrm{d}z_a(t_{k+1})$ represent small frequency differences accounting for the quantization of the continuous genotype frequency space, and where we have defined $d_a(\boldsymbol{z}(t_k)) := \frac{\bar{d}_a(\boldsymbol{z}(t_k))}{N}$ and

$$
(C(\boldsymbol{z}(t_k)))_{ab} := 2(\bar{C}(\boldsymbol{z}(t_k)))_{ab}. \tag{S20}
$$

The path integral expression for the probability of observing a trajectory of genotype frequencies $(\boldsymbol{z}(t_1), \boldsymbol{z}(t_2), \ldots, \boldsymbol{z}(t_K))$ is then given by

$$
\begin{aligned}
P\left((\boldsymbol{z}(t_k))_{k=1}^{K}|\boldsymbol{z}(t_0), N, \mu, \boldsymbol{h}\right) &= \prod_{k=0}^{K-1} P\left(\boldsymbol{z}(t_{k+1})|\boldsymbol{z}(t_k), N, \mu, \boldsymbol{h}\right) \\
&\approx \left(\prod_{k=0}^{K-1}\left[\frac{1}{\sqrt{\det C(\boldsymbol{z}(t_k))}}\left(\frac{N}{2\pi\Delta t_k}\right)^{M/2}\mathrm{d}\boldsymbol{z}(t_{k+1})\right]\right)\exp\left(-\frac{N}{2}S\left((\boldsymbol{z}(t_k))_{k=0}^{K}\right)\right)
\end{aligned} \tag{S21}
$$

where

$$S\left((\boldsymbol{z}(t_k))_{k=0}^K\right) = \sum_{k=0}^{K-1} \frac{1}{\Delta t_k} \sum_{a=1}^M \sum_{b=1}^M \left[z_a(t_{k+1}) - z_a(t_k) - \Delta t_k d_a(\boldsymbol{z}(t_k))\right] \left(C^{-1}(\boldsymbol{z}(t_k))\right)_{ab} \left[z_b(t_{k+1}) - z_b(t_k) - \Delta t_k d_b(\boldsymbol{z}(t_k))\right] .$$

## 1.4 Mutant allele-level path integral

Based on the genotype transition probability density (S19), we now provide an approximation for the mutant allele transition probability density. Let $\boldsymbol{x}(t) = (x_1(t), \ldots, x_L(t))$, for which

$$x_i(t) = \sum_{a=1}^M g_i^a z_a(t) \tag{S22}$$

is the mutant frequency at locus $i$ during generation $t$. Also, define the random mutant allele frequency vector $\boldsymbol{X}(t) = (X_1(t), \ldots, X_L(t))$, related to the random genotype frequency vector by

$$X_i(t) = \sum_{a=1}^M g_i^a Z_a(t).$$

The observed frequency vector $\boldsymbol{x}(t)$ is thus a realization of this random vector. The continuous process which characterizes the mutant allele frequencies is

$$\check{\boldsymbol{X}}(\tau) = \left(\check{X}_1(\tau), \ldots, \check{X}_L(\tau)\right) := \boldsymbol{X}(\lfloor N\tau \rfloor), \quad \tau \geq 0$$

taken as $N \to \infty$, which satisfies

$$\check{X}_i(\tau) = \sum_{a=1}^M g_i^a \check{Z}_a(\tau).$$

This allele frequency process is a diffusion process, whose probability density evolution is described by a solution to

$$\frac{\partial \phi}{\partial \tau} = \left[ -\sum_{i=1}^L \frac{\partial}{\partial \check{x}_i} \underline{d}_i(\check{\boldsymbol{x}}(\tau)) + \sum_{i=1}^L \sum_{j=1}^L \frac{\partial}{\partial \check{x}_i} \frac{\partial}{\partial \check{x}_j} \underline{C}_{ij}(\check{\boldsymbol{x}}(\tau)) \right] \phi, \tag{S23}$$

again fully characterized by the drift vector $\underline{d}(\check{\boldsymbol{x}}(\tau))$ and diffusion matrix $\underline{C}(\check{\boldsymbol{x}}(\tau))$. (Note that the drift and the diffusion equations given in the Methods section of the main text are the same as those given here; however, in the Methods we do not introduce notation to explicitly distinguish between continuous-time and discrete-time processes, for simplicity.)

For the genotype case, the transition probability density was approximated to have a Gaussian form (S19). As the mutant allele frequencies are linear combinations of the genotype frequencies (S22), this implies that the transition probability density of mutant alleles also has a Gaussian form. Therefore, the allele-level drift and diffusion terms in (S23) are also a linear combination of the genotype drift and diffusion terms. The drift vector thus has $i$th entry

$$\begin{aligned}
\underline{d}_i(\check{\boldsymbol{x}}(\tau)) &:= \sum_{a=1}^M g_i^a \bar{d}_a(\check{\boldsymbol{z}}(\tau)) \\
&= \sum_{a=1}^M g_i^a \left( \check{z}_a(\tau) \left( \bar{h}_a - \sum_{b=1}^M \bar{h}_b \check{z}_b(\tau) \right) + \bar{\mu} \left( -L\check{z}_a(\tau) + \sum_{b=1,d_{ab}=1}^M \check{z}_b(\tau) \right) - \bar{r}(L-1)\left( \check{z}_a(\tau) - \psi_a\left(\check{\boldsymbol{z}}(\tau)\right)\right) \right) \\
&= \check{x}_i(\tau)\left(1 - \check{x}_i(\tau)\right) \bar{s}_i + \sum_{j=1,j\neq i}^L \left(\check{x}_{ij}(\tau) - \check{x}_i(\tau)\check{x}_j(\tau)\right) \bar{s}_j + \bar{\mu}(1 - 2\check{x}_i(\tau)),
\end{aligned} \tag{S24}$$

where we have defined

$$\check{x}_{ij}(\tau) := x_{ij}(\lfloor N\tau \rfloor), \quad \tau \geq 0.$$

This result follows by applying (S2), noting

$$\sum_{a=1}^M g_i^a g_j^a \check{z}_a(\tau) = \check{x}_{ij}(\tau),$$

10

and recognizing that $\sum_{a=1}^{M} g_i^a \left( \check{z}_a(\tau) - \psi_a \left( \check{\boldsymbol{z}}(\tau) \right) \right) = 0$. Establishing this latter relation, which accounts for the effect of recombination, is non-trivial and requires algebraic development. To this end, let us first define

$$\theta_i^{cd} := \sum_{a=1}^{M} g_i^a R_{a,cd},$$

which is the probability that genotypes $c$ and $d$ recombine to form a genotype which has a mutation at locus $i$. Since the model is biallelic, the recombination event could involve allele pairs $(0,0)$, $(0,1)$, $(1,0)$, or $(1,1)$ at locus $i$ of genotypes $c$ and $d$. We thus have

$$
\begin{aligned}
\sum_{a=1}^{M} g_i^a \psi_a \left( \check{\boldsymbol{z}}(\tau) \right) &= \sum_{a=1}^{M} g_i^a \sum_{c=1}^{M} \sum_{d=1}^{M} R_{a,cd} \check{z}_c(\tau) \check{z}_d(\tau) \\
&= \sum_{c=1}^{M} \sum_{d=1}^{M} \theta_i^{cd} \check{z}_c(\tau) \check{z}_d(\tau) \\
&= \sum_{c=1}^{M} \left( \sum_{d=1}^{M} \theta_i^{cd} g_i^c g_i^d \check{z}_c(\tau) \check{z}_d(\tau) + \sum_{d=1}^{M} \theta_i^{cd} g_i^c (1 - g_i^d) \check{z}_c(\tau) \check{z}_d(\tau) \right) \\
&\quad + \sum_{c=1}^{M} \left( \sum_{d=1}^{M} \theta_i^{cd} (1 - g_i^c) g_i^d \check{z}_c(\tau) \check{z}_d(\tau) + \sum_{d=1}^{M} \theta_i^{cd} (1 - g_i^c)(1 - g_i^d) \check{z}_c(\tau) \check{z}_d(\tau) \right).
\end{aligned}
$$

Now by noting that

$$\theta_i^{cd} g_i^c g_i^d = g_i^c g_i^d$$

$$\theta_i^{cd} g_i^c (1 - g_i^d) = \frac{1}{2} g_i^c (1 - g_i^d)$$

$$\theta_i^{cd} (1 - g_i^c) g_i^d = \frac{1}{2} (1 - g_i^c) g_i^d$$

$$\theta_i^{cd} (1 - g_i^c)(1 - g_i^d) = 0$$

where the factor of $\frac{1}{2}$ arises because there is a 50% chance that genotype $c$ $(d)$ with a mutant at locus $i$ and genotype $d$ $(c)$ with WT at locus $i$ will recombine to a genotype with a mutant at locus $i$, we have

$$
\begin{aligned}
\sum_{a=1}^{M} g_i^a \psi_a \left( \check{\boldsymbol{z}}(\tau) \right) &= \sum_{c=1}^{M} \left( \sum_{d=1}^{M} g_i^c g_i^d \check{z}_c(\tau) \check{z}_d(\tau) + \frac{1}{2} \sum_{d=1}^{M} g_i^c (1 - g_i^d) \check{z}_c(\tau) \check{z}_d(\tau) \right) \\
&\quad + \frac{1}{2} \sum_{c=1}^{M} \sum_{d=1}^{M} (1 - g_i^c) g_i^d \check{z}_c(\tau) \check{z}_d(\tau) \\
&= \check{x}_i^2(\tau) + \frac{1}{2} \check{x}_i(\tau)(1 - x_i(\tau)) + \frac{1}{2} \check{x}_i(\tau)(1 - \check{x}_i(\tau)) \\
&= \check{x}_i(\tau),
\end{aligned}
$$

thus implying that $\sum_{a=1}^{M} g_i^a \left( \check{z}_a(\tau) - \psi_a \left( \check{\boldsymbol{z}}(\tau) \right) \right) = 0$, the quoted property.

