

Evolution of evolvability in rapidly adapting populations

In the format provided by the
authors and unedited

Supplementary Information

Contents

	Page
1 Comparison to deterministic modifier theory	3
2 Evolvability modifiers in the successive mutations regime	4
3 Evolvability modifiers in the clonal interference regime	7
3.1 Macroscopic epistasis expansion	7
3.2 Fitness wave formalism	11
3.3 Asymptotic solution in the “sharp shoulder” regime	13
4 Solution for a simple model of the distribution of fitness effects	16
4.1 Multiple mutations regime	17
4.1.1 Location of the interference threshold	17
4.1.2 Extending the shoulder solution to lower fitness values	18
4.1.3 First-order mutations and the rate of adaptation	20
4.1.4 Modifiers without direct costs or benefits	22
4.1.5 Modifiers with direct costs or benefits	25
4.1.6 Fixation of a dead-end modifier	27
4.1.7 Relation between the fixation probability and the long-term rate of adaptation	30
4.2 Quasi-sweep regime	31
4.2.1 Location of the interference threshold	31
4.2.2 Extending the shoulder solution to lower fitness values	32
4.2.3 Fixation probabilities of modifiers	34
5 Extension to continuous distributions of fitness effects	36
5.1 Location of the interference threshold for the wildtype population	37
5.2 Perturbative regime	39
5.2.1 Location of the interference threshold	39
5.2.2 Fixation probabilities of modifiers	41
5.3 Modifier-dominated regime (multiple mutations)	41
5.3.1 Location of the interference threshold	42
5.3.2 Extending the shoulder solution to lower fitness values.	42
5.3.3 Fixation probabilities of modifiers	47
5.4 Modifier-dominated regime (quasi-sweeps)	48
5.4.1 Location of the interference threshold	48
5.4.2 Extending the shoulder solution to lower fitness values	48
6 Incorporating deleterious mutations	50

7	Extensions to more general fitness landscapes	52
7.1	Weak Macroscopic Epistasis	52
7.2	Transient differences in evolvability	52
7.3	Global diminishing returns epistasis	53
8	Numerical Methods	55
8.1	Theoretical predictions	55
8.2	Empirical example from Ref. (8)	57

1 Comparison to deterministic modifier theory

In this section, we rederive a key result from classical modifier theory known as the ‘‘mean fitness principle’’ (20). This theory predicts that in an infinitely large asexual population, natural selection will favor modifiers that increase the long-term mean fitness of the population. Our derivation closely follows the one given in Ref. (21), as well as related work on mutators (81–83). We reproduce it here for completeness, using the same notation employed in our more general analysis below.

To establish this result, we note that in the absence of genetic drift ($N = \infty$), the deterministic dynamics of a well-mixed asexual population can be written in the general form,

$$\frac{\partial f(\vec{g})}{\partial t} = [X(\vec{g}) - \bar{X}(t)] f(\vec{g}) + \sum_{\vec{g}'} M(\vec{g}' \rightarrow \vec{g}) f(\vec{g}') - \sum_{\vec{g}} M(\vec{g} \rightarrow \vec{g}') f(\vec{g}), \quad (\text{S1})$$

where $f(\vec{g}, t)$ is the frequency of genotype \vec{g} , $X(\vec{g})$ is the (log) fitness of genotype \vec{g} , $\bar{X}(t) = \sum_{\vec{g}} X(\vec{g}) f(\vec{g}, t)$ is the mean fitness of the population, and $M(\vec{g} \rightarrow \vec{g}')$ is the mutation rate from genotype \vec{g} to \vec{g}' . The fate of a general evolvability modifier allele can be analyzed by introducing an analogous set of equations,

$$\frac{\partial f_m(\vec{g})}{\partial t} = [X_m(\vec{g}) - \bar{X}(t)] f_m(\vec{g}) + \sum_{\vec{g}'} M_m(\vec{g}' \rightarrow \vec{g}) f_m(\vec{g}') - \sum_{\vec{g}} M_m(\vec{g} \rightarrow \vec{g}') f_m(\vec{g}), \quad (\text{S2})$$

where $f_m(\vec{g}, t)$ is the frequency of individuals with the mutant allele at the modifier locus and genotype \vec{g} elsewhere, $X_m(\vec{g})$ is genotype-to-fitness map for the modifier, $M_m(\vec{g} \rightarrow \vec{g}')$ is the corresponding set of mutation rates, and the mean fitness in Eqs. (S1) and (S2) now sums over both the mutant and wildtype lineages,

$$\bar{X}(t) = \sum_{\vec{g}} X(\vec{g}) f(\vec{g}, t) + \sum_{\vec{g}} X_m(\vec{g}) f_m(\vec{g}, t). \quad (\text{S3})$$

In principle, this model allows for arbitrary changes to arbitrary non-linear fitness landscapes, encapsulating all possible forms of epistasis (84).

The solutions to this coupled system of equations can be written in the general form

$$f_m(\vec{g}, t) = f_m(t) h_m(\vec{g}, t), \quad f(\vec{g}, t) = [1 - f_m(t)] h_w(\vec{g}, t), \quad (\text{S4a})$$

where $f_m(t)$, $h_m(t)$, and $h_w(t)$ satisfy the related set of equations,

$$\frac{\partial h_w(\vec{g})}{\partial t} = [X(\vec{g}) - \bar{X}_w(t)] h_w(\vec{g}) + \sum_{\vec{g}'} M(\vec{g}' \rightarrow \vec{g}) h_w(\vec{g}') - \sum_{\vec{g}} M(\vec{g} \rightarrow \vec{g}') h_w(\vec{g}), \quad (\text{S4b})$$

$$\frac{\partial h_m(\vec{g})}{\partial t} = [X_m(\vec{g}) - \bar{X}_m(t)] h_m(\vec{g}) + \sum_{\vec{g}'} M_m(\vec{g}' \rightarrow \vec{g}) h_m(\vec{g}') - \sum_{\vec{g}} M_m(\vec{g} \rightarrow \vec{g}') h_m(\vec{g}), \quad (\text{S4c})$$

and

$$\frac{\partial f_m}{\partial t} = [\bar{X}_m(t) - \bar{X}_w(t)] f_m(1 - f_m), \quad (\text{S4d})$$

with

$$\bar{X}_w(t) \equiv \sum_{\vec{g}} X(\vec{g}) h_w(\vec{g}, t), \quad \bar{X}_m(t) \equiv \sum_{\vec{g}} X_m(\vec{g}) h_m(\vec{g}, t). \quad (\text{S4e})$$

In this notation, $h_m(\vec{g}, t)$ represents the re-normalized genotype distribution within the modifier lineage, $h_w(\vec{g}, t)$ denotes the corresponding distribution within the wildtype, $\bar{X}_m(t)$ and $\bar{X}_w(t)$ represent the mean fitnesses of each lineage, and $f_m(t)$ denotes the total frequency of the modifier.

This change-of-variables shows that the total frequency of the modifier lineage only depends on the relative values of the mean fitnesses, $\bar{X}_m(t)$ and $\bar{X}_w(t)$, yielding the time-dependent solution,

$$\log \left[\frac{f_m(t)}{1 - f_m(t)} \right] = \log \left[\frac{f_m(0)}{1 - f_m(0)} \right] + \int_0^t [\bar{X}_m(t') - \bar{X}_w(t')] dt'. \quad (\text{S5})$$

This shows that at long times ($t \rightarrow \infty$) the lineage that will dominate the population is the one with the higher value of $\int_0^\infty \bar{X}_i(t') dt'$. These mean fitnesses can be predicted from the dynamics of the intra-lineage frequencies, $h_m(\vec{g}, t)$ and $h_w(\vec{g}, t)$, which decouple from each other and from the total size of $f_m(t)$. In most cases of interest, the mean fitness of each lineage will approach an equilibrium value, $\bar{X}_i(t) \rightarrow \bar{X}_i$, which may be a complicated function of $X_i(\vec{g})$ and $M_i(\vec{g} \rightarrow \vec{g}')$, but is otherwise independent of the surrounding population. Equation (S5) then shows that the modifier will take over the population if and only if it increases the equilibrium mean fitness. This generalizes the ‘‘mean fitness principle’’ derived in previous work (20).

However, this deterministic calculation neglects two key factors that are relevant for any large but finite population. First, the mean-field dynamics in Eqs. (S1) and (S2) neglect the random occurrence of new mutations and the stochastic fluctuations they experience while rare. These fluctuations can dramatically influence the dynamics of the mean fitness – particularly when multiple beneficial mutations are available (51). In addition, the deterministic calculation neglects the possibility that the mutant or wildtype lineage may fix before their long-term benefits in Eq. (S5) are fully realized. As we will see below, both of these effects will become extremely important in the adapting populations that we analyze in this work. Interestingly, we will see that in some cases, it will be possible to account for these effects in an approximate manner by inserting an upper limit on the integral in Eq. (S5) (SI Section 4.1.7), providing a conceptual link between classical modifier theory and the more complex scenarios studied in this work. Identifying such cases and their appropriate time horizons is the goal of the next several sections.

2 Evolvability modifiers in the successive mutations regime

Another useful limit occurs in small populations, where the production of beneficial mutations is sufficiently rare that adaptation proceeds via a sequence of discrete selective sweeps. In this regime, the fate of a given modifier will strongly depend on its ability to generate the next beneficial mutation. We can formalize this idea by defining the local distribution of fitness effects (DFE),

$$\mu(s|\vec{g}) = \sum_{\vec{g}'} M(\vec{g} \rightarrow \vec{g}') \cdot \delta(s - X(\vec{g}') + X(\vec{g})), \quad (\text{S6})$$

which tabulates the fitness effects of all the mutations that can be accessed from a given genotype \vec{g} . Modifier individuals will have their own corresponding set of DFEs,

$$\mu_m(s|\vec{g}) = \sum_{\vec{g}'} M_m(\vec{g} \rightarrow \vec{g}') \cdot \delta(s - X_m(\vec{g}') + X_m(\vec{g})), \quad (\text{S7})$$

along with a direct cost or benefit

$$s_m(\vec{g}) = X_m(\vec{g}) - X(\vec{g}). \quad (\text{S8})$$

The distributions in Eqs. (S6) and (S7) can be computed for any epistatic fitness landscape, and will constitute the key input parameters in our analysis below.

In the absence of any modifier mutations, the evolutionary dynamics of the successive sweeps regime are well described by previous work (51). The wildtype population produces new beneficial mutations at a total rate $NU_b = N \cdot \int_0^\infty \mu(s|\vec{g})ds$ per generation. Most of these mutations will drift to extinction before they reach appreciable frequencies. However, with probability $p_{\text{est}} \approx 2s$, a lucky mutant will fluctuate to a sufficiently large frequency ($f \approx 1/2Ns$) that it begins to grow deterministically [$\partial_t f \approx sf(1-f)$], and will sweep through the population on a timescale of order $T_{\text{fix}} = \frac{2}{s} \log(Ns)$. The population will produce these successful mutations at a total rate $\lambda = \int_0^\infty N\mu(s|\vec{g}) \cdot 2s = 2NU_b\bar{s}$ per generation, where $\bar{s} \equiv \int s\mu(s|\vec{g}) ds / \int \mu(s|\vec{g}) ds$ is the average beneficial fitness effect that is accessible to the current genotype. This implies that the typical waiting time for the next sweep event is $T_c \approx 1/2NU_b\bar{s}$ generations. The condition that successive sweeps will not interfere with each requires that $T_{\text{fix}} \ll T_c$, which requires that $NU_b \log[N\bar{s}] \ll 1$ (51). This limit is also known as the ‘‘strong selection, weak mutation regime’’ (SSWM).

A general evolvability modifier ($\mu(s) \rightarrow \mu_m(s)$) will produce a change in the beneficial mutation rate ($U_b \rightarrow U'_b$) as well as the typical fitness benefit ($\bar{s} \rightarrow \bar{s}'$). (It may also change the spectrum of deleterious mutations, which we neglect for the time being; see SI Section 6). The mutator version of this model was previously analyzed in Ref. (26), as well as a number of related studies (29, 82, 85–94). We reproduce these calculations below, while generalizing them to allow for changes in the average fitness benefits of mutations (similar to Ref. 29) as well as a broader range of direct costs and benefits. These classical results will provide a useful baseline for understanding the impact of larger populations, which we analyze in SI Section 3 below.

In the successive mutations regime, a newly arising modifier lineage will initially compete with the wildtype population according to the single-locus dynamics

$$\frac{\partial f_m}{\partial t} = s_m f_m (1 - f_m) + \sqrt{\frac{f_m(1 - f_m)}{N}} \eta(t), \quad (\text{S9})$$

where $\eta(t)$ is a Brownian noise term (95) and s_m is the direct cost or benefit of the modifier. The next selective sweep will now be generated by a pair of competing Poisson processes with rates

$$\lambda_m(t) = 2NU'_b\bar{s}' \cdot f_m(t), \quad (\text{S10a})$$

$$\lambda_w(t) = 2NU_b\bar{s} \cdot [1 - f_m(t)], \quad (\text{S10b})$$

which correspond to the mutant and wildtype lineages, respectively. Since the mutation rate and fitness benefit both contribute linearly to $\lambda_i(t)$, this process is formally equivalent to the mutator scenario analyzed in (26), with $r = U'_b\bar{s}'/U_b\bar{s}$ replacing U'_b/U_b . This allows us to conclude that in absence of a direct cost ($s_m = 0$), the fixation probability of the modifier scales as

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s)) \approx \begin{cases} U'_b\bar{s}'/U_b\bar{s} & \text{if } \frac{U'_b\bar{s}'}{U_b\bar{s}} \ll NU_b \cdot N\bar{s}, \\ \sqrt{2 \left(\frac{U'_b\bar{s}'}{U_b\bar{s}} \right) NU_b \cdot N\bar{s}} & \text{if } \frac{U'_b\bar{s}'}{U_b\bar{s}} \gg NU_b \cdot N\bar{s}, \end{cases} \quad (\text{S11})$$

provided that $T_c \ll N$ (or $NU_b \cdot N\bar{s} \gg 1$). When $NU_b \cdot N\bar{s} \lesssim 1$, the modifier will either fix or go extinct neutrally ($p_{\text{fix}} \approx 1/N$) before the next sweep occurs. This shows that second-order selection is more efficient in larger populations, which is reminiscent of our results in Fig. 2. In this case, however, the benefits of second-order selection are capped by the ratio $U'_b\bar{s}'/U_b\bar{s}$, which implies that very large changes in U'_b or \bar{s}' are required to produce an appreciable fixation probability.

A similar application of Ref. (26) shows that the fixation probability of a modifier with a direct cost ($s_m < 0$) scales as

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{\int_{|s_m|}^{\infty} (s-s_m)\mu_m(s)ds}{U_b\bar{s}} \cdot \frac{1}{1 + \frac{|s_m|}{2NU_b\bar{s}}}, \quad (\text{S12})$$

provided that $U'_b\bar{s}'/U_b\bar{s} \ll NU_b \cdot N\bar{s}$. The integral in the numerator implies that modifiers with direct costs larger than $\sim\bar{s}'$ will have dramatically reduced fixation probabilities, since future beneficial mutations in these backgrounds will still be less fit than the current wildtype population. However, since $NU_b \ll 1$, the fixation probability of the modifier will be significantly reduced by the denominator term well before these “wasted opportunities” start to become relevant. This illustrates how even small direct costs – much smaller than a single driver mutation – can overwhelm the evolvability benefits of mutations in the successive sweeps regime.