(Note that above and also in the subsequent derivations, since we are focused on the mutant allele frequency dynamics, we adopt notation which only explicitly demonstrates dependencies on the mutant allele frequencies. It should be recognized, however, that these dynamics also depend on the pairwise variant frequencies. We do not explicitly show this for the sake of notational convenience.)

Now consider the diffusion matrix. This has $(i,j)$th entry

$$
\begin{aligned}
\underline{C}_{ij}(\check{\boldsymbol{x}}(\tau)) &:= \sum_{a=1}^{M} \sum_{b=1}^{M} g_i^a g_j^b \bar{C}_{ab}(\check{\boldsymbol{z}}(\tau)) \\
&= \frac{1}{2} \sum_{a=1}^{M} g_i^a g_j^a \check{z}_a(\tau) - \frac{1}{2} \left( \sum_{a=1}^{M} g_i^a \check{z}_a(\tau) \right) \left( \sum_{b=1}^{M} g_j^b \check{z}_b(\tau) \right) \\
&= \frac{1}{2} \left( \check{x}_{ij}(\tau) - \check{x}_i(\tau) \check{x}_j(\tau) \right).
\end{aligned}
$$

Note that if we concatenate the infinitesimal drift for all loci in vector form, it yields

$$\underline{d}(\check{\boldsymbol{x}}(\tau)) = 2\underline{C}(\check{\boldsymbol{x}}(\tau))\bar{\boldsymbol{s}} + \bar{\mu} \left( \mathbf{1} - 2\check{\boldsymbol{x}}(\tau) \right) \tag{S25}$$

11

where $\mathbf{1}$ is the vector of ones.

As for the genotype case, Eq. 4.109 of Risken [7] may be applied, which approximates the transition probability density of this diffusion process over $\Delta t$ generations by a Gaussian distribution, given as

$$\phi(\check{\boldsymbol{x}}(\tau + \delta\tau\Delta t)|\check{\boldsymbol{x}}(\tau)) \approx \frac{\exp\left(-\frac{1}{4\delta\tau\Delta t}\left(\check{\boldsymbol{x}}(\tau + \delta\tau\Delta t) - \check{\boldsymbol{x}}(\tau) - \underline{d}(\check{\boldsymbol{x}}(\tau))\delta\tau\Delta t\right)^{\mathrm{T}}[\underline{C}(\check{\boldsymbol{x}}(\tau))]^{-1}\left(\check{\boldsymbol{x}}(\tau + \delta\tau\Delta t) - \check{\boldsymbol{x}}(\tau) - \underline{d}(\check{\boldsymbol{x}}(\tau))\delta\tau\Delta t\right)\right)}{(4\pi\delta\tau\Delta t)^{L/2}\sqrt{\det(\underline{C}(\check{\boldsymbol{x}}(\tau)))}}$$

From this result, and recalling (S17), the transition probability from time $t_k$ to $t_{k+1}$ of the original discrete-time discrete-frequency WF process can (for large $\frac{N}{\Delta t_k}$) be approximated by

$$P(\boldsymbol{x}(t_{k+1})|\boldsymbol{x}(t_k)) \approx \phi(\boldsymbol{x}(t_{k+1})|\boldsymbol{x}(t_k))\mathrm{d}\boldsymbol{x}(t_{k+1})$$

$$= \frac{\left(\frac{N}{2\pi\Delta t_k}\right)^{L/2}\mathrm{d}\boldsymbol{x}(t_{k+1})}{\sqrt{\det C(\boldsymbol{x}(t_k))}}$$

$$\times \exp\left(-\frac{N}{2\Delta t_k}\sum_{i=1}^{L}\sum_{j=1}^{L}\left[x_i(t_{k+1}) - x_i(t_k) - d_i(\boldsymbol{x}(t_k))\Delta t_k\right]\left(C^{-1}(\boldsymbol{x}(t_k))\right)_{ij}\left[x_j(t_{k+1}) - x_j(t_k) - d_j(\boldsymbol{x}(t_k))\Delta t_k\right]\right)$$

where the $\mathrm{d}\boldsymbol{x}(t_{k+1}) = \prod_{i=1}^{L}\mathrm{d}x_i(t_{k+1})$ represents small frequency differences accounting for the quantization of the continuous mutant allele frequency space, and we have made the replacements $\underline{d}_i(\boldsymbol{x}(t)) = Nd_i(\boldsymbol{x}(t))$ and $(\underline{C}(\boldsymbol{x}(t)))_{ij} = (1/2)C(\boldsymbol{x}(t))_{ij}$. The path integral expression then follows by noting that the probability of observing a trajectory of mutant allele frequencies $(\boldsymbol{x}(t_1), \boldsymbol{x}(t_2), \ldots, \boldsymbol{x}(t_K))$ is given by

$$P\left((\boldsymbol{x}(t_k))_{k=1}^{K}|\boldsymbol{x}(t_0), N, \mu, \boldsymbol{s}\right) \approx \left(\prod_{k=0}^{K-1}\frac{1}{\sqrt{\det C(\boldsymbol{x}(t_k))}}\left(\frac{N}{2\pi\Delta t_k}\right)^{L/2}\prod_{i=1}^{L}\mathrm{d}x_i(t_{k+1})\right)\left(\prod_{k=0}^{K-1}\exp\left(-\frac{N}{2}S\left((\boldsymbol{x}(t_k))_{k=0}^{K}\right)\right)\right),$$

(S26)

where

$$S\left((\boldsymbol{x}(t_k))_{k=0}^{K}\right) = \sum_{k=0}^{K-1}\frac{1}{\Delta t_k}\sum_{i=1}^{L}\sum_{j=1}^{L}\left[x_i(t_{k+1}) - x_i(t_k) - \Delta t_k d_i(\boldsymbol{x}(t_k))\right]\left(C^{-1}(\boldsymbol{x}(t_k))\right)_{ij}\left[x_j(t_{k+1}) - x_j(t_k) - \Delta t_k d_j(\boldsymbol{x}(t_k))\right].$$

The quantity $S\left((\boldsymbol{x}(t_k))_{k=0}^{K}\right)$ is referred to as the "action" in physics.

## 1.5    The MPL estimator solution

We now present a proof of the MPL estimate, given in Eq. (11) of Methods. This is obtained as the MAP estimate of the selection coefficients, which is the solution to

$$\hat{\boldsymbol{s}} = \arg\max_{\boldsymbol{s}} \mathfrak{L}\left(\boldsymbol{s}|N, \mu, (\boldsymbol{x}(t_k))_{k=0}^{K}\right) P_{\mathrm{prior}}(\boldsymbol{s}),$$

where

$$\mathfrak{L}\left(\boldsymbol{s}|N, \mu, (\boldsymbol{x}(t_k))_{k=0}^{K}\right) = P\left((\boldsymbol{x}(t_k))_{k=1}^{K}|\boldsymbol{x}(t_0), N, \mu, \boldsymbol{s}\right) \tag{S27}$$

$$= \prod_{k=0}^{K-1} P\left(\boldsymbol{x}(t_{k+1})|\boldsymbol{x}(t_k), N, \mu, \boldsymbol{s}\right)$$

is the likelihood function, while $(\boldsymbol{x}(t_0), \boldsymbol{x}(t_1), \ldots, \boldsymbol{x}(t_K))$ is the observed trajectory of mutant allele frequencies.