The fixation probability of a modifier with a direct benefit ($s_m > 0$) can be computed using a similar procedure. When $s_m \ll 2NU_b\bar{s} \ll \bar{s}$, the next sweep will typically occur before the modifier lineage establishes, so an analogous version of Eq. (S12) still holds:

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{\frac{U'_b\bar{s}'}{U_b\bar{s}}}{1 - \frac{s_m}{2NU_b\bar{s}}}, \quad (s_m \ll 2NU_b\bar{s} \ll \bar{s}) \quad (\text{S13})$$

For stronger fitness benefits ($s_m \gg 2NU_b\bar{s} \gg 1/N$), the modifier will have the opportunity establish and grow deterministically before the next sweep occurs:

$$f_m(t) \approx \begin{cases} \frac{\frac{1}{2Ns_m}e^{s_mt}}{1 + \frac{1}{2Ns_m}(e^{s_mt}-1)} & \text{w/ prob } 2s_b, \\ 0 & \text{else.} \end{cases} \quad (\text{S14})$$

In this regime, the fixation process is similar to the “first-step” clonal interference analysis in Ref. (96). Provided that $s_m \ll \bar{s}$, the additional fitness benefit of the modifier will have a negligible impact the establishment probability of the next sweep in either genetic background, so that

$$\lambda_m(t) = 2NU'_b(\bar{s}' + s_m [1 - f_m(t)] \cdot f_m(t) \approx 2NU'_b\bar{s}' \cdot f_m(t), \quad (\text{S15a})$$

$$\lambda_w(t) = 2N \int_{s_m}^{\infty} [s - s_m f_m(t)] \mu(s) ds \cdot [1 - f_m(t)] \approx 2NU_b\bar{s} [1 - f_m(t)]. \quad (\text{S15b})$$

The fixation probability is determined by these competing Poisson processes, so that

$$\begin{aligned} \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &= N \left\langle \int_0^{\infty} dt \lambda_m(t) e^{-\int_0^{\infty} dt' [\lambda_m(t') + \lambda_w(t')]} \right\rangle, \\ &\approx 2NU'_b\bar{s}' \int_0^{\infty} dt e^{s_mt - 2NU_b\bar{s}t} \left[2N \left(\frac{U'_b\bar{s}' - U_b\bar{s}}{s_m} \right) + 1 \right] \cdot \log \left[1 + \frac{1}{2Ns_m} (e^{s_mt} - 1) \right], \\ &\approx 2Ns_m e^{-\frac{2NU_b\bar{s}}{s_m} \log[2Ns_m]}, \end{aligned} \quad (\text{S16})$$

which will be valid in the limit that $2NU_b\bar{s} \ll s_m \ll \bar{s}$. This result has a simple heuristic interpretation as the establishment probability of the modifier multiplied by the probability the wildtype population does not generate a sweep before the modifier fixes on its own (96). Since this fixation probability is independent of $\mu_m(s)$, it implies that even small direct benefits — much smaller than the size of a typical driver mutation

— can override the effects of second order selection, even if they reduce the long-term rate of adaptation to zero. Together with the direct cost results above, these calculations emphasize that natural selection can be extremely sensitive to the short-term costs or benefits of a mutation in the successive sweeps regime, in contrast to what we observe in larger populations like Fig. 1A.

3 Evolvability modifiers in the clonal interference regime

In larger populations ($NU_b \gtrsim 1$), the assumption of discrete selective sweeps will start to break down. Multiple beneficial lineages will segregate in the population at the same time, and will interfere with each other as they compete for dominance in the population (51). In this *clonal interference regime*, the fate of a given mutation will sensitively depend on the genetic background that it arises on, and the future mutations that its descendants produce before they fix or are driven to extinction. This requires a stochastic generalization of the full multi-locus dynamics in Eqs. (S1) and (S2),

$$\begin{aligned} \frac{\partial f(\vec{g})}{\partial t} = & \underbrace{[X(\vec{g}) - \bar{X}(t)] f(\vec{g})}_{\text{selection}} + \underbrace{\sum_{\vec{g}'} M(\vec{g}' \rightarrow \vec{g}) f(\vec{g}') - \sum_{\vec{g}} M(\vec{g} \rightarrow \vec{g}') f(\vec{g})}_{\text{mutation}} \\ & + \underbrace{\sqrt{\frac{f(\vec{g})}{N}} \eta(\vec{g}, t) - f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')}{N}} \eta(\vec{g}', t) - \sum_{\vec{g}'} \sqrt{\frac{f_m(\vec{g}')}{N}} \eta_m(\vec{g}', t)}_{\text{genetic drift}}, \end{aligned} \quad (\text{S17a})$$

and

$$\begin{aligned} \frac{\partial f_m(\vec{g})}{\partial t} = & \underbrace{[X_m(\vec{g}) - \bar{X}(t)] f_m(\vec{g})}_{\text{selection}} + \underbrace{\sum_{\vec{g}'} M_m(\vec{g}' \rightarrow \vec{g}) f_m(\vec{g}') - \sum_{\vec{g}} M_m(\vec{g} \rightarrow \vec{g}') f_m(\vec{g})}_{\text{mutation}} \\ & + \underbrace{\sqrt{\frac{f(\vec{g})}{N}} \eta(\vec{g}, t) - f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')}{N}} \eta(\vec{g}', t) - \sum_{\vec{g}'} \sqrt{\frac{f_m(\vec{g}')}{N}} \eta_m(\vec{g}', t)}_{\text{genetic drift}}, \end{aligned} \quad (\text{S17b})$$

where $\eta(\vec{g}, t)$ and $\eta_m(\vec{g}, t)$ are uncorrelated Brownian noise terms (97). While this model is straightforward to write down, there are no known solutions for arbitrary choices of $X(\vec{g})$ and $X_m(\vec{g})$. Further progress requires us to make specific assumptions about the shapes of these fitness landscapes, as well as the corresponding mutation kernels $M(\vec{g} \rightarrow \vec{g}')$ and $M_m(\vec{g} \rightarrow \vec{g}')$. We introduce one particularly convenient parameterization in the next section, which is motivated by the concept of *macroscopic epistasis* explored in previous work (69, 98).

3.1 Macroscopic epistasis expansion

There are many possible fitness landscapes one could consider. These are often parameterized as a power series involving different combinations of loci,

$$X(\vec{g}) = \sum_{\ell=1}^L s_\ell g_\ell + \sum_{\ell < \ell'} \epsilon_{\ell, \ell'} g_\ell g_{\ell'} + \sum_{\ell < \ell' < \ell''} \epsilon_{\ell, \ell', \ell''} g_\ell g_{\ell'} g_{\ell''} + \dots \quad (\text{S18})$$

where $g_\ell = 1$ if there is a mutation at site ℓ and 0 otherwise (84). This representation can be viewed as a Taylor expansion around a perfectly smooth fitness landscape, where all the $\{\epsilon_{\ell, \dots, \ell'}\}$ coefficients vanish. Nonzero values of the $\epsilon_{\ell, \dots, \ell'}$ coefficients correspond to epistatic interactions between loci; following previous work (69), we will refer to these interactions as *microscopic epistasis*, since they can in principle vary across all possible combinations of sites.

Modifier alleles can be expressed in this framework by designating an arbitrary site as the modifier locus (e.g. $\ell = m$), and recalculating the landscape for $g_m = 1$ (modifiers) and $g_m = 0$ (wildtype) separately. This notation makes it clear that a modifier that changes the fitness effects of other mutations must necessarily involve some microscopic epistasis, corresponding to terms like $\epsilon_{m, \ell}$, $\epsilon_{m, \ell, \ell'}$, and so on.

The space of epistatic fitness landscapes is enormous, and their impact on the evolutionary dynamics of large populations is not well understood (70). Some studies have attempted to navigate this complexity by truncating Eq. (S18) after the pairwise terms (99); others have focused on smaller landscapes containing just a handful of interacting loci (100, 101). In this work, we show that it will be useful to consider an alternative limit of Eq. (S18), in which the local DFEs defined in Eqs. (S6) and (S7) are approximately constant for different genotypes:

$$\mu(s|\vec{g}) \approx \mu(s), \quad \mu_m(s|\vec{g}) \approx \mu_m(s). \quad (\text{S19})$$

We can view this approximation as the lowest order term in an alternative expansion of the fitness landscape,

$$\mu(s|\vec{g}) \approx \mu(s) + \sum_{\ell} \delta\mu_{\ell}(s) \cdot g_{\ell} + \sum_{\ell < \ell'} \delta\mu_{\ell, \ell'}(s) \cdot g_{\ell} g_{\ell'} + \dots, \quad (\text{S20})$$

which works directly in the DFE basis. This genotype dependence of the DFE is sometimes known as *macroscopic epistasis* (69), since it aggregates over a large number of microscopic interactions in Eq. (S18). The approximation in Eq. (S19) can therefore be viewed as the simplest possible model that incorporates some amount of macroscopic epistasis, with a non-zero $\delta\mu_{\ell}(s)$ term at the modifier locus ($\ell = m$) and all other $\delta\mu_{\ell, \dots, \ell'}$ terms vanishing.

We note that such macroscopic epistasis can arise even in the absence of microscopic epistasis ($\epsilon_{\ell, \dots, \ell'} = 0$) if the modifier alters the mutation rates at other loci. For example, mutator alleles are often modeled as a simple change of scale,

$$M_m(\vec{g} \rightarrow \vec{g}') = r \cdot M_m(\vec{g} \rightarrow \vec{g}'), \quad (\text{S21})$$

leading to a proportional change in the DFE,

$$\mu_m(s|\vec{g}) = r \cdot \mu(s|\vec{g}), \quad (\text{S22})$$

which has been the starting point for many previous studies (26, 28, 94). In practice, mutator strains are typically biased toward specific types of mutations, so that the simple proportional model in Eq. (S21) does not necessarily apply (4, 6, 29). Other mutagenesis mechanisms that target specific genomic regions (102, 103) lead to conceptually similar complications. Both require us to consider changes in the shape of the DFE in addition to its overall scale (4, 6, 29).

In addition to mutator alleles, Eq. (S19) can also emerge from epistatic interactions between loci. The simplest example is a pairwise epistasis model, with nonzero values for the coefficients involving $\ell = m$, and zeros everywhere else. This is sufficient to recover the limit in Eq. (S19), but it is not the only possibility: any other landscape that gives rise to the same overall distributions in Eq. (S19) will exhibit similar dynamics, even if the individual fitness effects [$s_{\ell}(\vec{g}) \equiv X(\vec{g} + \ell) - \bar{X}(\vec{g})$] are undergoing more complicated rearrangements under the hood.

Moreover, our results will not require Eq. (S19) to hold across the entire fitness landscape, but only within a smaller region that is explored before the modifier either fixes or goes extinct. We determine the size of this local neighborhood in SI Section 7 and find that it is often modest, corresponding to just a handful of mutational steps for many empirically relevant parameter values (Fig. 4). We will also show that the assumption in Eq. (S19) is most sensitive to a narrow range of beneficial fitness effects, so that substantial deviations in other parts of the DFE can still have a negligible impact on the modifier lineage (SI Section 7). At present, it is difficult to enumerate all of the microscopic landscapes that are consistent with a given DFE function, $\mu(s|\vec{g})$. However, we can still identify several examples that satisfy Eq. (S19) – at least in the approximate sense required – that go beyond the mutator allele and pairwise epistasis examples above:

Branching epistatic landscapes. At a formal level, we can consider a “maximally epistatic” landscape of branching uphill paths of length K , where each step k of a given path can access $M \ll L$ other beneficial mutations (Extended Data Fig. 2A). If the uphill paths do not contain any loops, the fitness function can be expressed in the form

$$X(\vec{g}) = \sum_{k=1}^K \sum_{(\ell_1, \dots, \ell_k) \in \mathcal{L}_k} \left(\sum_{j=1}^k s_{\ell_j} \right) \left(\prod_{j=1}^k g_{\ell_j} \right) \prod_{\ell \notin \{\ell_j\}_{j=1}^k} (1 - g_{\ell}), \quad (\text{S23})$$

where \mathcal{L}_k denotes the set of all uphill sub-paths of length k , and the s_{ℓ} coefficients are independently drawn from the target DFE $\mu(s)$. This construction ensures that the DFEs calculated from Eq. (S23) will satisfy $\mu(s|\vec{g}) \approx \mu(s)$ when M is large, even though the landscape contains large amounts of microscopic epistasis. We can extend this construction to include a modifier allele by writing

$$X(\vec{g}) = \sum_{k=1}^K \sum_{(\ell_1, \dots, \ell_k) \in \mathcal{L}_k} \left(\sum_{j=1}^k s_{\ell_j} (1 - g_m) + s'_{\ell_j} g_m \right) \left(\prod_{j=1}^k g_{\ell_j} \right) \prod_{\ell \notin \{\ell_j\}_{j=1}^k} (1 - g_{\ell}), \quad (\text{S24})$$

where the s'_{ℓ} coefficients are independently drawn from the modifier DFE $\mu_m(s)$. This microscopic landscape satisfies the macroscopic epistasis approximation in Eq. (S19), but includes many non-pairwise interactions by construction.

Non-linear global phenotypes (e.g. stability-activity tradeoffs). The approximation in Eq. (S19) can also emerge in more concrete settings, when the fitness of the organism is a non-linear function of multiple global phenotypes. A prototypical example is the stability-activity tradeoff that is often observed in viral evolution (104, 105) and other protein-binding problems (106), where stabilizing mutations can potentiate the fitness benefits of mutations that would destabilize the protein on their own. The essential ingredients of this behavior can be captured in a simple model containing two global phenotypes: (i) activity, denoted by $\Psi(\vec{g})$, which contributes additively to the total fitness, and (ii) stability, denoted by $\Phi(\vec{g})$, which has a Boltzmann-like contribution,

$$X(\vec{g}) = -\log \left(1 + e^{-\Phi(\vec{g})} \right) + \Psi(\vec{g}) + \text{const}, \quad (\text{S25a})$$

We will assume that both phenotypes can be expressed as additive functions of the genotype,

$$\Phi(\vec{g}) = \phi_0 + \sum_{\ell}^L \phi_{\ell} \cdot g_{\ell}, \quad \Psi(\vec{g}) = \sum_{\ell}^L \psi_{\ell} \cdot g_{\ell}, \quad (\text{S25b})$$

where ϕ_0 denotes the stability of the reference strain. In this model, mutations that only increase the activity of the protein are always beneficial ($s_\ell \approx \psi_\ell$). However, mutations that increase activity while also decreasing stability can be either costly or beneficial depending on the stability of the background that they occur on.

Suppose that there are a large number (M) of such mutations with $|\phi_\ell| - \phi_0 \gg 1$, which implies that they will be deleterious on the wildtype background. In this scenario, any mutation that strongly increases stability will function like an evolvability modifier, by allowing these previously deleterious mutations to occur. In particular, if the stability enhancement ϕ_m is much larger than $|\phi_\ell| - \phi_0$, then the fitness benefits of the unleashed mutations will be approximately constant ($s_\ell \approx \psi_\ell$) until roughly $K \approx \phi_m / (|\phi_\ell| - \phi_0)$ such variants have accumulated. If $M \gg K \gg 1$, this example will satisfy the approximation in Eq. (S19), while relying on the higher-order epistatic interactions in Eq. (S25).

While this particular example was motivated by protein stability, similar behavior can occur for other combinations of phenotypes, as long as they include the appropriate non-linearities. For example, a simple model of stabilizing selection involving a one nearly optimized phenotype and another non-optimized one can be expressed as

$$X(\vec{g}) = -|\Phi(\vec{g})| + \Psi(\vec{g}), \quad (\text{S26a})$$

where

$$\Phi(\vec{g}) = \sum_{\ell}^L \phi_{\ell} \cdot g_{\ell}, \quad \Psi(\vec{g}) = \sum_{\ell}^L \psi_{\ell} \cdot g_{\ell}. \quad (\text{S26b})$$

A mutation that increases $\Psi(\vec{g})$ while displacing $\Phi(\vec{g})$ from its optimal value would enable $K \approx \phi_m / \phi_\ell$ previously deleterious mutations to accumulate — each providing a fitness benefit $s_\ell \approx \phi_\ell$ — before the phenotypic optimum is reattained. If K is sufficiently large, this evolvability modifier would also satisfy the approximation in Eq. (S19), while involving a distinct form of non-pairwise epistasis.

Chromosomal duplications / aneuploidy. The approximation in Eq. (S19) can also apply to scenarios that are difficult to capture with a traditional fitness landscape, because they involve changes to the structure of the genome itself. Classical examples include chromosomal duplications and other copy-number changes, which are frequently observed in cancer evolution (14, 107) and laboratory evolution experiments in eukaryotes (108). These copy number variants involve changes in ploidy — in addition to changes in target size — so that dominance effects start to become important. For example, a whole-genome duplication of a haploid genome would lead to a modified DFE of the form

$$\mu_m(s) = \sum_{\ell=1}^L 2\mu_{\ell} \cdot \delta(s - s_{\ell} h_{\ell}), \quad (\text{S27})$$

where h_{ℓ} denotes the dominance coefficient of the mutation at site ℓ . The successive sweeps picture in SI Section 2 predicts that the fixation probability of this copy-number variant is $\propto 2\overline{sh}/s$, which is neutral in purely semi-dominant case ($\overline{sh} = s/2 \implies \tilde{p}_{\text{fix}} = 1$) and only moderately beneficial for complete dominance ($\overline{sh} = s/2 \implies \tilde{p}_{\text{fix}} \approx 2$). In contrast, our results below suggest that in larger populations, these variants can be both strongly favored or disfavored by second-order selection depending on the joint distribution of (s_{ℓ}, h_{ℓ}) .