The MAP problem is equivalent to

$$\hat{\boldsymbol{s}} = \arg\max_{\boldsymbol{s}}\left(\log \mathfrak{L}\left(\boldsymbol{s}|N, \mu, (\boldsymbol{x}(t_k))_{k=0}^{K}\right) + \log P_{\mathrm{prior}}(\boldsymbol{s})\right). \tag{S28}$$

The likelihood $\mathfrak{L}\left(\boldsymbol{s}|N, \mu, (\boldsymbol{x}(t_k))_{k=0}^{K}\right)$ is given by (S27) and is approximated by (S26). Assuming a conjugate-prior distribution

$$P_{\mathrm{prior}}(\boldsymbol{s}) = \frac{1}{(2\pi\sigma^2)^{L/2}}\exp\left(-\frac{1}{2\sigma^2}\boldsymbol{s}^{\mathrm{T}}\boldsymbol{s}\right),$$

we take the vector derivative of the right-hand side of (S28) with respect to $\boldsymbol{s}$ (see equation 10, Chapter 10.2.1 of Lutkepohl[11]), equate to zero, and then solve for $\boldsymbol{s}$. This produces the MPL estimator

$$\hat{\boldsymbol{s}} = (C_{\text{int}} + \gamma I)^{-1} (\Delta \boldsymbol{x} - \boldsymbol{\mu}_{\text{fl}})$$

where

$$C_{\text{int}} = \sum_{k=0}^{K-1} \Delta t_k C(\boldsymbol{x}(t_k)),$$

$$\Delta \boldsymbol{x} = \boldsymbol{x}(t_K) - \boldsymbol{x}(t_0),$$

$$\boldsymbol{\mu}_{\text{fl}} = \mu \sum_{k=0}^{K-1} \Delta t_k (1 - 2\boldsymbol{x}(t_k)) ,$$

where $\gamma = 1/N\sigma^2$, which can be viewed as a parameter that regularizes the covariance matrix prior to inversion.

The individual elements can be written as

$$\hat{s}_i = \sum_{j=1}^{L} \left[ \sum_{k=0}^{K-1} \Delta t_k C(\boldsymbol{x}(t_k)) + \gamma I \right]_{ij}^{-1} \left[ x_j(t_K) - x_j(t_0) - \mu \sum_{k=0}^{K-1} \Delta t_k (1 - 2x_j(t_k)) \right] \tag{S29}$$

for $i = 1, \ldots, L$.

## 1.6 Equivalence of genotype- and allele-level analyses

While the MPL estimator was derived using a path integral expression for the probability of mutant allele frequency trajectories, the WF evolutionary process is defined for genotypes. Hence, one may naturally ask whether there is any loss in information in considering only the marginal frequency dynamics. As we now show, the answer is no. Specifically, we show that the estimate of the selection coefficients based on genotype frequencies is equivalent to the derived MPL estimate.

We begin by defining $G$ as a $M \times L$ matrix with $(a, j)$th entry $G_{aj} = g_j^a$. Let $\boldsymbol{h} = (h_1, h_2, \ldots, h_M)$ denote the column vector of genotype selection coefficients, which relates to the mutant allele selection coefficients via

$$\boldsymbol{h} = G\boldsymbol{s}. \tag{S30}$$

Given the observed genotype frequencies $(\boldsymbol{z}(t_0), \boldsymbol{z}(t_1), \ldots, \boldsymbol{z}(t_K))$, the MAP estimator of the selection coefficients admits

$$\hat{\boldsymbol{s}} = \arg \max_{\boldsymbol{s}} \left( \log \mathfrak{L} \left( \boldsymbol{s} | \mu, N, (\boldsymbol{z}(t_k))_{k=0}^K \right) + \log P_{\text{prior}}(\boldsymbol{s}) \right),$$

where the likelihood function is given as

$$\mathfrak{L} \left( \boldsymbol{s} | N, \mu, (\boldsymbol{z}(t_k))_{k=0}^K \right) = P \left( (\boldsymbol{z}(t_k))_{k=1}^K | \boldsymbol{z}(t_0), N, \mu, \boldsymbol{h} \right) .$$

The right-hand side is approximated by (S21). Using (S30), in addition, gives the likelihood of the mutant allele selection coefficients $\boldsymbol{s}$ for an observed genotype frequency path and parameters $N$, $\mu$. Differentiating the resulting expression with respect to $\boldsymbol{s}$ and equating to zero leads to

$$\boldsymbol{0} = \sum_{k=0}^{K-1} \left( G^{\text{T}} \boldsymbol{z}(t_{k+1}) - G^{\text{T}} \boldsymbol{z}(t_k) - G^{\text{T}} C(\boldsymbol{z}(t_k)) G \boldsymbol{s} - \mu G^{\text{T}} E \boldsymbol{z}(t_k) - rL \left( G^{\text{T}} \boldsymbol{z}(t_k) - G^{\text{T}} \boldsymbol{\psi} (\boldsymbol{z}(t_k)) \right) \right) + \frac{1}{N\sigma^2} \boldsymbol{s}, \tag{S31}$$

where $\boldsymbol{\psi} (\boldsymbol{z}(t_k)) = \{\psi_a (\boldsymbol{z}(t_k))\}_{a=1}^M$ and $C(\boldsymbol{z}(t_k))$ is the genotype covariance matrix as given by (S20), i.e., with elements

$$C_{ab}(\boldsymbol{z}(t_k)) = \begin{cases} z_a(t_k)(1 - z_a(t_k)) & a = b \\ -z_a(t_k)z_b(t_k) & a \neq b, \end{cases}$$

while matrix $E$ has elements

$$E_{ab} = \begin{cases} -L & a = b \\ 0 & d_{ab} > 1 \\ 1 & d_{ab} = 1. \end{cases}$$

Noting that the relation between allele and genotype frequencies (S22) can be expressed in vector form as $\boldsymbol{x}(t_k) = G^{\mathrm{T}}\boldsymbol{z}(t_k)$, and that $C(\boldsymbol{x}(t_k)) = G^{\mathrm{T}}C(\boldsymbol{z}(t_k))G$ and $G^{\mathrm{T}}\boldsymbol{\psi}\left(\boldsymbol{z}(t_k)\right) = \boldsymbol{x}(t_k)$, we can solve (S31) to obtain the MAP estimate of the allele selection coefficients

$$\hat{\boldsymbol{s}} = \left[\sum_{k=0}^{K-1} \Delta t_k C(\boldsymbol{x}(t_k)) + \gamma I\right]^{-1} \left[\boldsymbol{x}(t_K) - \boldsymbol{x}(t_0) - \mu \sum_{k=0}^{K-1} \Delta t_k \left(\mathbf{1} - 2\boldsymbol{x}(t_k)\right)\right], \tag{S32}$$

where $\gamma = 1/N\sigma^2$. This is the same as the MPL estimator (S29).

This equivalence is important. It implies that under the additive fitness model, only the single and pairwise mutational frequencies are needed for optimally estimating the selection coefficients, and higher order information (e.g., three-locus mutational frequencies) are irrelevant. The same equivalences can also be shown for extensions of the MPL estimator presented in the following section. Further extensions will be explored in future work.

## 1.7  Extension to multiple alleles per locus and asymmetric mutation probabilities

Here we present an extension of our inference framework to consider a further generalization of the model to incorporate multiple alleles per locus and asymmetric mutation probabilities. The same notation and definitions introduced in Section 1 will also be used, unless stated otherwise, but will be interpreted in terms of the extended model. For example, $\check{\boldsymbol{Z}}(\tau)$ was introduced in (S16) to refer to the binary genotype continuous process with symmetric mutation probabilities. This notation will also be used here, but it will now refer to the non-binary genotype continuous process with asymmetric mutation probabilities. We first describe the model in detail.

We assume there are $\ell$ alleles per locus, thus resulting in $M = \ell^L$ genotypes, with the multi-allelic sequence for genotype $a$ denoted by $\boldsymbol{g}^a = (g_1^a, g_2^a, \ldots, g_L^a)$, where $g_i^a$ is the allele at locus $i$. Nucleotide sequences, for example, will have $g_i^a \in \{A, C, G, T\}$. For convenience, we denote each allele with an integer index, and consider the fitness of each genotype with respect to a reference sequence comprising of allele $\ell$ at each locus, i.e., the reference sequence is a sequence of $\ell$s. For example, if $\ell = 4$ and $L = 3$, then the reference sequence is represented by $(4, 4, 4)$. We assume, once again, an additive model of fitness, and that the reference genotype has a fitness of one. Under these assumptions, the selection coefficient for genotype $a$ is given by

$$h_a = f_a - 1 = \sum_{i=1}^{L} \sum_{\alpha=1}^{\ell-1} \delta(g_i^a, \alpha)s_{i,\alpha},$$

where $s_{i,\alpha}$ is the selection coefficient of allele $\alpha$ at locus $i$, and where $\delta(\cdot, \cdot)$ is the Kronecker-delta function,

$$\delta(x, y) := \begin{cases} 1 & x = y \\ 0 & \text{otherwise.} \end{cases}$$

The assumptions above imply that $s_{i,\ell} = 0$ for all $i = 1, \ldots, L$.

The allele frequency vector, describing the frequency of all alleles except (the reference) allele one, is given by $\boldsymbol{x}(t) = (x_{1,1}(t), \ldots, x_{1,\ell-1}(t), \ldots, x_{L,1}(t), \ldots, x_{L,\ell-1}(t))$, where $x_{i,\alpha}(t)$ denotes the observed frequency of allele $\alpha$ at locus $i$ during generation $t$, and is related to the genotype frequencies by

$$x_{i,\alpha}(t) = \sum_{a=1}^{M} \delta(g_i^a, \alpha)z_a(t). \tag{S33}$$

Finally, we denote $\mu_{\alpha\beta}$ as the mutation probability per generation from allele $\alpha$ to allele $\beta$ and $\mu_{ab}$ as the mutation probability per generation from genotype $a$ to genotype $b$. These are related by

$$\mu_{ab} = \prod_{i=1}^{L} \left(\sum_{\alpha=1}^{\ell} \sum_{\beta=1}^{\ell} \mu_{\alpha\beta}\delta(g_i^a, \alpha)\delta(g_i^b, \beta)\right).$$

Given this extended model, the MAP estimate of the selection coefficients is the solution to

$$\hat{\boldsymbol{s}} = \arg\max_{\boldsymbol{s}} \mathfrak{L}\left(\boldsymbol{s}|\boldsymbol{\mu}, N, (\boldsymbol{x}(t_k))_{k=0}^T\right) P_{\text{prior}}(\boldsymbol{s}), \tag{S34}$$

where

$$\mathfrak{L}\left(\boldsymbol{s}|\boldsymbol{\mu}, N, (\boldsymbol{x}(t_k))_{k=0}^T\right) = P\left((\boldsymbol{x}(t_k))_{k=1}^K|\boldsymbol{x}(t_0), N, \boldsymbol{\mu}, \boldsymbol{s}\right) \tag{S35}$$

$$= \prod_{k=0}^{K-1} P\left(\boldsymbol{x}(t_{k+1})|\boldsymbol{x}(t_k), N, \boldsymbol{\mu}, \boldsymbol{s}\right)$$

14

is the likelihood of the selection coefficients $\boldsymbol{s} = (s_{1,1}, \ldots, s_{1,\ell-1}, \ldots, s_{L,1}, \ldots, s_{L,\ell-1})$ and $P_{\mathrm{prior}}(\boldsymbol{s})$ is their prior distribution.

As with the simpler biallelic and symmetric mutation probability scenario, the main challenge in solving (S34) is that it requires computing the likelihood (S35), which is complicated. This is simplified as before, by adopting the path integral approach outlined in Section 1; but now extending the analysis to account for multiple alleles per locus and asymmetric mutation probabilities. Other than these differences, we may follow the same approach, starting by giving asymptotic moment expansions. In the following, since much of the development involves the same algebraic steps as before, albeit with additional notational record-keeping, we present mainly a summary of the key formulas here. We do, however, provide extensive details for the more involved derivations; namely, those of the drift and covariance terms for the allele-level dynamics.

By direct analogy to (S4) and (S5), we will work under the assumption that $N$ is large, and that as $N \to \infty$,

$$s_{i,\alpha} = \frac{\bar{s}_{i,\alpha}}{N} + O\Big(\frac{1}{N^2}\Big), \quad h_a = \frac{\bar{h}_a}{N} + O\Big(\frac{1}{N^2}\Big), \quad \mu_{\beta\alpha} = \frac{\bar{\mu}_{\beta\alpha}}{N} + O\Big(\frac{1}{N^2}\Big), \quad \mu_{ab} = \frac{\bar{\mu}_{ab}}{N} + O\Big(\frac{1}{N^2}\Big), \quad r = \frac{\bar{r}}{N} + O\Big(\frac{1}{N^2}\Big) \tag{S36}$$

where $\bar{s}_{i,\alpha}$, $\bar{h}_a$, $\bar{\mu}_{\beta\alpha}$, $\bar{\mu}_{ab}$ and $\bar{r}$ are constants independent of $N$.

### 1.7.1 Conditional moment expansions

The mean frequency of genotype $a$ at generation $t + 1$, conditioned on $\boldsymbol{Z}(t) = \boldsymbol{z}(t)$, admits

$$p_a(\boldsymbol{z}(t)) := \mathrm{E}\left[Z_a(t+1)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right] \tag{S37}$$

$$= \frac{(1 + h_a)y_a(t) + \sum\limits_{b=1}^{M} \left(\mu_{ba}(1 + h_b)y_b(t) - \mu_{ab}(1 + h_a)y_a(t)\right)}{\sum\limits_{b=1}^{M}(1 + h_b)y_b(t)},$$

where recall that

$$y_a(t) = (1 - r)^L z_a(t) + \left(1 - (1 - r)^L\right)\psi_a(\boldsymbol{z}(t)).$$

Utilizing (S37) along with (S36), and applying a similar proof to that described in Section 1.2, the mean frequency of genotype $a$ at generation $t + \Delta t$, conditioned on $\boldsymbol{Z}(t) = \boldsymbol{z}(t)$, then admits

$$p_a(\boldsymbol{z}(t), \Delta t) = z_a(t) + \Delta t\left(z_a(t)\left(h_a - \sum_{b=1}^{M} h_b z_b(t)\right)\right. \tag{S38}$$

$$\left. + \sum_{b=1, d_{ab}=1}^{M} (\mu_{ba}z_b(t) - \mu_{ab}z_a(t)) - r(L-1)\left(z_a(t) - \psi_a(\boldsymbol{z}(t))\right)\right) + O\left(\frac{1}{N^2}\right)$$

$$= z_a(t) + O\left(\frac{1}{N^2}\right),$$

where $d_{ab}$ here denotes the number of loci for which the allele identity at genotypes $a$ and $b$ differ, i.e.,

$$d_{ab} := \sum_{i=1}^{L} \left(1 - \delta(g_i^a, g_i^b)\right).$$

The corresponding expressions for the variance and covariance of the frequency of genotype $a$ at generation $t + \Delta t$, conditioned on $\boldsymbol{Z}(t) = \boldsymbol{z}(t)$ are obtained analogously to (S14) and (S15) as

$$\mathrm{Var}\left(Z_a\left(t + \Delta t\right)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right) = \Delta t\frac{z_a(t)(1 - z_a(t))}{N} + O\left(\frac{1}{N^2}\right), \tag{S39}$$

and

$$\mathrm{Covar}\left(Z_a\left(t + \Delta t\right), Z_b\left(t + \Delta t\right)\Big|\boldsymbol{Z}(t) = \boldsymbol{z}(t)\right) = -\Delta t\frac{z_a(t)z_b(t)}{N} + O\left(\frac{1}{N^2}\right). \tag{S40}$$

### 1.7.2 Path integral expression for the likelihood function

Based on the above conditional moment expansions, we can derive a path integral expression for the likelihood function (S35). Starting by considering genotype evolution, the drift vector and diffusion matrix are given respectively by

$$\bar{d}_a(\check{\boldsymbol{z}}(\tau), \Delta t) = \lim_{\delta\tau\Delta t \to 0} \frac{1}{\delta\tau\Delta t} \mathrm{E}\left[\check{Z}_a\left(\tau + \delta\tau\Delta t\right) - \check{Z}_a(\tau)\middle|\check{\boldsymbol{Z}}(\tau) = (\check{z}_1(\tau), \ldots, \check{z}_M(\tau))\right] \tag{S41}$$

$$= \check{z}_a(\tau)\left(\bar{h}_a - \sum_{b=1}^{M} \bar{h}_b \check{z}_b(\tau)\right) + \sum_{b=1, d_{ab}=1}^{M} (\bar{\mu}_{ba}\check{z}_b(\tau) - \bar{\mu}_{ab}\check{z}_a(\tau)) - \bar{r}(L-1)(\check{z}_a(\tau) - \psi_a(\check{\boldsymbol{z}}(\tau)))$$

and

$$\bar{C}_{ab}(\check{\boldsymbol{z}}(\tau), \Delta t) = \frac{1}{2}\begin{cases} \check{z}_a(\tau)(1 - \check{z}_a(\tau)) & a = b \\ -\check{z}_a(\tau)\check{z}_b(\tau) & a \neq b \end{cases}, \tag{S42}$$

which follow from (S38), (S39), (S40), and by adopting the procedure described in Section 1.3. From these results, and again following the procedure in Section 1.3, the conditional probability of observing a trajectory of genotype frequencies $(\boldsymbol{z}(t_1), \boldsymbol{z}(t_2), \ldots, \boldsymbol{z}(t_K))$ is obtained as