Together, these examples illustrate that the limiting behavior in Eq. (S19) — while still capturing just a subset of all possible fitness landscapes — can nevertheless approximate many biologically relevant scenarios where

second-order selection is thought to play an important role. We will therefore use this model as a starting point for all of our mathematical derivations below. We will revisit this assumption in SI Section 7, where we discuss extensions to other possible fitness landscapes.

3.2 Fitness wave formalism

An advantage of the model in Eq. (S19) is that it allows us to exploit existing “fitness wave” methods for modeling clonal interference (26, 42, 43, 50, 51, 75, 109–112). This previous literature has shown that the multi-locus dynamics in Eq. (S17) can often be simplified by considering a coarse-grained picture, which groups together individuals with the same overall fitness:

$$f(X, t) = \sum_{\vec{g}} f(\vec{g}, t) \cdot \delta(X - X(\vec{g})). \quad (\text{S28})$$

We can extend this idea to our present context, defining a corresponding fitness distribution for the modifier lineage as well:

$$f_m(X, t) = \sum_{\vec{g}} f_m(\vec{g}, t) \delta(X - X_m(\vec{g})). \quad (\text{S29})$$

In the special case that $\mu(s|\vec{g}) = \mu(s)$ and $\mu_m(s|\vec{g}) = \mu_m(s)$ for all genotypes \vec{g} , Eq. (S17) can be written in the coarse-grained form,

$$\begin{aligned} \frac{\partial f(X)}{\partial t} = & \underbrace{[X - \bar{X}(t)] f(X)}_{\text{selection}} + \underbrace{\int ds \mu(s) [f(X-s) - f(X)]}_{\text{mutation}} \\ & + \underbrace{\int dX' [\delta(X - X') - f(X)] \sqrt{\frac{f(X')}{N}} \eta(X') - f(X) \int dX' \sqrt{\frac{f_m(X')}{N}} \eta_m(X')}_{\text{genetic drift}}, \end{aligned} \quad (\text{S30a})$$

$$\begin{aligned} \frac{\partial f_m(X)}{\partial t} = & \underbrace{[X - \bar{X}(t)] f_m(X)}_{\text{selection}} + \underbrace{\int ds \mu_m(s) [f_m(X-s) - f_m(X)]}_{\text{mutation}} \\ & + \underbrace{\int dX' [\delta(X - X') - f_m(X)] \sqrt{\frac{f_m(X')}{N}} \eta_m(X') - f_m(X) \int dX' \sqrt{\frac{f(X')}{N}} \eta(X')}_{\text{genetic drift}}, \end{aligned} \quad (\text{S30b})$$

where $\bar{X}(t) = \int X f(X, t) dX + \int X f_m(X, t) dX$ is the mean fitness of the population, and the $\eta_i(X)$ are uncorrelated Brownian noise terms (97). This allows us to generalize the notion of a mutator allele to handle more general differences in evolvability, while still bypassing the enormous complexity of the underlying fitness landscape.

In the absence of the modifier ($f_m = 0$), the wildtype distribution $f(X, t)$ approaches a traveling wave form that increases in fitness at an average rate $\langle \partial_t \bar{X}(t) \rangle \equiv v(\mu(s), N)$ (42, 51). Previous work has shown that the typical profile $f(x)$ is well-approximated by the deterministic equation,

$$-v \partial_x f(x) = x \cdot f(x) + \int \mu(s) [f(x-s) - f(x)] ds, \quad (\text{S31})$$

which is the expected value of Eq. (S30a) when $\bar{X}(t) \approx vt$ (42, 43).

We assume that the modifier mutation will arise in this steady-state population, on a genetic background drawn from $f(x)$. Its descendants will found a second fitness wave, $f_m(X, t)$, which competes with the wildtype population $f(X, t)$ as they continue to acquire additional mutations according to the dynamics in Eq. (S17) (Fig. 1).

Branching process approximation. The key approximation we will make in this work is that the fate of the modifier will often be determined while it is still at a low frequency in the population. Previous studies have shown that this is a good approximation for first-order mutations – including both neutral and deleterious mutations – in the clonal interference regime (42, 43, 47, 50, 112, 113). In this work, we use a combination of self-consistency arguments and comparisons with simulations to show that this same approximation holds for a broad range of modifier alleles as well. The main exceptions will occur for modifiers that dramatically reduce evolvability (e.g. the dead-end modifiers in Fig. 3D); we treat this case separately in SI Section 4.1.6.

When the frequency of the modifier is small, the mean fitness of the population can be approximated by the wildtype value $\bar{X}(t) \approx vt$, and the higher-order contributions in the drift term in Eq. (S30b) can be neglected. The dynamics of the modifier then reduce to the multitype linear branching process,

$$\frac{f_m(X)}{\partial t} = [X - vt] f_m(X) + \int ds \mu_m(s) [f_m(X - s) - f_m(X)] + \sqrt{\frac{f_m(X)}{N}} \eta_m(X), \quad (\text{S32})$$

with the initial condition $f_m(X, 0) = \frac{1}{N} \delta(X - x - s_m)$. The long-term non-extinction probability of this process, defined by

$$w_m(x) \equiv w(x | \mu(s) \rightarrow \mu_m(s)) \equiv \lim_{t \rightarrow \infty} \left(1 - e^{-\int f_m(X, t) dx} \right), \quad (\text{S33})$$

can be calculated using the formal procedures described in Ref. (50). This yields the standard branching process recursion listed in Eq. (1) in the main text:

$$0 = \underbrace{x \cdot w(x | \mu(s) \rightarrow \mu_m(s))}_{\text{selection}} + \underbrace{\int \mu_m(s) [w(x + s | \mu(s) \rightarrow \mu_m(s)) - w(x | \mu(s) \rightarrow \mu_m(s))] ds}_{\text{mutation}} - \underbrace{v(\mu(s), N) \cdot \partial_x w(x | \mu(s) \rightarrow \mu_m(s))}_{\text{competition with wildtype}} - \underbrace{\frac{1}{2} \cdot w(x | \mu(s) \rightarrow \mu_m(s))^2}_{\text{genetic drift while rare}}, \quad (\text{S34})$$

which is a straightforward generalization of the mutator version in Ref. (26). If the fate of the modifier is determined while it is rare, then this non-extinction probability must also coincide with the long-term probability of fixation. This implies that the overall fixation probability of the modifier can be obtained from Eq. (S34) by averaging over the random genetic background,

$$p_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) = \int f(x) \cdot w(x + s_m | \mu(s) \rightarrow \mu_m(s)) dx, \quad (\text{S35})$$

which is equivalent to Eq. (2) in the main text. In the limit that $\mu_m(s) \rightarrow \mu(s)$, Eq. (S35) reduces to the fixation probability of a first-order mutation which has been studied in previous work (42, 43, 50, 112). In particular, this previous work has shown that a self-consistency condition can be obtained from the fact that the overall fixation probability of a neutral mutation must always be equal to $1/N$ (this follows directly

from the normalizability of the full dynamics in Eq. S30). Combining this fact with Eq. (S35) yields the self-consistency condition,

$$p_{\text{fix}}(\mu(s) \rightarrow \mu(s), 0) = \int f(x) \cdot w(x|\mu(s) \rightarrow \mu(s)) dx = \frac{1}{N}. \quad (\text{S36})$$

Together with Eqs. (S31) and (S34), this completely determines the wildtype rate of adaptation $v(\mu(s), N)$ as a function of $\mu(s)$ and N (42, 43, 50, 112).

When $\mu_m(s) \neq \mu(s)$, the conditional fixation probability in Eq. (S34) will differ from that of a first-order mutation because there are different mutation spectra in the mean fitness and mutation terms. This is because the modifier lineage primarily competes against the wildtype population [whose mean fitness is controlled by $\mu(s)$], but acquires additional mutations from its own distribution, $\mu_m(s)$. We note, however, that since the competition with the wildtype population is completely mediated by $v(\mu(s), N)$ (a monotonic function of N), the conditional fixation probability of the modifier lineage can always be mapped to the conditional fixation probability of a first-order mutation in a population that is fixed for the modifier allele, but with a different population size N^* . In other words,

$$w(x|\mu(s) \rightarrow \mu_m(s), N) = w(x|\mu_m(s) \rightarrow \mu_m(s), N^*), \quad (\text{S37})$$

where N^* is defined by

$$v(\mu_m(s), N^*) = v(\mu(s), N). \quad (\text{S38})$$

This implies that the space of solutions for $w(x|\mu(s) \rightarrow \mu_m(s))$ will have the same general form as the “first-order” $w(x)$ function that has been studied in previous work (42, 43, 50, 111, 112). It also implies that the dynamics of non-extinction in Eq. (S32), which were previously described in Ref. (43), will be qualitatively similar as well. However, we will see below that actually using this result for evolvability modifiers will often require solutions to Eq. (S34) that go beyond the parameter regimes that have been examined in these earlier studies. We have therefore developed a new approach for deriving approximate analytical solutions for $w(x|\mu(s) \rightarrow \mu_m(s))$ that will apply across this broader range of parameters. We outline the general approach in SI Section 3.3 below, and apply it to different classes of distributions of fitness effects in SI Sections 4 and 5.

For notational convenience, we will suppress the explicit dependence on $\mu(s)$ and $\mu_m(s)$ in the following sections, writing $w(x) \equiv w(x|\mu(s) \rightarrow \mu(s))$ for the conditional fixation probability of a first-order mutation and $w_m \equiv w(x|\mu(s) \rightarrow \mu_m(s))$ for the conditional fixation probability of the modifier lineage. We will also let $v \equiv v(\mu(s), N)$ denote the rate of adaptation in the wildtype population and $v_m \equiv v(\mu_m(s), N)$ denote the long-term rate of adaptation that is achieved if the modifier takes over.

3.3 Asymptotic solution in the “sharp shoulder” regime

There are many parameter regimes of clonal interference that one can consider (43, 112). Here, we will primarily focus on a regime that is relevant for a broad range of naturally and experimentally evolving populations, which roughly corresponds to the case where the fitness benefits of a typical “driver mutation” (s_b) are much larger than the total rate at which they occur (U_b) (43, 50). We will define these conditions more precisely below.

Relative fitness distribution. Previous work has shown that in our regime of interest, the solution to Eq. (S31) can be approximated by a truncated Gaussian distribution,

$$f(x) \approx \begin{cases} \frac{1}{\sqrt{2\pi v}} e^{-\frac{x^2}{2v}} & x \leq x_c, \\ 0 & x > x_c, \end{cases} \quad (\text{S39})$$

where x_c denotes the location of the fittest individuals that are likely to exist within the population (42, 43, 50). This solution is valid when $x_c - s_b \gg \sqrt{v}$ and $s_b \gg \sqrt{v}$, which constitutes the formal definition of the regime that we consider in this work. The first condition ($x_c - s_b \gg \sqrt{v}$) implies that the parents of the most fit individuals are substantially more fit than the majority of the population (i.e. clonal interference is common). The second condition ($s_b \gg \sqrt{v}$) implies that the majority of the population is concentrated near the mean fitness. This is also known as the “moderate speeds” regime (43), since the second condition can be rewritten as $v \ll s_b^2$. We will assume that these conditions hold for the wildtype population in which the modifier arises.

Lineage fixation probability. Previous work (42, 50) has shown that in our parameter regime of interest, the solutions to Eq. (S34) can be decomposed into a high-fitness region, where the mutation term is subdominant:

$$0 \approx \underbrace{x \cdot w_m(x)}_{\text{selection}} - \underbrace{v \cdot \partial_x w_m(x)}_{\text{competition w/ wt}} - \underbrace{\frac{1}{2} \cdot w_m(x)^2}_{\text{genetic drift}}, \quad (\text{S40})$$

and a linearized region, where the mutation term is important but the drift term can be neglected:

$$0 \approx \underbrace{x \cdot w_m(x)}_{\text{selection}} + \int ds \underbrace{\mu_m(s) [w_m(x+s) - w_m(x)]}_{\text{mutation}} - \underbrace{v \cdot \partial_x w_m(x)}_{\text{competition w/ wt}}. \quad (\text{S41})$$

There is also a narrow region in the middle where both approximations are valid, so that asymptotic matching can be used to obtain a full global solution. This asymptotic decomposition will continue to be valid for our modifier case as well.

In the high fitness region, the solution to Eq. (S40) is given by the *shoulder solution*,

$$w_m(x) \approx \frac{2x_{cm} e^{(x^2 - x_{cm}^2)/2v}}{1 + (x_{cm}/x) e^{(x^2 - x_{cm}^2)/2v}}, \quad (\text{S42})$$

where x_{cm} is a constant of integration that will be determined self-consistently below. When $x_{cm} \gg \sqrt{v}$, this shoulder solution develops a narrow boundary layer around $x = x_{cm} \pm O(v/x_{cm})$, such that the fixation probability can be well-approximated by the piecewise form,

$$w_m(x) \approx \begin{cases} 2x & \text{if } x - x_{cm} \gtrsim O(v/x_{cm}), \\ 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} & \text{if } x_{cm} - x \lesssim O(v/x_{cm}). \end{cases} \quad (\text{S43})$$

This piecewise function has a simple interpretation in terms of the dominant balances in Eq. (S40) (42). The linear scaling above x_{cm} emerges from a balance between the selection and genetic drift terms; this implies that lineages with $x \gtrsim x_{cm}$ will fix provided that they survive genetic drift. Conversely, the exponential scaling at lower x emerges from a balance between the selection and mean fitness terms; this implies that

lineages with $x \lesssim x_{cm}$ are strongly influenced by competition with the surrounding population (i.e. clonal interference). In the special case where $\mu_m(s) = \mu(s)$, previous work has shown that x_{cm} approximately coincides with the nose of the fitness distribution in Eq. (S39) (42, 43, 50). This makes some intuitive sense: lineages with relative fitness $\gtrsim x_c$ do not experience clonal interference, and will fix if they survive genetic drift; this suggests that the steady-state fitness distribution should not contain such individuals, since they should already have fixed. We will therefore also refer to x_c as the *interference threshold* for the wildtype population. We will therefore also refer to x_c as the *interference threshold* for the wildtype population. When $\mu_m(s) \neq \mu(s)$ the interference threshold of the modifier (x_{cm}) will generally differ from the location of the nose (x_c) — this will be critically important for our analysis in SI Sections 4 and 5 below.

The mutation terms must eventually become important for smaller values of x , since the shoulder solution in Eq. (S43) starts to increase for $x < 0$. To avoid this unphysical behavior, and ensure $w_m(x)$ decreases as $x \rightarrow \infty$, the shoulder solution must eventually map on to the correct branch of the solution to Eq. (S34). Previous work has used saddle point methods to identify the relevant solutions to Eq. (S41) (43, 112), while other studies have utilized thresholding approximations that set $w_m(x) \approx 0$ below a critical fitness value (42, 50). Here we take a slightly different approach, by recasting Eq. (S34) as an integral equation.