$$P\left((\boldsymbol{z}(t_k))_{k=1}^{K}|\boldsymbol{z}(t_0)\right) = \prod_{k=0}^{K-1} P\left(\boldsymbol{z}(t_{k+1})|\boldsymbol{z}(t_k)\right)$$

$$\approx \left(\prod_{k=0}^{K-1}\left[\frac{1}{\sqrt{\det \bar{C}(\boldsymbol{z}(t_k))}}\left(\frac{N}{4\pi\Delta t_k}\right)^{M/2}\prod_{a=1}^{M} \mathrm{d}z_a(t_{k+1})\right]\right)\exp\left(-\frac{N}{4}S\left((\boldsymbol{z}(t_k))_{k=0}^{K}\right)\right)$$

where $\Delta t_k = t_{k+1} - t_k$ and

$$S\left((\boldsymbol{z}(t_k))_{k=0}^{K}\right) = \sum_{k=0}^{K-1}\frac{1}{\Delta t_k}\sum_{a=1}^{M}\sum_{b=1}^{M}\left[z_a(t_{k+1}) - z_a(t_k) - \frac{\Delta t_k \bar{d}_a(\boldsymbol{z}(t_k))}{N}\right]$$

$$\times \left(\bar{C}^{-1}(\boldsymbol{z}(t_k))\right)_{ab}\left[z_b(t_{k+1}) - z_b(t_k) - \frac{\Delta t_k \bar{d}_b(\boldsymbol{z}(t_k))}{N}\right].$$

In this last equation we have dropped the second argument of the drift vector and diffusion matrix, as both quantities turn out to be independent of $\Delta t$, as seen from (S41) and (S42).

Now turn to mutant allele evolution. From the linearity of expectation and the genotype-to-allele mapping (S33), the mean frequency of allele $\alpha$ at locus $i$ during time $\tau + \delta\tau\Delta t$, conditioned on $\check{\boldsymbol{X}}(\tau) = \check{\boldsymbol{x}}(\tau)$, is given by

$$\sum_{a=1}^{M}\delta(g_i^a, \alpha)\check{z}_a(\tau) + \underline{d}_{i,\alpha}(\check{\boldsymbol{x}}(\tau))\delta\tau\Delta t = \check{x}_{i,\alpha}(\tau) + \underline{d}_{i,v}(\check{\boldsymbol{x}}(\tau))\delta\tau\Delta t$$

where

$$\underline{d}_{i,\alpha}(\check{\boldsymbol{x}}(\tau)) = \sum_{a=1}^{M}\delta(g_i^a, \alpha)\,\bar{d}_a(\check{\boldsymbol{z}}(\tau), \Delta t)$$

$$= \sum_{a=1}^{M}\delta(g_i^a, \alpha)\left(\check{z}_a(\tau)\left(\bar{h}_a - \sum_{b=1}^{M}\bar{h}_b\check{z}_b(\tau)\right) + \sum_{b=1, d_{ab}=1}^{M}(\bar{\mu}_{ba}\check{z}_b(\tau) - \bar{\mu}_{ab}\check{z}_a(\tau)) - \bar{r}(L-1)(\check{z}_a(\tau) - \psi_a(\check{\boldsymbol{z}}(\tau)))\right)$$

$$= I_1 + I_2 + I_3$$

with the second line following from (S41). For $I_1$, we have

$$I_1 = \sum_{a=1}^{M}\delta(g_i^a, \alpha)\check{z}_a(\tau)\left(\bar{h}_a - \sum_{b=1}^{M}\bar{h}_b\check{z}_b(\tau)\right)$$

$$= \sum_{a=1}^{M}\delta(g_i^a, \alpha)\check{z}_a(\tau)\sum_{j=1}^{L}\sum_{\beta=1}^{\ell-1}\delta(g_j^a, \beta)\bar{s}_{j,\beta} - \sum_{a=1}^{M}\delta(g_i^a, \alpha)\check{z}_a(\tau)\sum_{b=1}^{M}\sum_{j=1}^{L}\sum_{\beta=1}^{\ell-1}\delta(g_j^b, \beta)\bar{s}_{j,\beta}\check{z}_b(\tau)$$

$$= \sum_{j=1}^{L}\sum_{\beta=1}^{\ell-1}(\check{x}_{ij,\alpha\beta}(\tau) - \check{x}_{i,\alpha}(\tau)\check{x}_{j,\beta}(\tau))\bar{s}_{j,\beta}$$

where $\check{x}_{ij,\alpha\beta}(\tau)$ denotes the frequency of allele $\alpha$ and $\beta$ occurring respectively at loci $i$ and $j$ during time $\tau$, and given by

$$\check{x}_{ij,\alpha\beta}(\tau) = \sum_{a=1}^{M} \delta(g_i^a, \alpha)\delta(g_j^a, \beta)\check{z}_a(\tau).$$

For $I_2$, we have

$$I_2 = \sum_{a=1}^{M} \delta(g_i^a, \alpha) \sum_{b=1, d_{ab}=1}^{M} (\bar{\mu}_{ba}\check{z}_b(\tau) - \bar{\mu}_{ab}\check{z}_a(\tau))$$

$$= \sum_{a=1}^{M} \delta(g_i^a, \alpha) \sum_{b=1, d_{ab}=1}^{M} \left(\delta(g_i^b, \alpha) + 1 - \delta(g_i^b, \alpha)\right) (\bar{\mu}_{ba}\check{z}_b(\tau) - \bar{\mu}_{ab}\check{z}_a(\tau))$$

$$= I_{2a} + I_{2b}$$

where

$$I_{2a} = \sum_{a=1}^{M} \sum_{b=1, d_{ab}=1}^{M} \delta(g_i^a, \alpha)\delta(g_i^b, \alpha) (\bar{\mu}_{ba}\check{z}_b(\tau) - \bar{\mu}_{ab}\check{z}_a(\tau))$$

$$I_{2b} = \sum_{a=1}^{M} \sum_{b=1, d_{ab}=1}^{M} \delta(g_i^a, \alpha) \left(1 - \delta(g_i^b, \alpha)\right) (\bar{\mu}_{ba}\check{z}_b(\tau) - \bar{\mu}_{ab}\check{z}_a(\tau)).$$

Consider now $I_{2a}$, where we observe that (i) $\delta(g_i^a, \alpha)\delta(g_i^b, \alpha)$ is non-zero only when genotypes $a$ and $b$ both have allele $\alpha$ at locus $i$, while (ii) the summation $\sum_{b=1, d_{ab}=1}^{M}$ is over all genotypes $b$ which have a different allele from genotype $a$ at only one locus. These two observations imply that $\delta(g_i^a, \alpha)\delta(g_i^b, \alpha)$ is non-zero, for $b = 1, \ldots, M$, with $d_{ab} = 1$, only if $a$ and $b$ have a different allele at a single locus, but where the position of this locus is different from $i$. To illustrate this, consider a simple example with $L = 2$ and $\ell = 3$, in which case

$$\boldsymbol{g}^1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \boldsymbol{g}^2 = \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \quad \boldsymbol{g}^3 = \begin{pmatrix} 1 \\ 3 \end{pmatrix},$$

$$\boldsymbol{g}^4 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \quad \boldsymbol{g}^5 = \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \quad \boldsymbol{g}^6 = \begin{pmatrix} 2 \\ 3 \end{pmatrix},$$

$$\boldsymbol{g}^7 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, \quad \boldsymbol{g}^8 = \begin{pmatrix} 3 \\ 2 \end{pmatrix}, \quad \boldsymbol{g}^9 = \begin{pmatrix} 3 \\ 3 \end{pmatrix}.$$

Then if $i = 1$ and $\alpha = 1$, the only genotype-pairs $(a, b)$ which result in a non-zero $\delta(g_i^a, \alpha)\delta(g_i^b, \alpha)$ (subject to the condition $d_{ab} = 1$) are $(1, 2)$, $(2, 1)$, $(1, 3)$, $(3, 1)$, $(2, 3)$ and $(3, 2)$, as these genotype-pairs differ only at the second locus, while having allele one at locus $i = 1$.

Observe that every genotype-pair which result in a non-zero $\delta(g_i^a, \alpha)\delta(g_i^b, \alpha)$ (subject to the condition $d_{ab} = 1$) occur in conjugates, i.e., if $(1, 2)$ is such a genotype-pair, then so is $(2, 1)$. This implies that $I_{2a} = 0$, as for every genotype pair $(a, b)$ resulting in a non-zero $\delta(g_i^a, \alpha)\delta(g_i^b, \alpha)\bar{\mu}_{ba}\check{z}_b(\tau)$, there always exist one other conjugate pair $(b, a)$ which cancels this term through the $-\delta(g_i^a, \alpha)\delta(g_i^b, \alpha)\bar{\mu}_{ab}\check{z}_a(\tau)$ term.

Consider now $I_{2b}$, and observe that the $\delta(g_i^a, \alpha)\left(1 - \delta(g_i^b, \alpha)\right)$ term is non-zero only when genotype $a$, but not genotype $b$, has allele $\alpha$ at locus $i$. Combining the above, we can thus write $I_2$ as

$$I_2 = I_{2b} = \sum_{\beta=1, \beta \neq \alpha}^{\ell} (\bar{\mu}_{\beta\alpha}\check{x}_{i,\beta}(\tau) - \bar{\mu}_{\alpha\beta}\check{x}_{i,\alpha}(\tau))$$

$$= \sum_{\beta=1}^{\ell} (\bar{\mu}_{\beta\alpha}\check{x}_{i,\beta}(\tau) - \bar{\mu}_{\alpha\beta}\check{x}_{i,\alpha}(\tau))$$

$$= \sum_{\beta=1}^{\ell-1} \bar{\mu}_{\beta\alpha}\check{x}_{i,\beta}(\tau) + \bar{\mu}_{\ell\alpha} \left(1 - \sum_{\beta=1}^{\ell-1} \check{x}_{i,\beta}(\tau)\right) - \sum_{\beta=1}^{\ell} \bar{\mu}_{\alpha\beta}\check{x}_{i,\alpha}(\tau)$$

$$= \bar{\mu}_{\ell\alpha} + \sum_{\beta=1}^{\ell-1} (\bar{\mu}_{\beta\alpha} - \bar{\mu}_{\ell\alpha}) \check{x}_{i,\beta}(\tau) - \check{x}_{i,\alpha}(\tau) \sum_{\beta=1}^{\ell} \bar{\mu}_{\alpha\beta}.$$

17

In the third line above we have used the fact that the frequency of the reference allele is one minus the sum of the frequency of all other alleles.

Finally, we have

$$I_3 = -\bar{r}(L-1)\sum_{a=1}^{M}\delta(g_i^a,\alpha)\left(\check{z}_a(\tau) - \psi_a\left(\check{\boldsymbol{z}}(\tau)\right)\right).$$

Similar to the biallelic case, we define

$$\theta_{i,\alpha}^{cd} := \sum_{a=1}^{M}\delta(g_i^a,\alpha)R_{a,cd}$$

where $\theta_{i,\alpha}^{cd}$ is the probability that genotypes $c$ and $d$ recombine to form a genotype which has an allele $\alpha$ at locus $i$. Next we split this summation into four summations, one for each of the four possible recombination scenarios. Namely, when both genotypes $c$ and $d$ have allele $\alpha$ at locus $i$, when both do not have allele $\alpha$ at locus $i$ and when only of the genotypes has allele $\alpha$ at locus $i$. We thus have

$$\sum_{a=1}^{M}\delta(g_i^a,\alpha)\psi_a\left(\check{\boldsymbol{z}}(\tau)\right) = \sum_{a=1}^{M}\delta(g_i^a,\alpha)\sum_{c=1}^{M}\sum_{d=1}^{M}R_{a,cd}\check{z}_c(\tau)\check{z}_d(\tau)$$

$$= \sum_{c=1}^{M}\sum_{d=1}^{M}\theta_{i,\alpha}^{cd}\check{z}_c(\tau)\check{z}_d(\tau)$$

$$= \sum_{c=1}^{M}\left(\sum_{d=1}^{M}\theta_i^{cd}\delta(g_i^c,\alpha)\delta(g_i^d,\alpha)\check{z}_c(\tau)\check{z}_d(\tau) + \sum_{d=1}^{M}\theta_i^{cd}\delta(g_i^c,\alpha)(1-\delta(g_i^d,\alpha))\check{z}_c(\tau)\check{z}_d(\tau)\right)$$

$$+ \sum_{c=1}^{M}\left(\sum_{d=1}^{M}\theta_i^{cd}(1-\delta(g_i^c,\alpha))\delta(g_i^d,\alpha)\check{z}_c(\tau)\check{z}_d(\tau) + \sum_{d=1}^{M}\theta_i^{cd}(1-\delta(g_i^c,\alpha))(1-\delta(g_i^d,\alpha))\check{z}_c(\tau)\check{z}_d(\tau)\right).$$

Moreover, noting that

$$\theta_i^{cd}\delta(g_i^c,\alpha)\delta(g_i^d,\alpha) = \delta(g_i^c,\alpha)\delta(g_i^d,\alpha)$$

$$\theta_i^{cd}\delta(g_i^c,\alpha)(1-\delta(g_i^d,\alpha)) = \frac{1}{2}\delta(g_i^c,\alpha)(1-\delta(g_i^d,\alpha))$$

$$\theta_i^{cd}(1-\delta(g_i^c,\alpha))\delta(g_i^d,\alpha) = \frac{1}{2}(1-\delta(g_i^c,\alpha))\delta(g_i^d,\alpha)$$

$$\theta_i^{cd}(1-\delta(g_i^c,\alpha))(1-\delta(g_i^d,\alpha)) = 0$$

where the factor of $\frac{1}{2}$ arises because there is a 50% chance that genotype $c$ ($d$) with allele $v$ at locus $i$ and genotype $d$ ($c$) which does not have allele $\alpha$ at locus $i$ will recombine to a genotype with allele $\alpha$ at locus $i$, we have

$$\sum_{a=1}^{M}\delta(g_i^a,\alpha)\psi_a\left(\check{\boldsymbol{z}}(\tau)\right) = \sum_{c=1}^{M}\left(\sum_{d=1}^{M}\delta(g_i^c,\alpha)\delta(g_i^d,\alpha)\check{z}_c(\tau)\check{z}_d(\tau) + \frac{1}{2}\sum_{d=1}^{M}\delta(g_i^c,\alpha)(1-\delta(g_i^d,\alpha))\check{z}_c(\tau)\check{z}_d(\tau)\right)$$

$$+ \frac{1}{2}\sum_{c=1}^{M}\sum_{d=1}^{M}(1-\delta(g_i^c,\alpha))\delta(g_i^d,\alpha)\check{z}_c(\tau)\check{z}_d(\tau)$$

$$= \check{x}_{i,\alpha}^2(\tau) + \frac{1}{2}\check{x}_{i,\alpha}(\tau)(1-\check{x}_{i,\alpha}(\tau)) + \frac{1}{2}\check{x}_{i,\alpha}(\tau)(1-\check{x}_{i,\alpha}(\tau))$$

$$= \check{x}_{i,\alpha}(\tau)$$

thus implying that $\sum_{a=1}^{M}\delta(g_i^a,\alpha)\left(\check{z}_a(\tau) - \psi_a\left(\check{\boldsymbol{z}}(\tau)\right)\right) = 0$ and hence $I_3 = 0$.

The final expression for the drift vector is thus

$$\underline{d}_{i,\alpha}(\check{\boldsymbol{x}}(\tau),\Delta t) = \sum_{j=1}^{L}\sum_{\beta=1}^{\ell-1}\left(\check{x}_{ij,\alpha\beta}(\tau) - \check{x}_{i,\alpha}(\tau)\check{x}_{j,\beta}(\tau)\right)\bar{s}_{j,\beta} + \bar{\mu}_{\ell\alpha} + \sum_{\beta=1}^{\ell-1}\left(\bar{\mu}_{\beta\alpha} - \bar{\mu}_{\ell\alpha}\right)\check{x}_{i,\beta}(\tau) - \check{x}_{i,\alpha}(\tau)\sum_{\beta=1}^{\ell}\bar{\mu}_{\alpha\beta}. \quad \text{(S43)}$$

As expected, we observe that when there are two alleles, the mutation probability is symmetric and the reference allele is the WT, the expression in (S43) reduces to the expression in (S24).