Multiplying both sides of Eq. (S34) by $e^{-\frac{(x-U_0)^2}{2v}}$ and integrating from $-\infty$ to x , we can rewrite Eq. (S34) in the recursive form,

$$w_m(x) = e^{\frac{(x-U_0)^2}{2v}} \int_0^\infty ds \int_{-\infty}^{x+s} \frac{\mu_m(s)}{v} \cdot e^{-\frac{(y-s-U_0)^2}{2v}} w_m(y) dy - e^{\frac{(x-U_0)^2}{2v}} \int_{-\infty}^x e^{-\frac{(y-U_0)^2}{2v}} \cdot \frac{w_m(y)^2}{2v} dy, \quad (\text{S44})$$

where $U_0 \equiv \int \mu_m(s) ds$ is the total mutation rate. Since $w_m(x)$ rapidly declines below the interference threshold, the contributions from the second term will become negligible when $x \lesssim x_{cm} - \mathcal{O}(v/x_{cm})$, and Eq. (S44) will reduce to the simpler form

$$w_m(x) = e^{\frac{(x-U_0)^2}{2v}} \int_0^\infty ds \int_{-\infty}^{x+s} \frac{\mu_m(s)}{v} \cdot e^{-\frac{(y-s-U_0)^2}{2v}} \cdot w_m(y) dy. \quad (\text{S45})$$

This recursive formula has a simple intuitive interpretation illustrated by the schematic in Fig. 1. A modifier lineage founded at a relative fitness $x < x_{cm}$ will start as a single clone, whose average size will evolve as

$$\langle n(t) \rangle = \exp \left[xt - \frac{vt^2}{2} - U_0 t \right]. \quad (\text{S46})$$

These growth dynamics account for the steady increase in the wildtype mean fitness, as well as the outflow of individuals that acquire further mutations as the clone is growing. The total production rate of new mutations is $\langle n(t) \rangle \cdot \mu(s) ds$, and each of these events will found a new modifier lineage with relative fitness $x + s - vt$; the original lineage will fix if one of these descendant lineages is ultimately successful, yielding the recursive formula,

$$w_m(x) = \int_0^\infty dt \int_0^\infty ds \mu_m(s) \cdot e^{xt - \frac{vt^2}{2} - U_0 t} \cdot w_m(x + s - vt). \quad (\text{S47})$$

The linearity of this expression implies that successful clones are highly unlikely to give rise to multiple successful mutants, which is consistent with the assumption that clonal interference is very strong ($w_m(x) \ll 2x$) when $x < x_{cm}$. Equation (S45) can be recovered from Eq. (S47) by changing variables from t to $y \equiv x - v \cdot t + s$, which represents the time-dependent landing fitness of the clone's mutant offspring.

Note that while the recursion in Eq. (S45) is also a solution to the linearized version of Eq. (S34), our present derivation shows that the interpretation is slightly different here. In particular, since the upper limit of the y -integral in Eq. (S45) is larger than x , the mutation term can in principle depend on the behavior of $w_m(x)$ for $x > x_{cm}$, where the effects of the nonlinear $w_m(x)^2$ term start to become important. We can account for this non-locality using the shoulder solution in Eq. (S43), by rewriting Eq. (S45) as a piecewise integral,

$$w_m(x) = e^{\frac{(x-U_0)^2}{2v}} \int_0^\infty ds \int_{-\infty}^{x+s} \frac{\mu_m(s)}{v} \cdot e^{-\frac{(y-s-U_0)^2}{2v}} \cdot \left[\frac{2x_{cm} e^{\frac{y^2-x_{cm}^2}{2v}} \theta(y-x_{\min})}{1 + \frac{x_{cm}}{y} e^{\frac{y^2-x_{cm}^2}{2v}}} + w_m(y) \theta(x_{\min} - y) \right] dy, \quad (\text{S48})$$

where $\theta(\cdot)$ is the unit step function, and x_{\min} represents the point at which the shoulder solution starts to break down. When $x_{\min} \lesssim x_{cm} - \mathcal{O}(v/x_{cm})$, we can use this expression to determine x_{cm} , by noting that Eq. (S48) must also match the shoulder solution in Eq. (S43) in the overlap region where both approximations are valid [$x_{\min} \lesssim x \lesssim x_c - \mathcal{O}(x_c)$].

While this integral formulation is similar to the transform methods (43, 112) that have previously been used to analyze $w_m(x)$, it offers practical advantages that will become important for our analysis of modifier mutations below. In particular, we will see that in our parameter regime of interest, we can use Eq. (S48) to analytically extend the shoulder solution to progressively lower fitness values. These analytical approximations will be critically important for treating modifier mutations with large direct costs (Fig. 3).

In the following sections, we use this framework to derive explicit solutions for $w_m(x)$ for different choices of $\mu_m(s)$. We begin by considering a simple model, where the driver mutations all share the same characteristic fitness benefit (SI Section 4). This will allow us to verify the self-consistency conditions assumed above, and to obtain explicit predictions for the overall fixation probabilities of different modifier mutations. We then extend these calculations to continuous distributions of fitness effects in SI Section 5, and show that more general distributions can often be understood using the simple model in SI Section 4.

4 Solution for a simple model of the distribution of fitness effects

To make analytical progress, we first consider a simple scenario in which the wildtype population produces mutations with a single fitness benefit s_b at a total rate U_b . The distribution of fitness effects can then be written as a point mass,

$$\mu(s) = U_b \cdot \delta(s - s_b), \quad (\text{S49})$$

where $\delta(z)$ is the Dirac delta function. Previous work has shown that in our parameter regime of interest, the rate of adaptation (v) and interference threshold (x_c) follow the approximate scaling,

$$v \approx \frac{2s_b^2 \log(Ns_b)}{\log^2(s_b/U_b)}, \quad x_c \approx \frac{2s_b \log(Ns_b)}{\log(s_b/U_b)}, \quad (\text{S50})$$

which are valid in the limit that $x_c \gg s_b$ (42, 43, 51). In terms of these parameters, the conditions we assumed for our solution in SI Section 3.3 ($x_c \gg \sqrt{v}$ and $s_b \gg \sqrt{v}$) become

$$1 \ll \log(Ns_b) \ll \log^2(s_b/U_b), \quad (\text{S51})$$

which require that $Ns_b \gg 1$ and $s_b \gg U_b$.

For distributions of fitness effects in the same class as Eq. (S49), the most general evolvability modifier is one that produces mutations with a different benefit s'_b and different rate U'_b , so that

$$\mu_m(s) = U'_b \cdot \delta(s - s'_b). \quad (\text{S52})$$

In the analysis below, we will assume that U'_b and s'_b are chosen such that $U'_b \ll s'_b$, $x_{cm} \gg \sqrt{v}$, and $s'_b \gg \sqrt{v}$. In this regime, the outflow of mutations can be neglected in Eq. (S48). We will also assume that the fold change in the mutation rate may be large [$\log(U'_b/U_b) \gg 1$], but the fold change in the selection strength will always be comparatively modest [$\log(s'_b/s_b) \lesssim \log(10)$]. This will be sufficient to derive all of the results in the main text.

Substituting Eq. (S52) into Eq. (S48) yields an integral relation for $w_m(x)$,

$$w_m(x) = 2x_{cm} e^{\frac{x^2}{2v} - \frac{x_{cm}^2}{2v}} \times \frac{U'_b}{v} \int_{-\infty}^{x+s'_b} \left(\frac{e^{\frac{ys'_b}{v} - \frac{s'^2_b}{2v}} \theta(y - x_{\min})}{1 + \frac{x_{cm}}{y} e^{\frac{y^2}{2v} - \frac{x_{cm}^2}{2v}}} + \frac{w_m(y) \theta(x_{\min} - y)}{2x_{cm} e^{\frac{(y-s'_b)^2}{2v} - \frac{x_{cm}^2}{2v}}} \right) dy, \quad (\text{S53})$$

where we have neglected the outflow due to new mutations since $U'_b \ll s'_b \lesssim x$. The interference threshold is determined by matching Eq. (S53) to the shoulder solution when $x \lesssim x_{cm} - \mathcal{O}(v/x_{cm})$, which yields a related integral,

$$1 \approx \frac{U_b}{v} \int_{-\infty}^{x_{cm}+s'_b} \left(\frac{e^{\frac{ys'_b}{v} - \frac{s'^2_b}{2v}} \theta(y - x_{\min})}{1 + \frac{x_{cm}}{y} e^{\frac{y^2}{2v} - \frac{x_{cm}^2}{2v}}} + \frac{w_m(y) \theta(x_{\min} - y)}{2x_{cm} e^{\frac{(y-s'_b)^2}{2v} - \frac{x_{cm}^2}{2v}}} \right) dy. \quad (\text{S54})$$

The integrals in Eqs. (S53) and (S54) will sensitively depend on the relative values of x_{cm} and s'_b . When $x_{cm} - s'_b \gg \sqrt{v}$, the integrand in Eq. (S54) will be maximized for relative fitnesses near x_{cm} . We will refer to this limit as the *multiple mutations regime*, since it implies that successful lineages must always possess anomalously high relative fitnesses ($x_{cm} - s'_b \gg \sqrt{v}$). In the opposite extreme ($s'_b - x_{cm} \gg \sqrt{v}$), the integrand in Eq. (S54) will be maximized for relative fitnesses near s'_b , where the shoulder solution has already transitioned to the Haldane limit [$w(x) \approx 2x$]. We will refer to this limit as the *quasi-sweep regime*, since it implies that lineages will be guaranteed to fix if they produce a single mutation that survives genetic drift. Our assumptions imply that the wildtype population will always fall in the multiple mutations regime, but the modifier lineage may differ depending on the relative values of (U'_b, s'_b) and (U_b, s_b) . We consider each case separately below.

4.1 Multiple mutations regime

4.1.1 Location of the interference threshold

When $x_{cm} - s'_b \gg \sqrt{v}$, the integral in Eq. (S53) will be dominated by fitnesses close to x_{cm} . In particular, the contributions from $x < x_{cm}$ will be dominated by fitnesses within $\mathcal{O}(v/s'_b)$ of x_{cm} , while the contributions from $x > x_{cm}$ will be dominated by fitness within $\mathcal{O}(v/(x_{cm} - s'_b))$ of x_{cm} . The dominant contribution will therefore depend on the relative magnitudes of s'_b and $x_{cm} - s'_b$. When $x_{cm} - s'_b \gg s'_b \gg \sqrt{v}$, the dominant contribution will come from the exponential region of Eq. (S43), so that the auxiliary condition becomes

$$1 \approx \left(\frac{U'_b}{s'_b} \right) \exp \left[\frac{x_{cm} s'_b}{v} - \frac{s'^2_b}{v} \right]. \quad (\text{S55})$$

Note that this derivation implicitly assumes that the shoulder solution extends at least $\mathcal{O}(v/s'_b)$ below x_{cm} ; we will validate this assumption below. In the opposite case where $s'_b \gg x_{cm} - s'_b \gg \sqrt{v}$, the dominant contribution to the integral in Eq. (S54) will come from the linear region of Eq. (S43), so that the auxiliary condition becomes

$$1 \approx \left(\frac{U'_b}{x_{cm} - s'_b} \right) \exp \left[\frac{x_{cm}s'_b}{v} - \frac{s'^2_b}{2v} \right]. \quad (\text{S56})$$

For convenience, we will summarize these two equations with the common expression,

$$1 \approx \left(\frac{U'_b}{s'_b} \right) \exp \left[\frac{x_{cm}s'_b}{v} - \frac{s'^2_b}{2v} \right] \left(\frac{x_{cm}}{x_{cm} - s'_b} \right), \quad (\text{S57})$$

which reduces to the correct scaling in the corresponding limits. Since we have assumed that the wildtype population always lies in the multiple mutations regime ($x_c - s_b \gg \sqrt{v}$), the background interference threshold must satisfy an analogous condition,

$$1 \approx \left(\frac{U_b}{s_b} \right) \exp \left[\frac{x_c s_b}{v} - \frac{s_b^2}{2v} \right] \left(\frac{x_c}{x_c - s_b} \right), \quad (\text{S58})$$

which matches the expression derived in Ref. (50). These expressions allow us to solve for x_{cm} and x_c as a function of (U'_b, s'_b) , (U_b, s_b) , and v . In particular, by dividing Eq. (S57) by Eq. (S58), we can solve for x_{cm} as a function of the fold changes in U_b and s_b :

$$x_{cm} \approx \frac{s_b}{s'_b} \left[x_c + \frac{s_b}{2} \left(\frac{s'^2_b}{s_b^2} - 1 \right) - \frac{v}{s_b} \log \left(\frac{U'_b}{U_b} \right) \right], \quad (\text{S59})$$

where we have neglected logarithmic corrections in x_{cm}/x_c and s'_b/s_b . This shows that the $x_{cm} - s'_b \gg \sqrt{v}$ assumption will be valid provided that

$$\frac{s'_b}{s_b} \lesssim \sqrt{\frac{2x_c - s_b}{s_b} - \frac{2v}{s_b^2} \log \left(\frac{U'_b}{U_b} \right)}. \quad (\text{S60})$$

In particular, Eq. (S60) shows that the multiple mutations regime will apply for pure mutation rate modifiers as large as $\log(U'_b/U_b) \lesssim (x_c - s_b)s_b/v$, while pure selection strength modifiers require the more stringent condition, $s'_b/s_b \lesssim \sqrt{(2x_c - s_b)/s_b}$. Violations of these conditions are considered in SI Section 4.2 below.

4.1.2 Extending the shoulder solution to lower fitness values

With the location of x_{cm} fixed by Eq. (S57), we will now use the integral recursion in Eq. (S53) to extend the shoulder solution for $w_m(x)$ to progressively lower fitness values. When $x + s'_b \gtrsim x_{cm} + \mathcal{O}(v/(x_{cm} - s'_b))$, the dominant contribution to Eq. (S54) will also be contained within the region of integration in Eq. (S53). We can therefore substitute Eq. (S54) to obtain

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}}. \quad (\text{S61})$$

This derivation makes it clear that the shoulder solution will continue to be valid for fitnesses as low as

$$x_{\min}(\mu_m(s)) \equiv x_{cm} - s'_b, \quad (\text{S62})$$

which will serve as our definition of x_{\min} in Eqs. (S54) and (S53) above. This validates our assumption that the mutation term in Eq. (S34) is negligible for $x \gtrsim x_{\min} + O(v/(x_{cm} - s'_b))$. Moreover, since $x_{cm} - x_{\min} = s'_b \gg v/s'_b$, it also validates our assumption that the location of x_{cm} in Eq. (S57) is completely determined by regions where the shoulder solution is valid ($x > x_{\min}$). For fitnesses below x_{\min} , the finite upper limit in Eq. (S53) will start to become important. However, as long as $x + s'_b \gtrsim x_{\min} + O(v/s'_b)$, the integral will continue to be dominated by regions where the shoulder solution is valid. When $x \lesssim x_{\min} - O(v/x_{cm})$, the upper limit of integration will fall within the exponential region of the shoulder solution, yielding

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \cdot \frac{U'_b}{v} \int_{x_{\min}}^{x+s'_b} e^{\frac{ys'_b}{v} - \frac{s'^2_b}{2v}} dy \approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \cdot \frac{U'_b}{s'_b} e^{\frac{xs'_b}{v} + \frac{s'^2_b}{2v}}. \quad (\text{S63})$$

After dividing this expression by the auxiliary condition for x_{cm} in Eq. (S57), we obtain

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \cdot \left(\frac{x_{cm} - s'_b}{x_{cm}} \right) e^{-\frac{(x_{cm} - s'_b - x)s'_b}{v}}, \quad (\text{S64})$$

which will be valid for relative fitnesses in the range

$$x_{cm} - s'_b - O(v/x_{cm}) \gtrsim x \gtrsim x_{cm} - 2s'_b + O(v/s'_b). \quad (\text{S65})$$

By comparing Eq. (S61) and Eq. (S64), we can see that $w_m(x)$ declines by a factor of $\sim e^{s'^2_b/v}$ across this region. Since we have assumed that $s'_b \gg \sqrt{v}$, this provides a natural justification for the thresholding approximation employed in Ref. (50), which assumed that $w(x)$ was negligible for fitnesses below x_{\min} . Equation (S64) constitutes a more quantitative version of this approximation, which will be useful for the analysis below.

We can continue these calculations to recursively extend $w_m(x)$ to progressively lower fitness values. When $x \lesssim x_{cm} - 2s'_b + O(v/x_{cm})$, the shoulder solution will not contribute at all to Eq. (S53), but the region in Eq. (S64) will dominate, yielding

$$\begin{aligned} w_m(x) &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \left(\frac{x_{cm} - s'_b}{x_{cm}} \right) \frac{U'_b}{v} \int^{x+s'_b} e^{\frac{2ys'_b}{v} - \frac{s'^2_b}{2v} - \frac{(x_{cm} - s'_b)s'_b}{v}} dy, \\ &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \cdot \frac{1}{2} \left(\frac{x_{cm} - s'_b}{x_{cm}} \right)^2 e^{-\frac{2s'_b}{v} \left(x_{cm} - \frac{3s'_b}{2} - x \right)}, \end{aligned} \quad (\text{S66})$$

for $x \gtrsim x_{cm} - 3s'_b + O(v/s'_b)$. More generally, for relative fitnesses in the range

$$x_{cm} - ns'_b - O(v/s'_b) \gtrsim x \gtrsim x_{cm} - (n+1)s'_b + O(v/s'_b), \quad (\text{S67})$$

one can continue this argument to show that

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \cdot \frac{1}{n!} \left(\frac{x_{cm} - s'_b}{x_{cm}} \right)^n e^{-\frac{ns'_b(x_{cm} - ns'_b - x)}{v} - \frac{n(n-1)s'^2_b}{2v}}. \quad (\text{S68})$$

This shows that each additional ‘‘step’’ reduces the original shoulder solution by factor of $\sim e^{n(n-1)s'^2_b/2v} \gg 1$. For fitness values that are many multiples of s'_b below x_{cm} , we can substitute $n \approx (x_{cm} - x)/s'_b$ to show that the leading order contribution to $w_m(x)$ scales as

$$\log w_m(x) \sim \text{const} - \frac{(x_{cm} - x)(x_{cm} - s'_b/2)}{v}, \quad (\text{S69})$$

which obeys the required boundary condition that $w_m(x) \rightarrow 0$ as $x \rightarrow -\infty$. Thus, Eq. (S68) provides an asymptotic solution for $w_m(x)$ that is valid across the full range of relative fitnesses below x_{cm} .