18

We now move on to the diffusion matrix. This can be partitioned into $L^2$ blocks each of size $(\ell-1)\times(\ell-1)$, where the $(i,j)$th block has $(\alpha,\beta)$th entry

$$
\begin{aligned}
\underline{C}_{ij,\alpha\beta}(\check{\boldsymbol{x}}(\tau)) &= \sum_{a=1}^{M}\sum_{b=1}^{M}\delta(g_i^a,\alpha)\delta(g_j^b,\beta)\bar{C}_{ab}(\check{\boldsymbol{z}}(\tau)) \\
&= \frac{1}{2}\sum_{a=1}^{M}\delta(g_i^a,\alpha)\delta(g_j^a,\beta)\check{z}_a(\tau)(1-\check{z}_a(\tau)) - \frac{1}{2}\sum_{a=1}^{M}\sum_{b=1,b\neq a}^{M}\delta(g_i^a,\alpha)\delta(g_i^b,\beta)\check{z}_a(\tau)\check{z}_b(\tau) \\
&= \frac{1}{2}\sum_{a=1}^{M}\delta(g_i^a,\alpha)\delta(g_j^a,\beta)\check{z}_a(\tau) - \frac{1}{2}\left(\sum_{a=1}^{M}\delta(g_i^a,\alpha)\check{z}_a(\tau)\right)\left(\sum_{b=1}^{M}\delta(g_i^b,\beta)\check{z}_b(\tau)\right) \\
&= \frac{1}{2}\left(\check{x}_{ij,\alpha\beta}(\tau) - \check{x}_{i,\alpha}(\tau)\check{x}_{j,\beta}(\tau)\right).
\end{aligned}
$$

Following along the lines of the proof in Section 1.4, the probability of observing mutant allele frequencies $(\boldsymbol{x}(t_1), \boldsymbol{x}(t_2), \ldots, \boldsymbol{x}(t_K))$, and hence the likelihood function (S35), is approximated by the path integral

$$
\mathfrak{L}\left(\boldsymbol{s}|\boldsymbol{\mu}, N, (\boldsymbol{x}(t_k))_{k=0}^{K}\right) \approx \left[\prod_{k=0}^{K-1}\frac{1}{\sqrt{\det C(\boldsymbol{x}(t_k))}}\left(\frac{N}{2\pi\Delta t_k}\right)^{L(\ell-1)/2}\mathrm{d}\boldsymbol{x}(t_{k+1})\right]\exp\left(-\frac{N}{2}S\left((\boldsymbol{x}(t_k))_{k=0}^{K}\right)\right), \quad \text{(S44)}
$$

where $\mathrm{d}\boldsymbol{x}(t_{k+1}) = \prod_{i=1}^{L}\prod_{\alpha=1}^{\ell-1}\mathrm{d}x_{i,\alpha}(t_{k+1})$ and

$$
\begin{aligned}
&S\left((\boldsymbol{x}(t_k))_{k=0}^{K}\right) \\
&= \sum_{1}\frac{[x_{i,\alpha}(t_{k+1}) - x_{i,\alpha}(t_k) - \Delta t_k d_{i,\alpha}(\boldsymbol{x}(t_k))]\left(C^{-1}(\boldsymbol{x}(t_k))\right)_{ij,\alpha\beta}[x_{j,\beta}(t_{k+1}) - x_{j,\beta}(t_k) - \Delta t_k d_{j,\beta}(\boldsymbol{x}(t_k))]}{\Delta t_k}.
\end{aligned}
$$

Here we have adopted the shorthand notation $\sum_{1} = \sum_{k=0}^{K-1}\sum_{i=1}^{L}\sum_{\alpha=1}^{\ell-1}\sum_{j=1}^{L}\sum_{\beta=1}^{\ell-1}$, while $(\boldsymbol{X})_{ij,\alpha\beta}$ indicates the $(\alpha,\beta)$th entry of the $(i,j)$th block, with each block having dimension $(\ell-1)\times(\ell-1)$. Moreover, we have

$$
d_{i,\alpha}(\boldsymbol{x}(t_k)) = \sum_{j=1}^{L}\sum_{\beta=1}^{\ell-1}(x_{ij,\alpha\beta}(t_k) - x_{i,\alpha}(t_k)x_{j,\beta}(t_k))s_{j,\beta} + \mu_{\ell\alpha} + \sum_{\beta=1}^{\ell-1}(\mu_{\beta\alpha} - \mu_{\ell\alpha})x_{i,\beta}(t_k) - x_{i,\alpha}(t_k)\sum_{\beta=1}^{\ell}\mu_{\alpha\beta}
$$

and

$$
C_{ij,\alpha\beta}(\boldsymbol{x}(t_k)) = x_{ij,\alpha\beta}(t_k) - x_{i,\alpha}(t_k)x_{j,\beta}(t_k).
$$

### 1.7.3 The MPL estimator solution

Substituting the likelihood approximation (S44) and the prior

$$
P_{\mathrm{prior}}(\boldsymbol{s}) = \frac{1}{(2\pi\sigma^2)^{L(\ell-1)/2}}\exp\left(-\frac{1}{2\sigma^2}\boldsymbol{s}^{\mathrm{T}}\boldsymbol{s}\right)
$$

into the equivalent MAP problem

$$
\hat{\boldsymbol{s}} = \arg\max_{\boldsymbol{s}}\left(\log\mathfrak{L}\left(\boldsymbol{s}|\boldsymbol{\mu}, N, (\boldsymbol{x}(t_k))_{k=0}^{K}\right) + \log P_{\mathrm{prior}}(\boldsymbol{s})\right),
$$

we take the vector derivative with respect to $\boldsymbol{s}$ and equate to zero. This yields the MPL estimate

$$
\begin{aligned}
\hat{s}_{i,\alpha} = \sum_{j=1}^{L}\sum_{\beta=1}^{\ell-1}&\left[\sum_{k=0}^{K-1}\Delta t_k C(\boldsymbol{x}(t_k)) + \gamma I\right]_{ij,\alpha\beta}^{-1} \\
&\times\left[x_{j,\beta}(t_K) - x_{j,\beta}(t_0) - \sum_{k=1}^{K}\Delta t_k\left(\mu_{\ell\alpha} + \sum_{l=1}^{\ell-1}\left(\mu_{l\alpha} - \mu_{\ell\alpha}\right)x_{j,l}(t_k) - x_{j,\alpha}(t_k)\sum_{l=1}^{\ell}\mu_{\alpha l}\right)\right] \quad \text{(S45)}
\end{aligned}
$$

for $i = 1, \ldots, L$ and $\alpha = 1, \ldots, \ell-1$, which is the result quoted in Eq. (12) in Methods.

# 2 Data analysis

This section provides details about the synthetic data and its analysis used to obtain results presented in the paper. Section 2.1 describes the data and analysis in Figs. 1 and 2 and Extended Data Figs. 1-4. This data in these figures was designed to demonstrate the performance of MPL for a dynamic but generic system where linkage is pervasive. Section 2.2 describes the data pre-processing performed to implement different methods in literature. Section 2.3 describes the data and analysis of a system designed to mimic the HIV-1 data set described in Supplementary Table 1.

## 2.1 Simulated data generation and performance analysis

The results in Figs. 1 and 2 and Extended Data Figs. 1-4 were obtained for a 50-locus biallelic system. We generated a population of $N = 1000$ sequences, each of length $L = 50$, and evolved the population according to the WF model with selection and mutation. We assumed loci to be biallelic with 0 representing the WT and 1 representing the mutant allele. Parameters for the simulations are included in the figure captions. Scripts for generating and analyzing these data are located in the GitHub repository.

We compared MPL with seven methods from the literature: WFABC [12], FIT [13], ApproxWF [14], LLS [15], CLEAR [16], EandR [17], and IM [18]. We compared the accuracy of these methods by measuring the AUROC for classification of beneficial and deleterious selection coefficients, as well as the normalized root-mean-squared error (NRMSE) of the inferred selection coefficients,

$$\text{NRMSE} := \sqrt{\frac{\sum_{i=1}^{L} \left(\hat{s}_i - s_i\right)^2}{\sum_{i=1}^{L} s_i^2}} . \tag{S46}$$

We also recorded the run time of all algorithms in the same computational environment. All values were averaged over samples from 100 WF simulations with the same underlying parameters.

## 2.2 Implementation and data preprocessing specifications

We wrote custom codes for the FIT and IM methods, while for the remaining methods we used software implementations provided by the respective authors. Scripts used to generate the comparison results in the paper are available at the GitHub repository.

**Pre-processing**: Some methods required the input sequence or allele frequency data to be modified in order to produce reasonable results. This *pre-processing* was performed according to one or two rules, as indicated below. In describing these, we define a *trajectory* as a set of mutant allele frequencies at consecutively sampled time points. The start of a trajectory is marked by the first observation of a polymorphism, while the end is marked by either the fixation or loss of the mutant allele, or by the last sampled time point. Thus, over an entire observation period, it is possible to observe multiple "trajectories" at the same locus.

1. *Pre-filtering:* Trajectories were filtered such that (a) the minimum length of a trajectory was two time points and (b) the maximum (minimum) frequency was at least 0.05 (0.95) for a trajectory that started from a frequency of zero (one).

2. *Bit-flipping:* The definitions of WT and mutant were reversed at loci for which trajectories had initial frequency greater than 0.95.