4.1.3 First-order mutations and the rate of adaptation

When $\mu_m(s) \rightarrow \mu(s)$, our solution for $w_m(x)$ can be used to derive the fixation probabilities of first-order mutations, which have been studied in previous work (42, 43, 50). We reproduce these results here for completeness, using the new expressions for $w(x)$ that we have derived above. This will allow us to make comparisons to the second-order selection pressures analyzed in SI Section 4.1.4 below. (Readers who are familiar with this material may skip directly to SI Section 4.1.4.)

For a completely neutral mutation ($s_m = 0$), we can substitute our solution for $f(x)$ in Eq. (S39) into the self-consistency condition in Eq. (S36) to rewrite it in the convenient form:

$$\frac{1}{N} = \int f(x)w(x) dx = \frac{2x_c e^{-\frac{x_c^2}{2v}}}{\sqrt{2\pi v}} \int_{-\infty}^{x_c} \frac{w(x)}{2x_c e^{-\frac{x^2-x_c^2}{2v}}} dx, \quad (\text{S70})$$

which depends on the ratio between the true value of $w(x)$ and the exponential shoulder solution in Eq. (S43). Substituting our solution for $w(x)$ in Eq. (S68) then yields

$$\frac{1}{N} = \frac{2x_c e^{-\frac{x_c^2}{2v}}}{\sqrt{2\pi v}} \left(\int_{x_c-s_b}^{x_c} dx + \int_{x_c-2s_b}^{x_c-s_b} e^{-\frac{s_b(x_c-s_b-x)}{v}} dx + \dots \right), \quad (\text{S71})$$

where we have only included the $n = 0$ and $n = 1$ regions from Eq. (S68); all of the other terms are smaller by additional factors of $\exp(n(n-1)s_b^2/2v)$, and will therefore provide a negligible contribution when $s_b \gg \sqrt{v}$. Note that the first term in Eq. (S71) contributes equally over the whole range of $x \in (x_c - s_b, s_b)$, while the second term is dominated by fitnesses within $O(v/s_b)$ of $x_c - s_b$. Since $s_b \gg v/s_b$, the first term dominates, yielding

$$1 \approx \frac{2Nx_c s_b}{\sqrt{2\pi v}} e^{-\frac{x_c^2}{2v}}, \quad (\text{S72})$$

which matches the condition previously derived in Ref. (50). When combined with the auxiliary condition for x_c in Eq. (S58), this allows us to solve for v and x_c as a function of the underlying parameters N , s_b , and U_b . An iterative solution of Eqs. (S72) and (S58) assuming that $x_c \gg s_b \gg \sqrt{v}$ yields the asymptotic expressions for v and x_c in Eq. (S50). The fixation probabilities of non-neutral mutations can be calculated using a similar procedure. We first use the self-consistency condition in Eq. (S72) to rewrite the scaled fixation probability in a similar form as Eq. (S70),

$$\tilde{p}_{\text{fix}}(s) = \int N f(x-s)w(x) dx \approx \int_{-\infty}^{x_c+s} e^{\frac{xs}{v} - \frac{s^2}{2v}} \cdot \frac{w(x)}{2x_c e^{-\frac{x^2-x_c^2}{2v}}} \cdot \frac{dx}{s_b}. \quad (\text{S73})$$

This differs from the previous integral in Eq. (S70) by the presence of the exponential factor $e^{xs/v}$ and the shifting of the upper limit of integration from x_c to $x_c + s$. The behavior of this integral will therefore depend on the sign of s .

Beneficial mutations. When $s > 0$, the $e^{xs/v}$ term will enhance the contributions from higher fitness values, so the $n \geq 1$ terms in Eq. (S68) will continue to be negligible. Moreover, since the upper limit of integration is now larger than x_c , there will also be a new contribution from the Haldane region of the

shoulder solution, so that

$$\begin{aligned}\tilde{p}_{\text{fix}}(s) &\approx e^{\frac{x_c s}{v} - \frac{s^2}{2v}} \int_{x_c - s_b}^{x_c} e^{-\frac{(x_c - x)s}{v}} \frac{dx}{s_b} + \int_{x_c}^{x_c + s} \frac{2x e^{-\frac{(x-s)^2}{2v}}}{2x_c s_b e^{-\frac{x_c^2}{2v}}}, \\ &\approx e^{\frac{x_c s}{v} - \frac{s^2}{2v}} \cdot \frac{1 - e^{-\frac{s_b s}{v}}}{\frac{s_b s}{v}} + 2Ns \cdot \Phi\left(\frac{s - x_c}{\sqrt{v}}\right),\end{aligned}\quad (\text{S74})$$

where $\Phi(z) \equiv \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$ is the Gaussian cumulative distribution function. The leading order scaling is given by

$$\log \tilde{p}_{\text{fix}}(s) \approx \begin{cases} \frac{x_c s}{v} & \text{if } 0 < s \ll x_c, \\ \log(2Ns) & \text{if } s \gg x_c, \end{cases}\quad (\text{S75})$$

which transitions between a regime of strong clonal interference for $s \ll x_c$ and the Haldane limit when $s \gg x_c$ (42).

The fixation probabilities in Eq. (S75) have a simple heuristic interpretation. To be a successful, a moderately beneficial mutation ($s \ll x_c$) will typically arise on a background near $x = x_c - s$. This mutation will found a new lineage which experiences the effects of clonal interference and must generate multiple more fit descendants to take over and fix. Over this fixation time, $T_{\text{sw}} = x_c/v$, the beneficial mutation will provide an exponential advantage to its founding lineage. On the other hand, successful mutations that are extremely beneficial ($s \gg x_c$) will typically land above the interference threshold ($x > x_c - s$) and fix if they survive genetic drift.

Deleterious mutations. The fixation probabilities of deleterious mutations ($s < 0$) can be calculated using a similar procedure. In this case, however, the $e^{-x|s|/v}$ term becomes increasingly large at lower fitness values, and will need to be cut off by one of the $n \geq 1$ terms in Eq. (S68). To build intuition, let us first consider the case where $|s| < s_b$, so that

$$\tilde{p}_{\text{fix}}(s) \approx e^{-\frac{s^2}{2v}} \left[\int_{x_c - s_b}^{x_c - |s|} e^{-\frac{x|s|}{v}} \frac{dx}{s_b} + \left(\frac{x_c - s_b}{x_c}\right) \int_{x_c - 2s_b}^{x_c - s_b} e^{-\frac{(x_c - s_b - x)s_b}{v} - \frac{x|s|}{v}} \frac{dx}{s_b} + \dots \right].\quad (\text{S76})$$

The contribution from the $n = 0$ term is now peaked at its lower limit of integration ($x \approx x_c - s_b$). However, when $|s|$ is smaller than s_b , the contribution from the $n = 1$ term is still peaked at its upper limit of integration ($x \approx x_c - s_b$), which implies that the total contribution to the fixation probability will also be peaked at $x \approx x_c - s_b$. [The contributions from the $n \geq 2$ terms will be smaller by additional factors of $\exp(n(n-1)s_b^2/2v)$, and will therefore be negligible in the limit that $s_b \gg \sqrt{v}$.] Evaluating the two integrals then yields

$$\tilde{p}_{\text{fix}}(s) \approx e^{-\frac{(x_c - s_b)|s|}{v} - \frac{|s|^2}{2v}} \left(\frac{1 - e^{-\frac{(s_b - |s|)|s|}{v}}}{\frac{s_b |s|}{v}} + \left(\frac{x_c - s_b}{x_c}\right) \frac{1 - e^{-\frac{s_b(s_b - |s|)}{v}}}{\frac{s_b(s_b - |s|)}{v}} \right),\quad (\text{S77})$$

which is valid for fitness costs $s \in (-s_b + \mathcal{O}(v/s_b), 0)$. When $|s| \ll s_b$, this matches the fixation probability derived using the thresholding approximation in Ref. (50), but it starts to deviate from Ref. (50) for $|s| \sim \mathcal{O}(s_b)$. We expect that the expression in Eq. (S77) will be more accurate in this case, since it better captures the behavior of $w(x)$ below $x \approx x_c - s_b$.

The fixation probabilities of more strongly deleterious mutations can be calculated in a similar manner. In this case, it will be convenient to write

$$|s| = ks_b + \Delta,\quad (\text{S78})$$

where k is the largest integer such that $\Delta > 0$. The upper limit of integration at $x_c - s \equiv x_c - ks_b - \Delta$ implies that the first term with nonzero contribution to p_{fix} will be the $n = k$ term, so that

$$\begin{aligned} \tilde{p}_{\text{fix}}(s) \approx & \frac{e^{-\frac{s^2}{2v}}}{k!} \left(1 - \frac{s_b}{x_c}\right)^k e^{-\frac{k(k-1)s_b^2}{2v}} \left[\int_{x_c - (k+1)s_b}^{x_c - ks_b - \Delta} e^{-\frac{x(ks_b + \Delta)}{v} - \frac{ks_b(x_c - ks_b - x)}{v}} \frac{dx}{s_b} \right. \\ & \left. + \frac{e^{-\frac{ks_b^2}{v}}}{k+1} \left(\frac{x_c - s_b}{x_c}\right) \int_{x_c - (k+2)s_b}^{x_c - (k+1)s_b} e^{-\frac{x(ks_b + \Delta)}{v} - \frac{(k+1)s_b(x_c - (k+1)s_b - x)}{v}} \frac{dx}{s_b} + \dots \right]. \end{aligned} \quad (\text{S79})$$

As above, the contribution from the $n = k$ term is peaked at the lower limit of integration ($x \approx x_c - (k+1)s_b$), while the contribution from the $n = k+1$ term is peaked at its upper limit ($x \approx x_c - (k+1)s_b$). The terms with $n \geq k+2$ are smaller by exponential factors of $\mathcal{O}(s_b^2/v)$, and will therefore provide a negligible contribution to the fixation probability when $s_b \gg \sqrt{v}$. Evaluating the two integrals then yields,

$$\begin{aligned} \tilde{p}_{\text{fix}}(s) \approx & \frac{e^{-\frac{s^2}{2v}}}{k!} \left(1 - \frac{s_b}{x_c}\right)^k e^{-\frac{k(k-1)s_b^2}{2v}} \left[e^{-\frac{(x_c - (k+1)s_b)\Delta}{v} - \frac{ks_b(x_c - ks_b)}{v}} \left(\frac{1 - e^{-\frac{\Delta(s_b - \Delta)}{v}}}{\frac{s_b\Delta}{v}}\right) \right. \\ & \left. + \frac{e^{-\frac{ks_b^2}{v}}}{k+1} \left(\frac{x_c - s_b}{x_c}\right) e^{-\frac{(x_c - (k+1)s_b)(ks_b + \Delta)}{v}} \left(\frac{1 - e^{-\frac{s_b(s_b - \Delta)}{v}}}{\frac{s_b(s_b - \Delta)}{v}}\right) \right], \\ \approx & e^{-\frac{(x_c - \frac{s_b}{2})|s|}{v}} \frac{e^{-\frac{\Delta^2}{2v} + \frac{s_b\Delta}{2v}}}{k!} \left(1 - \frac{s_b}{x_c}\right)^k \left[\left(\frac{1 - e^{-\frac{\Delta(s_b - \Delta)}{v}}}{\frac{s_b\Delta}{v}}\right) + \frac{\left(\frac{x_c - s_b}{x_c}\right)}{k+1} \left(\frac{1 - e^{-\frac{s_b(s_b - \Delta)}{v}}}{\frac{s_b(s_b - \Delta)}{v}}\right) \right], \end{aligned} \quad (\text{S80})$$

whose leading order behavior simplifies to

$$\log \tilde{p}_{\text{fix}}(s) \approx \begin{cases} -\frac{(x_c - s_b)|s|}{v} & \text{if } s < 0 \text{ and } |s| \ll s_b, \\ -\frac{(x_c - \frac{s_b}{2})|s|}{v} & \text{if } s < 0 \text{ and } |s| \gg s_b. \end{cases} \quad (\text{S81})$$

For the purposes of numerical evaluation, it is useful to employ a modified version of Eq. (S80),

$$\begin{aligned} \tilde{p}_{\text{fix}}(s) \approx & e^{-\frac{(x_c - \frac{s_b}{2})|s|}{v}} \frac{e^{-\frac{\Delta^2}{2v} + \frac{s_b\Delta}{2v}}}{k!} \left(1 - \frac{s_b}{x_c}\right)^k \\ & \times \left[\left(\frac{1 - e^{-\frac{\Delta(s_b - \Delta)}{v}}}{\frac{s_b\Delta}{v}}\right) \left(1 - e^{-\frac{s_b(s_b - \Delta)}{v}}\right) + \frac{\left(\frac{x_c - s_b}{x_c}\right)}{k+1} \left(\frac{1 - e^{-\frac{s_b(s_b - \Delta)}{v}}}{\frac{s_b(s_b - \Delta)}{v}}\right) \left(1 - e^{-\frac{\Delta s_b}{v}}\right) \right], \end{aligned} \quad (\text{S82})$$

which has the same asymptotic limit when $v \ll s_b^2$, but enforces strict continuity at $\Delta = s_b$.

4.1.4 Modifiers without direct costs or benefits

We now are in a position to calculate the fixation probabilities of modifier mutations ($\mu_m(s) \neq \mu(s)$). In the absence of a direct cost or benefit ($s_m = 0$), this will only be a slight generalization of the neutral fixation probability calculation in SI Section 4.1.3. In the case of a modifier mutation, Eq. (S70) becomes

$$\tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) = \int Nf(x)w_m(x) \approx \frac{x_c m s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \int_{-\infty}^{x_c} \frac{w_m(x)}{2x_c m e^{\frac{x^2 - x_{cm}^2}{2v}}} \frac{dx}{s'_b}, \quad (\text{S83})$$

where the primary difference is the presence of the prefactor term, and the fact that the upper limit of integration is now equal to $x_c \neq x_{cm}$. The behavior of this integral will sensitively depend on the relative magnitudes of x_{cm} and x_c .

Positively selected modifiers. If the modifier interference threshold is less than that of the wildtype ($x_{cm} < x_c$), then the integral in Eq. (S83) will contain both the $n \geq 0$ terms as well as a portion of the Haldane region,

$$\begin{aligned} \tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) &\approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \left[\int_{-\infty}^{x_{cm}} \frac{w_m(x)}{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} s'_b} dx + \int_{x_{cm}}^{x_c} \frac{2x}{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} s'_b} dx \right], \\ &\approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \left[1 + \frac{v}{x_{cm}s'_b} \left(1 - e^{\frac{x_{cm}^2 - x_c^2}{2v}} \right) \right]. \end{aligned} \quad (\text{S84})$$

The second term is negligible in our regime of interest where $x_{cm}s'_b/v \gg 1$, so the fixation probability reduces to

$$\tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) \approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}}. \quad (\text{S85})$$

Since $x_{cm} < x_c$, Eq. (S85) implies that $\tilde{p}_{\text{fix}} > 1$ (i.e., the modifier is favored by natural selection).

Negatively selected modifiers. In the opposite case, where the interference threshold of the modifier is greater than that of the wildtype ($x_{cm} > x_c$), the upper limit of the integral in Eq. (S83) will occur somewhere in the interference region of $w_m(x)$. The precise behavior will depend on how far x_c extends into this region. If $x_c > x_{cm} - s'_b + \mathcal{O}(v/s'_b)$, then the dominant contribution will still come from the $n = 0$ region of Eq. (S68), so that

$$\tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) \approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \left(1 - \frac{x_{cm} - x_c}{s'_b} \right). \quad (\text{S86})$$

Since $x_{cm} > x_c$, this implies that $p_{\text{fix}}/p_0 < 1$ (i.e., the modifier is disfavored by natural selection). More generally, if

$$x_c = x_{cm} - ks'_b - \Delta, \quad (\text{S87})$$

for some $\Delta \in (0, s'_b)$, then the primary contribution will come from the $n = k$ region of Eq. (S68). This yields

$$\begin{aligned} \tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) &\approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \frac{\left(\frac{x_{cm} - s'_b}{x_{cm}} \right)^k}{k!} e^{-\frac{k(k-1)s'^2_b}{2v}} \int_{x_{cm} - (k+1)s'_b}^{x_{cm} - ks'_b - \Delta} e^{-\frac{ks'_b(x_{cm} - ks'_b - x)}{v}} \frac{dx}{s'_b}, \\ &\approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \frac{\left(\frac{x_{cm} - s'_b}{x_{cm}} \right)^k}{k!} e^{-\frac{k(k-1)s'^2_b}{2v}} \frac{v}{ks'^2_b} e^{-\frac{ks'_b\Delta}{v}}. \end{aligned} \quad (\text{S88})$$

Predictions for specific values of s'_b and U'_b . We can express the fixation probability of a modifier in terms of the underlying parameters, s'_b , U'_b , and N by substituting Eq. (S59) into Eq. (S85) to obtain,

$$\begin{aligned} \log \tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) \approx & \left(1 - \frac{s_b^2}{s'^2_b}\right) \left(\frac{x_c^2}{2v} - \frac{x_c s_b}{2v}\right) + \frac{x_c}{s_b} \frac{s_b^2}{s'^2_b} \log\left(\frac{U'_b}{U_b}\right) - \frac{v}{2s'^2_b} \log^2\left(\frac{U'_b}{U_b}\right) \\ & + \frac{1}{2} \left(1 - \frac{s_b^2}{s'^2_b}\right) \log\left(\frac{U'_b}{U_b}\right), \end{aligned} \quad (\text{S89})$$

where we have assumed that $x_{cm} - s'_b \gg \sqrt{v}$. In the case of a selection strength modifier ($U'_b = U_b$), the leading order contributions simplify to

$$\log \tilde{p}_{\text{fix}}(s_b \rightarrow s'_b) \approx \left(1 - \frac{s_b^2}{s'^2_b}\right) \left(\frac{x_c^2}{2v} - \frac{x_c s_b}{2v}\right) \approx \left(1 - \frac{s_b^2}{s'^2_b}\right) \left(\log\sqrt{(Ns_b)(NU_b)}\right), \quad (\text{S90})$$

where we have substituted $x_c^2/2v \approx \log(Ns_b)$ and $x_c s_b/v \approx \log(s_b/U_b)$ on the right hand side.

The fixation probability of a selection strength modifier has a simple heuristic explanation. For small changes in the selection coefficient ($s'_b - s_b \ll s'_b$), successful mutations arise in the high-fitness ‘‘nose’’ of $f(x)$ ($x \approx x_c \gg s_b$) and must acquire $\sim x_c/s_b$ additional mutations before they reach $\mathcal{O}(1)$ frequencies. In each of these steps, a selection-strength modifier produces $\sim \exp\left[\frac{s_b}{v} \cdot j(s'_b - s_b)\right]$ more mutations than a wildtype individual with the same fitness, leading to the exponential scaling observed in Eq. (3).

Similarly, the fixation probability of a mutation rate modifier is given by setting $s'_b = s_b$ in Eq. (S89), which yields the leading order contribution

$$\log \tilde{p}_{\text{fix}}(U_b \rightarrow U'_b) \approx \frac{x_c}{s_b} \log\left(\frac{U'_b}{U_b}\right) - \frac{v}{2s_b^2} \log^2\left(\frac{U'_b}{U_b}\right) \approx \frac{x_c}{s_b} \log\left(\frac{U'_b}{U_b}\right), \quad (\text{S91})$$

where we have assumed that $x_c/s_b \approx 2\log(Ns_b)/\log(s_b/U_b)$ and $\log(U'_b/U_b) \ll x_c s_b/v \approx \log(s_b/U_b) \gg 1$. This matches the result previously derived in Ref. (26). The heuristic picture is similar to the one described in that work: a mutation rate modifier must also arise in the high-fitness nose of the fitness distribution ($x \approx x_c$) and accrue x_c/s_b additional mutations to fix. In each of these x_c/s_b steps before a coalescent event, the mutation rate modifier produces U'_b/U_b more offspring, leading to the observed scaling in Eq. (S91).

More generally, we see that the fixation probability of a joint selection strength and mutation rate modifier can be expressed in terms of the fixation probabilities of these stand-alone modifiers,

$$\log \tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) \approx \log \tilde{p}_{\text{fix}}(s_b \rightarrow s'_b) + \alpha \cdot \log \tilde{p}_{\text{fix}}(U_b \rightarrow U'_b), \quad (\text{S92})$$

where $\alpha = s_b^2/s'^2_b$. The presence of this factor α indicates that a change to the selection strength of mutations can modulate the effect of a linked mutation rate change. The effect of this interplay can be extremely important to determining the fate a modifier mutation, causing modifiers that would be disfavored by natural selection without this modulation to be favored. Interestingly, the modulation effect only works in this direction: modifiers that would be favored without the effect cannot actually be deleterious.

4.1.5 Modifiers with direct costs or benefits

Finally, we can extend these calculations to modifiers with direct costs or benefits by combining our results from SI Sections 4.1.3 and 4.1.4. Generalizing Eqs. (S73) or (S83) to this case yields

$$\begin{aligned} \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &= \int N f(x - s_m) w_m(x) dx \\ &\approx \frac{x_{cm} s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v} - \frac{s_m^2}{2v}} \int_{-\infty}^{x_c + s_m} e^{\frac{x s_m}{v}} \cdot \frac{w_m(x)}{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}}} \cdot \frac{dx}{s'_b}. \end{aligned} \quad (\text{S93})$$

Our analysis above suggests that the two key considerations will be (i) the sign of s_m , and (ii) where $x_c + s_m$ falls relative to the internal scales of $w_m(x)$ (e.g. x_{cm} , $x_{cm} - s'_b$, etc.).

Modifiers with direct fitness benefits. For a modifier with a direct fitness benefit ($s_m > 0$), the exponential term will once again amplify the contributions from higher fitness values. The behavior of Eq. (S93) will therefore strongly depend on the location of the upper limit of integration. If $s_m - (x_{cm} - x_c) > \mathcal{O}(v/x_{cm})$, there will be a new contribution from the Haldane region of $w_m(x)$, so that

$$p_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{x_{cm} s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v} - \frac{s_m^2}{2v}} \int_{x_{cm} - s_b}^{x_{cm}} e^{\frac{x s_m}{v}} \cdot \frac{dx}{s'_b} + \int_{x_{cm}}^{x_{cm} + s_m - (x_{cm} - x_c)} \frac{2N x e^{-\frac{(x - s_m)^2}{2v}}}{\sqrt{2\pi v}} dx. \quad (\text{S94})$$

Since $x_c \gg \sqrt{v}$, the upper limit from the Haldane region can be taken to infinity, so that Eq. (S94) reduces to

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{x_{cm} s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \left(e^{\frac{x_{cm} s_m}{v} - \frac{s_m^2}{2v}} \cdot \frac{1 - e^{-\frac{s'_b s_m}{v}}}{\frac{s'_b s_m}{v}} \right) + 2N s_m \Phi \left(\frac{s_m - x_{cm}}{\sqrt{v}} \right). \quad (\text{S95})$$

In the opposite case, where $s_m \lesssim x_{cm} - x_c - \mathcal{O}(v/x_{cm})$, the upper limit of integration in Eq. (S93) will fall below the interference threshold of $w_m(x)$. When $s_m > 0$, this will only occur when $x_{cm} > x_c$ (i.e., when the modifier would be disfavored on its own). Similar to the pure modifier case in Eqs. (S86) and (S88), the behavior of Eq. (S93) will depend on where $x_c + s_m$ falls relative to the internal scales of $w_m(x)$. If $x_{cm} \gtrsim x_c + s_m \gtrsim x_{cm} - s'_b$, then the dominant contribution will come from the $n = 0$ region of Eq. (S68), so that

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{x_{cm} s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \cdot e^{\frac{x_{cm} s_m}{v} - \frac{s_m^2}{2v}} \cdot \frac{e^{-\frac{(x_{cm} - x_c - s_m) s_m}{v}} \left(1 - e^{-\frac{(s'_b + x_c - x_{cm} - s_m) s_m}{v}} \right)}{\frac{s'_b s_m}{v}}. \quad (\text{S96})$$

More generally, if we have

$$x_c + s_m = x_{cm} - k s'_b - \Delta, \quad (\text{S97})$$

for some $\Delta \in (0, s'_b)$, then the primary contribution will come from the $n = k$ region of Eq. (S68). This

yields

$$\begin{aligned}
\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &\approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \frac{\left(\frac{x_{cm} - s'_b}{x_{cm}}\right)^k}{k!} e^{-\frac{k(k-1)s_b'^2}{2v}} \int_{x_{cm} - (k+1)s'_b}^{x_{cm} - ks'_b - \Delta} e^{\frac{xsm}{v} - \frac{s_m^2}{2v} - \frac{ks'_b(x_{cm} - ks'_b - x)}{v}} \frac{dx}{s'_b}, \\
&\approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v} + \frac{x_{cm}s_m}{v} - \frac{s_m^2}{2v}} \frac{\left(\frac{x_{cm} - s'_b}{x_{cm}}\right)^k e^{-\frac{k(k-1)s_b'^2}{2v} - \frac{(ks'_b + \Delta)s_m}{v} - \frac{ks'_b\Delta}{v}} \left(1 - e^{-\frac{(s_m + ks'_b)(s'_b - \Delta)}{v}}\right)}{k! \left(\frac{(s_m + ks'_b)s'_b}{v}\right)}.
\end{aligned} \tag{S98}$$

Modifiers with direct fitness costs. For a modifier with a direct fitness cost ($s_m < 0$), the exponential term will once again amplify contributions from lower fitness values, and must eventually be cut off by the $n = k + 1$ term in Eq. (S68), where k is defined by

$$|s_m| = ks'_b + \Delta, \tag{S99}$$

for some $\Delta \in (0, s'_b)$. The major difference from the first-order selection scenario in Eq. (S82) is that the upper limit of integration includes an additional term, $x_c - x_{cm}$. If $x_{cm} < x_c$ (i.e. a positively selected modifier), then the upper limit of integration is larger than in Eq. (S82). This can in principle lead to contributions from the $n < k$ terms in Eq. (S68). However, the $n < k$ terms will all be smaller by exponential factors of $\mathcal{O}(s_b^2/v)$ so they will provide a negligible contribution when $s_b \gg \sqrt{v}$. The larger upper limit of integration will therefore only alter the $n = k$ integral in Eq. (S79), shifting it from $x_{cm} - ks'_b - \Delta \rightarrow \min\{x_c - ks'_b - \Delta, x_{cm} - ks'_b\}$. This implies that the modifier fixation probability will be given by a slight generalization of Eq. (S82),

$$\begin{aligned}
\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &\approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} e^{-\frac{(x_{cm} - \frac{s'_b}{2})|s_m|}{v}} \frac{e^{-\frac{\Delta^2}{2v} + \frac{s'_b\Delta}{2v}}}{k!} \left(1 - \frac{s'_b}{x_{cm}}\right)^k \\
&\times \left[\left(\frac{1 - e^{-\frac{\Delta \cdot \min(s'_b, s'_b - \Delta + x_c - x_{cm})}{v}}}{\frac{s'_b\Delta}{v}} \right) \cdot \left(1 - e^{-\frac{s'_b(s'_b - \Delta)}{v}}\right) \right. \\
&\quad \left. + \frac{\left(\frac{x_{cm} - s'_b}{x_{cm}}\right) \left(1 - e^{-\frac{s'_b(s'_b - \Delta)}{v}}\right) \left(1 - e^{-\frac{\Delta s'_b}{v}}\right)}{k + 1} \right],
\end{aligned} \tag{S100}$$

which will be valid for $x_{cm} < x_c$. The opposite scenario ($x_{cm} > x_c$) corresponds to a negatively selected modifier that has an additional a direct cost. We will not consider such mutations here, as they will have a negligible chance of fixing.

Leading order scaling. Combining these expressions, we see that the leading-order solution for the fixation probability of a modifier with a direct cost or benefit can be summarized as,

$$\log \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \begin{cases} \frac{x_c^2}{2v} - \frac{x_{cm}^2}{2v} + \frac{(x_{cm} - s'_b)s_m}{v} & s'_b \gg |s_m|; s_m < 0, \\ \frac{x_c^2}{2v} - \frac{x_{cm}^2}{2v} + \frac{(x_{cm} - s'_b/2)s_m}{v} & s'_b \ll |s_m|, s_m < 0, \\ \frac{x_c^2}{2v} - \frac{x_{cm}^2}{2v} + \frac{x_{cm}s_m}{v} & x_c \gg |s_m|, s_m > 0. \end{cases} \tag{S101}$$

Interestingly, we see that this fixation probability decomposes into contributions from first- and second-order selection,

$$\log \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \log \tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) + \gamma \cdot \log \tilde{p}_{\text{fix}}(s_m), \quad (\text{S102})$$

where $\tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b))$ is the fixation probability of a modifier without direct costs or benefits, $\tilde{p}_{\text{fix}}(s_m)$ is the fixation probability of a first order mutation, and the weighting factor γ satisfies $\gamma = \min(x_{cm}/x_c, \sqrt{2s_b/x_c})$. The presence of this factor γ demonstrates that second order selection modulates the effects of first order selection; the fixation probability of a deleterious or beneficial mutation arising in the background population scales with x_c/v , while the contribution of a direct cost or benefit scales with x_{cm}/v . This means that a modifier without direct costs or benefits that is favored by natural selection ($x_{cm} < x_c$) suppresses the contribution of a direct cost, while a modifier without direct costs or benefits that is disfavored by natural selection ($x_{cm} > x_c$) amplifies the contribution of a direct benefit. This effect can be extremely large, affecting the contribution of a direct cost or benefit and overall favorability of a modifier mutation by orders of magnitude.

For sufficiently large direct benefits ($|s_m| \gg x_c$), we expect that even the most deleterious modifiers must fix deterministically upon surviving genetic drift, irrespective of the background on which they arise ($p_{\text{fix}}/p_0 \approx 2Ns_m$). This expectation follows from the fact that a modifier with a large direct benefit will jump far out ahead of the 'traveling wave' and take over the population long before the background population is able to 'catch-up'. This is true even in the extreme case of a dead-end modifier, which would drive the rate of adaptation to zero upon fixing. Our current theoretical approach, however, does not enable us to predict the fixation probability of a modifier with direct benefits $s_m \sim O(x_c)$. In particular, to derive Eq. (S34) we assumed that the fate of a modifier mutation is determined at small frequency while competing against the mean fitness set by the rate of adaptation of the background population. However, if the rate of adaptation does not change, our theory predicts the background population will always "catch up" and pass a deleterious modifier in fitness, no matter how large the direct benefit. This pathology can be seen in Eq. (S57), where $\lim_{s'_b \rightarrow 0} x_{cm} = \infty$.

This limitation prevents us from understanding a key aspect of the trade-off between short-term fitness and long-term evolvability. How large a direct benefit is sufficient to drive a deleterious modifier to be favored? In the next section, we focus on the extreme case of an evolutionary dead-end ($U'_b = 0$ or $s'_b = 0$), showing that even modest direct benefits can drive these modifiers to be favored by selection.

4.1.6 Fixation of a dead-end modifier

When a dead-end modifier arises with an extremely large direct benefit ($|s_m| \gg x_c$), it will sweep long before the background population is able to catch-up in fitness. On the other hand, if the direct benefit is not sufficiently large, the background population will surpass the modifier lineage in fitness. We can understand this cross-over by analyzing the dynamics of a modifier lineage arising with landing fitness, x , as it transitions from small to large frequency.