**MPL and SL.** MPL and its single locus variant, the independent model described in the main text, were implemented in C++. No preprocessing of the sequence data was required for MPL and SL.
**WFABC.** Trajectories were converted to WFABC's format using custom Matlab scripts and then analyzed using the code provided at http://jjensenlab.org/software (accessed on September 16, 2017). Initial tests revealed that we had to pre-process the raw sampled trajectories before applying the WFABC method in order to get meaningful results. All trajectories were pre-filtered except the ones that were rendered monomorphic due to pre-filtering. Bit-flipping was also applied. These processing steps removed spurious noisy blips from the trajectories and redefined WT/mutant, resulting in improved performance of the WFABC method.
**ApproxWF.** Trajectories were converted to an input format compatible with ApproxWF using custom Matlab scripts and then analyzed using the code provided at https://bitbucket.org/wegmannlab/approxwf/wiki/Home (commit fcc7964 dated: 2016-10-04 accessed on September 26, 2017). All trajectories were pre-filtered as outlined above, except the ones that were rendered monomorphic due to filtering. Bit-flipping was not used.
**IM.** We developed a custom MATLAB script to implement the IM method [18]. Trajectories were pre-filtered, but not flipped. We used the filtering threshold of 0.1 (0.9) as specified by the authors of this method. The method

was implemented according to the description given in ref. [18], which applies a simulated annealing algorithm to estimate the selection coefficients. We used the simulated annealing algorithm implementation provided in the Global Optimization Toolbox of MATLAB 2017a. We let the algorithm run for a maximum of $100,000$ iterations and stopped the algorithm when the average change in value of the objective function in the previous $10,000$ iterations was less than a tolerance value (default value of $10^{-8}$).

In all simulation runs, the optimization stopped in less than $100,000$ iterations, indicating that the algorithm had converged. To further validate our in-house implementation of IM, we applied it to data sets generated using the system parameters as in ref.[18] and under the same simulation setup; e.g., continuous long runs of several hundred thousand generations. The results were qualitatively similar to those reported in ref. [18]. We found that IM's procedure of calling trajectories (described in ref. [18]), was sensitive to the evolutionary parameters used, e.g., population size, mutation probability ($N = 10,000$ and $\mu = 5 \times 10^{-7}$ were used in ref. [18]), as well as to the simulation condition of not allowing a locus to mutate after it has reached fixation (in the simulation setup of IM, a locus that reaches fixation remains at fixation for the next 3200 generations). For the simulation results presented in this paper, where $N = 1000$, $\mu = 10^{-4}$ and there are no restrictions on mutations at a locus reaching fixation, we had to change the parameters of the procedure for calling trajectories so that the algorithm returned meaningful results. Using the default parameters of ref. [18] for calling trajectories resulted in the IM method missing the start/end of trajectories. The modified parameters are specified in Matlab scripts in the GitHub repository.

**LLS.** Trajectories were fed directly into the LLS code provided by the authors at https://github.com/ThomasTaus/poolSeq. A small number of loci resulted in a "N/A" selection coefficient estimate, which occurred when the loci had a frequency of zero or one for the vast majority of time points, and having a frequency close to zero or one for the other (small number of) time points. These loci were excluded from the analysis.

**FIT.** No additional preprocessing was required. The trajectories were analyzed using custom Matlab scripts.

**CLEAR.** Simulation data was converted into a format readable by CLEAR using custom Python scripts. No preprocessing of the data was necessary.

**EandR.** Simulation data was converted into a format readable by EandR using custom Python scripts. No preprocessing of the data was required.

## 2.3 Simulated data generation and performance analysis with sampling conditions similar to HIV-1 dataset

To test robustness to different sampling parameters, we analyzed the performance of MPL under sampling conditions similar to the HIV-1 dataset (Supplementary Table 1), where the number of sequences $n_s$ and the time between samples $\Delta t$ varies at each time point, and where the trajectory length $T$ can vary across simulations. Specifically, the number of sequences $n_s$ at each time point was drawn from a binomial distribution with $n = 1000$ and $p = 0.0139$ (Extended Data Fig. 4a). This resulted in $n_s$ values that varied between 5 and 26, with the mean equal to the average number of sequences per time point in our data. The time between samples, $\Delta t$, was drawn from a gamma mixture distribution selected to match with the distribution of $\delta t$ values in the data, resulting in $\Delta t$ values that varied between 2 and 131 (Extended Data Fig. 4b). Time gaps for each trajectory were rearranged so that shorter gaps appeared at the beginning of the trajectories and larger time gaps at the end (Extended Data Fig. 4d), similar to what was observed in the HIV-1 data (Supplementary Table 1). The length of the trajectories used for inference was also drawn from a gamma mixture distribution resulting in trajectory lengths between 17 and 330 generations (Extended Data Fig. 4c). The simulation parameters consisted of population size $N = 1000$, $L = 50$ loci with two alleles at each locus (mutant and WT), 10 beneficial mutants with selection coefficients $s$ uniformly distributed over the range [0.075, 0.125], 30 neutral mutants with $s = 0$, and 10 deleterious mutants with $s$ uniformly distributed over the range [-0.125, -0.075], mutation probability per site per generation $\mu = 10^{-4}$, and recombination probability per site per generation $r = 10^{-4}$. Results for the perfect sampling case are evaluated for those sites that are polymorphic in the heterogeneous sampling case.

These simulations show that MPL readily identifies loci with beneficial selection coefficients (Extended Data Fig. 4e). Our ability to detect beneficial selection with heterogeneous sampling (i.e., imperfect sampling with parameters chosen according to the distributions in Extended Data Fig. 4a-c) is only slightly worse than in the case with perfect sampling, where the complete genotype distribution in the population at every generation is known. Under these sampling conditions, deleterious selection coefficients are more difficult to identify than beneficial ones.

# References

[1] Doria-Rose, N. A. *et al.* Developmental pathway for potent V1V2-directed HIV-neutralizing antibodies. *Nature* **509**, 55 (2014).

[2] Kimura, M. Diffusion models in population genetics. *Journal of Applied Probability* **1**, 177–232 (1964).

[3] Ewens, W. J. *Mathematical Population Genetics 1: Theoretical Introduction* (Springer Science & Business Media, 2012).

[4] Tataru, P., Bataillon, T. & Hobolth, A. Inference under a Wright-Fisher model using an accurate beta approximation. *Genetics* **201**, 1133–1141 (2015).

[5] He, Z., Beaumont, M. & Yu, F. Effects of the ordering of natural selection and population regulation mechanisms on Wright-Fisher models. *G3: Genes, Genomes, Genetics* **7**, 2095–2106 (2017).

[6] Tataru, P., Simonsen, M., Bataillon, T. & Hobolth, A. Statistical inference in the Wright-Fisher model using allele frequency data. *Systematic Biology* **66**, e30–e46 (2017).

[7] Risken, H. *The Fokker-Planck Equation: Methods of Solution and Applications* (Springer-Verlag, 1989), 2nd edn.

[8] Mustonen, V. & Lässig, M. Fitness flux and ubiquity of adaptive evolution. *Proceedings of the National Academy of Sciences* **107**, 4248–4253 (2010).

[9] Schraiber, J. G. A path integral formulation of the Wright–Fisher process with genic selection. *Theoretical Population Biology* **92**, 30–35 (2014).

[10] Illingworth, C. J., Parts, L., Schiffels, S., Liti, G. & Mustonen, V. Quantifying selection acting on a complex trait using allele frequency time series data. *Molecular Biology and Evolution* **29**, 1187–1197 (2011).

[11] Lütkepohl, H. *Handbook of Matrices*, vol. 1 (Wiley Chichester, 1996).

[12] Foll, M., Shim, H. & Jensen, J. D. WFABC: a Wright–Fisher ABC–based approach for inferring effective population sizes and selection coefficients from time-sampled data. *Molecular Ecology Resources* **15**, 87–98 (2015).

[13] Feder, A. F., Kryazhimskiy, S. & Plotkin, J. B. Identifying signatures of selection in genetic time series. *Genetics* **196**, 509–522 (2014).

[14] Ferrer-Admetlla, A., Leuenberger, C., Jensen, J. D. & Wegmann, D. An approximate Markov model for the Wright–Fisher diffusion and its application to time series data. *Genetics* **203**, 831–846 (2016).

[15] Taus, T., Futschik, A. & Schlötterer, C. Quantifying selection with pool-seq time series data. *Molecular Biology and Evolution* **34**, 3023–3034 (2017).

[16] Iranmehr, A., Akbari, A., Schlötterer, C. & Bafna, V. Clear: Composition of likelihoods for evolve and resequence experiments. *Genetics* **206**, 1011–1023 (2017).

[17] Terhorst, J., Schlötterer, C. & Song, Y. S. Multi-locus analysis of genomic time series data from experimental evolution. *PLoS Genetics* **11**, e1005069 (2015).

[18] Illingworth, C. J. & Mustonen, V. Distinguishing driver and passenger mutations in an evolutionary history categorized by interference. *Genetics* **189**, 989–1000 (2011).