A newly established modifier clone with relative fitness x will initially start to grow deterministically as

$$f_m(t) \approx \frac{e^{xt - \frac{vt^2}{2}}}{2Nx}, \quad (\text{S103})$$

where $1/2Nx$ corresponds to the size of the clone immediately after it establishes. This will be a good approximation as long as the modifier frequency remains small ($f_m \ll 1$), so that the mean fitness of the population is still primarily determined by the wildtype ($\partial_t \bar{X}(t) \approx v$). If $x \ll \sqrt{2v \log(N\sqrt{v})}$, then the

modifier clone will remain small throughout its entire lifetime, and the wildtype population will eventually pass it by. However, for larger initial fitnesses, the modifier clone will eventually grow to $O(1)$ frequencies, and will start to exert its own effect on the population mean fitness. At this point, the dynamics in Eq. (S103) will start to break down.

The ultimate fate of a dead-end modifier will depend on how much fitness the wildtype population has gained during this time. To understand this process, let t^* denote the time required for the modifier clone to reach a reference frequency f_m^* . If $f_m^* \ll 1$, then t^* can be calculated from Eq. (S103):

$$\log f_m^* = -\log(2Nx) + xt^* - \frac{vt^{*2}}{2}. \quad (\text{S104})$$

During this time, the wildtype population will have increased in fitness by a total amount $\Delta x = vt^*$, so that the current relative fitness of the modifier is $x - vt^*$. If the nose of the wildtype population has surpassed the modifier in this time ($x_c \gtrsim x - vt^*$), then regardless of how much the modifier grows in the short-term, it will be destined for extinction in the long-term, since it is unable to produce additional mutations.

On the other hand, if $x - vt^* > x_c$, then it is possible for the modifier to sweep through the rest of population extremely rapidly, and “freeze” the wildtype population in its place. If further adaptation of the wildtype can be neglected, then the modifier will transition from frequency f_m^* to $1 - f_m^*$ according to the logistic dynamics,

$$\frac{\partial f_m}{\partial t} = (x - vt^*)f_m(1 - f_m), \quad (\text{S105})$$

which requires an additional time interval

$$\Delta t \approx \frac{2}{x - vt^*} \log\left(\frac{1}{f_m^*}\right). \quad (\text{S106})$$

Since we have assumed that $x - vt^* \geq x_c$, it is always possible to choose a reference frequency f_m^* that is $\ll 1$, but large enough that Eq. (S106) is much smaller than the time required for the wildtype population to acquire one additional mutation (s_b/v). For example, choosing $f_m^* = \exp(-\sqrt{x_c s_b/v}) \ll 1$ yields

$$\frac{v\Delta t}{s_b} \leq \frac{2v}{x_c s_b} \sqrt{\frac{x_c s_b}{v}} \ll 1, \quad (\text{S107})$$

since $x_c s_b/v \gg 1$. This shows that the wildtype population is effectively frozen in place while the modifier transitions from f_m^* to $1 - f_m^*$. At this point, the mean fitness of the population is dominated by the modifier ($\bar{X} \approx x - vt^*$). Since we have assumed that $x - vt^* > x_c$, this implies that even the fittest individuals in the wildtype population now have a negative relative fitness, so their further adaptation will be effectively halted. The modifier clone will therefore continue to sweep through the population, and will fix in a time of order $\sim (x - vt^*)^{-1} \log(N(x - vt^*)f_m^*)$.

Based on this reasoning, we conclude that a newly established modifier will fix if its initial fitness exceeds a critical threshold x^* defined by

$$x^* - vt^* \equiv x_c + \Delta x, \quad (\text{S108})$$

where Δx is a small correction term [$\lesssim O(s_b)$] that accounts for potential ambiguities in the definition of the nose, as well as the possibility that a leapfrogged nose may produce one additional mutation via stochastic tunneling (100) before it goes extinct. It also accounts for the fact that a small subclass of individuals in the

background population will need to have fitness greater than x^* to out-compete the modifier after the mean fitness changes. The size of this subclass can be determined by requiring that the sum of the establishment probabilities of the individuals that have crossed x^* be equal to 1,

$$1 \approx \int_{x^*}^{x_c + vt^*} \frac{2N(x - x^*)}{\sqrt{2\pi v}} e^{-\frac{(x - vt^*)^2}{2v}}. \quad (\text{S109})$$

Assuming $v/x_c \ll x^* - vt^* - x_c$ and $x_c - (x^* - vt^* - x_c) \gg \sqrt{v}$, we can evaluate this integral with a Laplace approximation to obtain the self-consistent solution,

$$\Delta x \equiv x^* - vt^* - x_c \approx -\frac{v}{x_c} \log \left(\frac{x_c s_b}{v} \cdot \frac{x_c^2}{v} \right). \quad (\text{S110})$$

This constitutes a small correction to x^* (as expected), but we will see that it provides an important contribution to the fixation probability below. To solve for x^* , we can substitute Eq. (S110) into Eq. (S104) to obtain the condition,

$$\frac{x_c^2}{2v} \left[\left(\frac{x^*}{x_c} \right)^2 - 2 - \frac{2\Delta x}{x_c} \right] = \log f_m^* - \log \left(\frac{s_b}{\sqrt{2\pi v}} \right) + \log \left(\frac{x^*}{x_c} \right), \quad (\text{S111})$$

where we have also used Eq. (S72) to substitute for N . In our regime of interest where $x_c \gg \sqrt{v}$, the terms on the left will be asymptotically larger than those on the right, and the dominant balance will be given by the terms in the square brackets, yielding

$$x^* \approx \sqrt{2} \cdot x_c \left(1 + \frac{\Delta x}{2x_c} \right) \approx \sqrt{2} \cdot x_c \left[1 - \frac{v}{2x_c^2} \log \left(\frac{x_c s_b}{v} \cdot \frac{x_c^2}{v} \right) \right]. \quad (\text{S112})$$

This result implies that a modifier lineage with relative fitness greater than $x^* \approx \sqrt{2} \cdot x_c$ will be guaranteed to fix as long as it survives genetic drift while rare — regardless of its ability to produce additional mutations. This will always be consistent with the wildtype $w(x)$ function from the branching process approximation in SI Section 4.1.3, since $x^* > x_c$. However, it suggests that the interference threshold for the modifier lineage must be capped at a maximum value,

$$x_{cm} \approx \min \{ x_{cm}^0, x^* \}, \quad (\text{S113})$$

where x_{cm}^0 is our original solution for the interference threshold derived in SI Section 4.1.1 above. This constitutes a simple modification of the original branching process model in Eq. (1) that can capture some of these non-linear feedback effects.

In the extreme case of a dead-end modifier ($\mu(s) \rightarrow 0$), this ansatz leads to a conditional fixation probability of the form

$$w_m(x) \approx \begin{cases} 2x & \text{if } x \geq x^*, \\ 0 & \text{else.} \end{cases} \quad (\text{S114})$$

Substituting this result into Eq. (S93) yields a prediction for the total fixation probability of a dead-end modifier with a direct benefit s_m :

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow 0, s_m) \approx \begin{cases} \int_{x^*}^{x_c + s_m} \frac{x e^{-\frac{(x - sm)^2}{2v} + \frac{x_c^2}{2v}}}{x_c s_b} dx & \text{if } s_m > x^* - x_c, \\ 0 & \text{else.} \end{cases} \quad (\text{S115})$$

Evaluating the integral in the upper branch yields,

$$\int_{x^*}^{x_c+s_m} \frac{x e^{-\frac{(x-s_m)^2}{2v} + \frac{x_c^2}{2v}}}{x_c s_b} dx = \frac{v}{x_c s_b} \left[e^{\frac{x_c^2 - (x^* - s_m)^2}{2v}} - 1 \right] + 2N s_m \left[\Phi\left(\frac{s_m - x^*}{\sqrt{v}}\right) - \Phi\left(\frac{-x_c}{\sqrt{v}}\right) \right], \quad (\text{S116})$$

where $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$ is the Gaussian cumulative distribution function. The leading-order scaling is given by

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow 0, s_m) \approx \begin{cases} 2N s_m & \text{if } s_m \gtrsim x^* + \mathcal{O}(\sqrt{v}), \\ \frac{v x^*}{x_c s_b (x^* - s_m)} e^{\frac{x_c (s_m - x^* + x_c)}{v} - \frac{(s_m - x^* + x_c)^2}{2v}} & \text{if } x^* - \mathcal{O}(\sqrt{v}) \gtrsim s_m \gtrsim x^* - x_c + \mathcal{O}(\sqrt{v}), \\ 0 & \text{else.} \end{cases} \quad (\text{S117})$$

These results suggests that modest direct benefits of size $s_m \gtrsim x^* - x_c \approx (\sqrt{2} - 1)x_c$ are sufficient to cause an evolutionary dead-end to be favored by natural selection. While this critical size is larger than the advantage of a single mutation (s_b), it is smaller than the total fitness variation in the population (x_c), and is only weakly dependent on NU_b . This contrasts with the traditional linear scaling observed in models with small numbers of competing loci (101), highlighting the unique features that arise in genome-wide models of adaptation.

Together, this analysis shows that the fate of a dead-end modifier is still determined by chance events that occur while the modifier lineage is at low frequency (e.g. its random genetic background and whether it survives genetic drift while rare) even though its ability to take over critically relies on non-linear feedbacks that occur while it is common.

4.1.7 Relation between the fixation probability and the long-term rate of adaptation

It is useful to compare these results with the deterministic modifier theory in SI Section 1, which predicts that modifiers will be favored if they increase the long-term mean fitness of the population. While this result was originally derived for non-adapting populations at mutation-selection balance, a natural extension of this idea to adapting populations might suggest that natural selection would favor modifiers that increase the long-term rate of adaptation (i.e. $v_m > v$).

Our results above show that this simple heuristic clearly breaks down for modifiers with direct costs or benefits (Fig. 3). This discrepancy is most dramatic for the ‘‘dead-end’’ modifiers in SI Section 4.1.6, which can be strongly favored by selection even though they lower the rate of adaptation to zero. Our analysis in SI Section 4.1.6 showed that the origin of this effect could be traced to the early fixation of the modifier lineage before its long-term costs are fully realized. These fixation events are completely neglected by the deterministic theory in SI Section 1: beneficial variants can grow to arbitrarily large frequencies in the short-term, but their long-term fitness gains will always override their initial costs or benefits. This highlights how finite population sizes can be critically important, even when the population size is very large (51).

To understand this relationship in more detail, we can combine our results in Eqs. (S59), (S72), and S101 to derive an approximate formula connecting p_{fix} , v_m , and s_m . In the limit that $x_c \gg s_b$ and $x_{cm} \gg s'_b$, we find that the leading-order contributions satisfy the approximate scaling,

$$\log \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{x_{cm}^2}{2v} \left(\frac{v_m}{v} - 1 \right) + \frac{x_{cm} s_m}{v}, \quad (\text{S118})$$

to leading order in the logarithm of \tilde{p}_{fix} . This result shows that in the absence of direct costs or benefits ($s_m = 0$), the sign of selection in our simple model is directly related to the modifier’s effect on the long-term

rate of adaptation ($v_m - v$). This is reminiscent of the mean fitness principle in SI Section 1. Equation (S118) allows us to generalize this result to modifiers with direct costs or benefits ($s_m \neq 0$). Interestingly, we find that Eq. (S118) has the same form as the deterministic prediction in Eq. (S5) if we take $\bar{X}_m \approx s_m + v_m t$ and $\bar{X}_w \approx vt$, as expected, but impose a finite integration time $T_{\max} = x_{cm}/v$. This time limit roughly coincides with the fixation time of a successful mutation. This allows us to recover a version of the mean fitness principle in SI Section 1: within our simple clonal interference model, natural selection will favor modifiers that produce a higher mean fitness within a typical fixation time.

We note, however, that Eq. (S118) is only an approximate formula, which captures the leading-order scaling of the logarithm of p_{fix} (rather than p_{fix} itself). The sub-leading contributions to the logarithm are often important in practice, since they can be large in an absolute sense ($\gtrsim 1$) and must be exponentiated to obtain p_{fix} . Thus, while Eq. (S118) can provide a useful rule of thumb, we use the full expressions in SI Section 4.1.5 when calculating our theoretical predictions in Figs. 2 and 3 in the main text.

4.2 Quasi-sweep regime

4.2.1 Location of the interference threshold

As we consider larger values of s'_b , we will eventually reach a regime where a single established driver mutation will be sufficient to drive a modifier lineage to fixation ($s'_b - x_{cm} \gg \sqrt{v}$). To solve for the shape of $w_m(x)$ in this regime, we first revisit the auxiliary condition for x_{cm} in Eq. (S54). When $s'_b - x_{cm} \gg \sqrt{v}$, this integral is dominated by fitnesses within $\mathcal{O}(\sqrt{v})$ of $y = s'_b$, where the shoulder solution has already transitioned to the linear Haldane limit. Performing a Gaussian Laplace approximation around this maximum yields a new auxiliary condition for x_{cm} ,

$$1 \approx \frac{\sqrt{2\pi} U'_b s'_b}{x_{cm} \sqrt{v}} e^{\frac{x_{cm}^2}{2v}}, \quad (\text{S119})$$

which is valid when $s'_b - x_{cm} \gg \sqrt{v}$ and $x_{cm} \gg \sqrt{v}$. Solving for x_{cm} , we obtain the leading order solution,

$$x_{cm} = \sqrt{2v \log\left(\frac{v}{U'_b s'_b}\right)}, \quad (\text{S120})$$

which is valid when

$$\frac{U'_b s'_b}{U_b s_b} \ll \frac{v}{U_b s_b} \sim \frac{s_b}{U_b} \frac{\log(N s_b)}{\log^2(s_b/U_b)}. \quad (\text{S121})$$

This is a broad regime when $U_b \ll s_b$, but violations of this condition can become important when we consider continuous distributions of fitness effects below (SI Section 5). It will also be useful to derive an alternative expression for x_{cm} by substituting the auxiliary condition for x_c in Eq. (S58) into Eq. (S119), yielding

$$\frac{x_{cm}^2}{2v} - \frac{x_c s_b}{v} + \frac{s_b^2}{2v} + \log\left(\frac{U'_b}{U_b} \cdot \frac{s'_b}{s_b}\right) + \log\left(\frac{\sqrt{2\pi} s_b^2}{x_{cm} \sqrt{v}}\right) \approx 0. \quad (\text{S122})$$

The leading order solution is given by

$$x_{cm} \approx \sqrt{(2x_c - s_b)s_b - 2v \log(U'_b/U_b)}, \quad (\text{S123})$$

where we have included the potential contribution from $\log(U'_b/U_b)$, but assumed that the corresponding $\log(s'_b/s_b)$ term is subdominant. The condition that $s'_b - x_{cm} \gg \sqrt{v}$ will be valid provided that

$$\frac{s'_b}{s_b} \gtrsim \sqrt{\frac{2x_c - s_b}{s_b} - \frac{2v}{s_b^2} \log(U'_b/U_b)}, \quad (\text{S124})$$

which complements the analogous condition for the multiple mutations regime in Eq. (S60). Thus, for a pure selection strength modifier, the quasi-sweep regime will occur for $s'_b \gg \sqrt{(2x_c - s_b)s_b}$, while a pure mutation rate modifier requires the more stringent condition, $\log(U'_b/U_b) \gtrsim (x_c - s_b)s_b/v$. This suggests that relatively modest changes to the selection strength will require the analysis described below.

4.2.2 Extending the shoulder solution to lower fitness values

Having fixed the location of x_{cm} , we can repeat the procedure in SI Section 4.1.2 to compute the shape of $w_m(x)$ for fitness values below x_{cm} . If $x + s'_b \gtrsim s_b + \mathcal{O}(\sqrt{v})$, then the dominant contribution to Eq. (S54) will be contained within the region of integration, so that we can substitute the auxilliary condition to obtain

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}}. \quad (\text{S125})$$

This once again coincides with the shoulder solution in Eq. (S43) in the region where $x < x_{cm} - \mathcal{O}(v/x_{cm})$. This derivation shows that the shoulder solution will continue to be valid for fitnesses as low as

$$x_{\min}(\mu_m(s)) = \mathcal{O}(\sqrt{v}), \quad (\text{S126})$$

which is significantly smaller than both s'_b and x_{cm} .

For fitnesses that are less than $-\mathcal{O}(\sqrt{v})$, the integral in Eq. (S53) will start to be dominated by the upper limit of integration at $x + s'_b$. Provided that $x + s'_b - \mathcal{O}(v/(s'_b - x_{cm})) \gtrsim x_{cm} + \mathcal{O}(v/x_{cm})$, then the dominant contribution to the integral will still come from the Haldane region of the shoulder solution. An exponential Laplace approximation then yields

$$w_m(x) \approx \frac{2U'_b(s'_b - |x|)}{|x|} \approx \frac{2x_{cm}\sqrt{v}}{\sqrt{2\pi}|x|} \left(1 - \frac{|x|}{s'_b}\right) e^{-\frac{x^2}{2v}}, \quad (\text{S127})$$

where we have substituted the auxilliary condition in Eq. (S119). To capture the behavior in the intermediate region around $x = \pm\mathcal{O}(\sqrt{v})$, we can turn to the full Gaussian integral,

$$\begin{aligned} w_m(x) &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \left[\frac{U'_b s'_b \sqrt{2\pi} v e^{\frac{x_{cm}^2}{2v}}}{v x_{cm}} \cdot \Phi\left(\frac{x}{\sqrt{v}}\right) - \frac{U'_b e^{\frac{x_{cm}^2 - x^2}{2v}}}{x_{cm}} \right], \\ &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \left[\Phi\left(\frac{x}{\sqrt{v}}\right) - \frac{\sqrt{v}}{s'_b \sqrt{2\pi}} e^{-\frac{x^2}{2v}} \right], \end{aligned} \quad (\text{S128})$$

where $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$ is the Gaussian cumulative distribution function. This reduces to Eqs. (S125) and (S127) in the appropriate limits, but captures the intermediate region where $x = \pm\mathcal{O}(\sqrt{v})$.

For fitness below $x_{cm} - s'_b - \mathcal{O}(v/x_{cm})$, the upper limit of the integral in Eq. (S53) will start to fall within the exponential region of the shoulder solution in Eq. (S43). However, since the latter smoothly maps on to

the Gaussian integral in Eq. (S128), we can use this solution to extend $w_m(x)$ down to $x \approx x_{cm} - 2s_b$:

$$\begin{aligned} w_m(x) &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \frac{U'_b}{v} \int_{-\infty}^{x+s'_b} e^{\frac{ys'_b}{v} - \frac{s'^2_b}{2v}} \left[\Phi\left(\frac{y}{\sqrt{v}}\right) - \frac{\sqrt{v}}{s'_b \sqrt{2\pi}} e^{-\frac{y^2}{2v}} \right] dy, \\ &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \left(\frac{U'_b}{s'_b}\right) \left[e^{\frac{xs'_b}{v} + \frac{s'^2_b}{2v}} \cdot \Phi\left(\frac{x+s'_b}{\sqrt{v}}\right) - \frac{2v}{|x|\sqrt{2\pi v}} e^{-\frac{x^2}{2v}} \right]. \end{aligned} \quad (\text{S129})$$

Iterating this procedure again yields

$$\begin{aligned} w_m(x) &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \left(\frac{U'_b}{s'_b}\right)^2 \frac{s'_b}{v} \left[\int_{-\infty}^{x+s'_b} e^{\frac{2ys'_b}{v}} \Phi\left(\frac{y+s'_b}{\sqrt{v}}\right) dy - \int_{-\infty}^{x+s'_b} \frac{2v}{y\sqrt{2\pi v}} e^{-\frac{(y-s'_b)^2}{2v}} dy \right], \\ &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \left(\frac{U'_b}{s'_b}\right)^2 \left[\frac{1}{2} e^{\frac{2xs'_b}{v} + \frac{2s'^2_b}{v}} \Phi\left(\frac{x+2s'_b}{\sqrt{v}}\right) - \frac{ve^{-\frac{x^2}{2v}}}{2|x|\sqrt{2\pi v}} - \frac{2s'_b v e^{-\frac{x^2}{2v}}}{|x+s_b||x|\sqrt{2\pi v}} \right], \end{aligned} \quad (\text{S130})$$

which will be valid for $x \in (x_{cm} - 3s'_b, x_{cm} - 2s'_b)$. More generally, for $x \in (x_{cm} - (n+1)s'_b, x_{cm} - ns'_b)$, one can show that the solution for $w_m(x)$ is given by:

$$\begin{aligned} w_m(x) &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \left(\frac{U'_b}{s'_b}\right)^n \left[\frac{1}{n!} e^{\frac{2nxs'_b}{2v} + \frac{n^2 s'^2_b}{2v}} \Phi\left(\frac{x+ns'_b}{\sqrt{v}}\right) \right. \\ &\quad \left. - \frac{ve^{-\frac{x^2}{2v}}}{s'_b \sqrt{2\pi v}} \left(\sum_{m=1}^n \frac{1}{m!} \prod_{j=0}^{n-m} \frac{s'_b}{|x+js'_b|} + \prod_{j=0}^{n-1} \frac{s'_b}{|x+js'_b|} \right) \right], \end{aligned} \quad (\text{S131})$$

which simplifies to,

$$\begin{aligned} w_m(x) &\approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \left(\frac{U'_b}{s'_b}\right)^n \left[\frac{1}{n!} e^{\frac{2nxs'_b}{2v} + \frac{n^2 s'^2_b}{2v}} \Phi\left(\frac{x+ns'_b}{\sqrt{v}}\right) \right. \\ &\quad \left. - \frac{ve^{-\frac{x^2}{2v}}}{s'_b \sqrt{2\pi v}} \left(\frac{-\Gamma\left(-\frac{x}{s'_b} - n\right) \left(n+1 + \frac{x}{s'_b}\right)}{\Gamma\left(1 - \frac{x}{s'_b}\right)} - \frac{s'_b}{n!(x+ns'_b)} \right) \right], \end{aligned} \quad (\text{S132})$$

and approaches the asymptotic limit,

$$w_m(x) \approx 2x_{cm} e^{\frac{-x^2_{cm}}{2v}} \left(\frac{U'_b}{s'_b}\right)^n \frac{v}{s'_b \sqrt{2\pi v}} \frac{1}{\Gamma\left(1 + \frac{|x|}{s'_b}\right)} \frac{x + (n+1)s'_b}{|x+ns'_b|} \approx \frac{2s'_b \left(\frac{U'_b}{s'_b}\right)^{n+1}}{\Gamma\left(1 + \frac{|x|}{s'_b}\right)} \frac{x + (n+1)s'_b}{|x+ns'_b|} \quad (\text{S133})$$

when $x \lesssim -ns'_b - \mathcal{O}(\sqrt{v})$. Each additional "step" reduces the original shoulder by a factor of $(U'_b/ns'_b)^n$. For fitness values that are many multiples of s'_b below x_{cm} , we can substitute $n \approx (x_{cm} - x)/s'_b$ to show that the leading order contribution to $w_m(x)$ scales as

$$\log w_m(x) \sim \text{const} + (x_{cm} - x)/s'_b \log(U'_b/(x_{cm} - x)), \quad (\text{S134})$$

which obeys the required boundary condition that $w_m(x) \rightarrow 0$ as $x \rightarrow \infty$. Thus Eq. (S132) provides an asymptotic solution for $w_m(x)$ that is valid across the full range of relative fitness below x_{cm} when $s'_b - x_{cm} \gg \sqrt{v}$.

4.2.3 Fixation probabilities of modifiers

Using the solution for $w_m(x)$ in Eq. (S132), we can repeat our calculations from SI Sections 4.1.4 and 4.1.5 to predict the fixation probabilities of modifier mutations in the quasi-sweep regime.

Modifiers without direct costs or benefits. In the absence of direct costs and benefits, the primary contribution to the fixation probability in Eq. (S70) will come from the $n = 0$ region of Eq. (S132). This yields

$$\tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) \approx \frac{x_{cm}s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \int_0^{\min\{x_c, x_{cm}\}} \frac{dx}{s_b} \approx \frac{x_{cm} \min\{x_c, x_{cm}\}}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}}. \quad (\text{S135})$$

Using Eqs. (S72) and (S119), we can rewrite this expression in the compact form,

$$\frac{2NU'_b \min\{x_c, x_{cm}\} s'_b}{v}. \quad (\text{S136})$$

This expression has a simple heuristic explanation. A modifier lineage founded at relative fitness x will produce a total of $NU'_b \int_0^\infty e^{xt - \frac{vt^2}{2}} dt$ offspring before it is purged from the population. The modifier lineage will produce the bulk of these offspring within $\sim 1/\sqrt{v}$ generations of time $t = x/v$, when it will have reached size $n(t = x/v) = e^{\frac{x^2}{2v}}$ and be near the mean fitness of the population. These offspring will fix provided they survive genetic drift, which occurs with probability $p_{\text{est}} \approx 2s'_b$. Multiplying these terms together, we see that the advantage of arising on a more fit background and hitchhiking to exponentially larger frequency is exactly balanced by the exponentially smaller probability of arising on one of these privileged backgrounds $f(x) \sim e^{-\frac{x^2}{2v}}$. Consequently, a successful modifier is equally likely to have arisen on a background with relative fitness in the range $\mathcal{O}(\sqrt{v}) < x < \min\{x_{cm}, x_c\}$.

By comparing Eq. (S91), Eq. (S90), and Eq. (S136) and Eq. (S89), we can see that there is a limit to the degree in which changes to the selection strength of mutations can modulate the effect of linked changes to mutation rate ($\log(U'_b/U_b) < (x_c - s_b)s_b/v$): $\alpha = \max\{(s'_b/s_b)^2, s_b/x_c\}$. When $s'_b - x_{cm} \gg \sqrt{v}$, a mutation rate modification influences only a single subsequent mutation compared to multiple mutations when $x_{cm} - s'_b \gg \sqrt{v}$. We also note that Eq. (S136) implies that large changes to the mutation rate ($\log(U'_b/U_b) < (x_c - s_b)s_b/v < (x_c - s_b)s_b/v + s_b^2/2v$) can modulate the effect of a selection-strength modifier, limiting its ability to compound.

Modifiers with direct fitness benefits. For a modifier with a direct fitness benefit ($s_m > 0$), the fixation probability in Eq. (S93) will once again depend on how the upper limit of integration ($x_c + s_m$) compares to the interference threshold x_{cm} . If $s_m \gtrsim x_{cm} - x_c + \mathcal{O}(v/x_{cm})$, there will again be a contribution from the Haldane region of $w_m(x)$, so that

$$\begin{aligned} \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &\approx \frac{x_{cm}^2}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v} - \frac{s_m^2}{2v}} \int_0^{x_{cm}} e^{\frac{x s_m}{v}} \cdot \frac{dx}{x_{cm}} + \int_{x_{cm}}^{x_{cm} + s_m - (x_{cm} - x_c)} \frac{2Nx e^{-\frac{(x-s_m)^2}{2v}}}{\sqrt{2\pi v}} dx, \\ &\approx \frac{x_{cm}^2}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2 - s_m^2}{2v}} \left(\frac{e^{\frac{x_{cm} s_m}{v}} - 1}{\frac{x_{cm} s_m}{v}} \right) + 2Ns_m \cdot \Phi\left(\frac{s_m - x_{cm}}{\sqrt{v}}\right), \end{aligned} \quad (\text{S137})$$

where $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$ is the Gaussian cumulative distribution function.

In the opposite case, where $s_m \lesssim x_{cm} - x_c - O(v/x_{cm})$, the upper limit of integration in Eq. (S93) will fall below the interference threshold of $w_m(x)$. Repeating the calculation in Eq. (S137) in this case yields

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{x_{cm}^2}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2 - s_m^2}{2v}} e^{\frac{(x_c + s_m)s_m}{v} - 1} \frac{1}{\frac{x_{cm}s_m}{v}}. \quad (\text{S138})$$

Modifiers with direct fitness costs. For a modifier with a beneficial effect on evolvability ($x_c > x_{cm}$) but a direct fitness cost ($s_m < 0$), the derivation will be similar to Eq. (S100) above. The dominant contribution to the fixation probability in Eq. (S93) will again come from the $n = k$ term in Eq. (S132) (as well as the first portion of the $k + 1$ term), where k is defined by

$$|s_m| = k s'_b + \Delta, \quad (\text{S139})$$

for some $\Delta \in (0, s'_b)$. The fixation probability of the modifier can therefore be approximated as

$$\begin{aligned} \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &\approx \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), 0) \cdot \left(\frac{U'_b}{s'_b}\right)^k \left\{ \frac{e^{-\frac{ks'_b\Delta}{v} - \frac{\Delta^2}{2v}}}{k!} \int_{x_{cm} - (k+1)s'_b}^{x_{cm} - ks'_b - \zeta} e^{-\frac{x\Delta}{v}} \Phi\left(\frac{x + ks'_b}{\sqrt{v}}\right) \frac{dx}{x_{cm}} \right. \\ &+ \frac{v}{s'_b \sqrt{2\pi v}} \int_{x_{cm} - (k+1)s'_b}^{x_{cm} - ks'_b - \zeta} e^{-\frac{(x-s_m)^2}{2v}} \cdot \frac{(x + (k+1)s'_b) \Gamma\left(-\frac{x}{s'_b} - k\right)}{s'_b \Gamma\left(1 - \frac{x}{s'_b}\right)} \frac{dx}{x_{cm}} \\ &+ \frac{v}{s'_b \sqrt{2\pi v}} \int_{x_{cm} - (k+1)s'_b}^{x_{cm} - ks'_b - \zeta} e^{-\frac{(x-s_m)^2}{2v}} \cdot \frac{1}{k!} \frac{s'_b}{(x + ks'_b)} \frac{dx}{x_{cm}} \\ &\left. + \frac{e^{-\frac{ks'_b\Delta}{v} - \frac{\Delta^2}{2v} + \frac{(2k+1)s'_b{}^2}{2v} - \frac{x_{cm}^2}{2v}}}{(k+1)!} \frac{x_{cm} \sqrt{v}}{s'_b{}^2 \sqrt{2\pi}} \int_{x_{cm} - (k+2)s'_b}^{x_{cm} - (k+1)s'_b} e^{\frac{x(s'_b - \Delta)}{v}} \Phi\left(\frac{x + (k+1)s'_b}{\sqrt{v}}\right) \frac{dx}{x_{cm}} \right\}, \quad (\text{S140}) \end{aligned}$$

where $\zeta \equiv \max\{0, \Delta - (x_c - x_{cm})\}$. When $k = 0$ and $\Delta \lesssim O(\sqrt{v})$, the fixation probability is dominated by the portion of the first term where $\Phi(z) \approx 1$. This yields the asymptotic approximation

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \tilde{p}_{\text{fix}}(\mu \rightarrow \mu_m, 0) \cdot \frac{v}{x_{cm} \Delta} \left(1 - e^{-\frac{x_{cm} \Delta}{v}}\right), \quad (|s_m| \lesssim \sqrt{v}) \quad (\text{S141})$$

On the other hand, when $O(\sqrt{v}) \lesssim \Delta \lesssim s'_b - x_{cm} - O(\sqrt{v})$, the fixation probability will be dominated by the sum of the first three terms. These contributions will be peaked around a narrow interval near $x^* = -ks'_b - \Delta$, where the sum of the first three terms transitions to the asymptotic expression in Eq. (S133):

$$\begin{aligned} \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &\approx \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), 0) \left(\frac{U'_b}{s'_b}\right)^k \frac{v}{s'_b \sqrt{2\pi v}} \int_{x_{cm} - (k+1)s'_b}^{x_{cm} - ks'_b - \zeta} \frac{e^{-\frac{(x-s_m)^2}{2v}}}{\Gamma\left(1 + \frac{|x|}{s'_b}\right)} \frac{x + (k+1)s'_b}{|x + ks'_b|} \frac{dx}{x_{cm}}, \\ &\approx \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), 0) \left(\frac{U'_b}{s'_b}\right)^k \frac{v}{s'_b x_{cm}} \frac{1}{\Gamma\left(k + 1 + \frac{\Delta}{s'_b}\right)} \frac{s'_b - \Delta}{\Delta}. \quad (\text{S142}) \end{aligned}$$

which is valid for $O(\sqrt{v}) \lesssim \Delta \lesssim s'_b - x_{cm} - O(\sqrt{v})$. Finally, in the last region where when $s'_b - x_{cm} + O(\sqrt{v}) \lesssim \Delta \lesssim s'_b - O(s'_b)$, the fixation probability will be dominated by the sum of all four terms, but the first three

will be peaked near the lower boundary. This yields

$$\begin{aligned}
\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &\approx \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), 0) \left(\frac{U'_b}{s'_b}\right)^k \left\{ \frac{v}{s'_b \sqrt{2\pi v}} \int_{x_{cm} - (k+1)s'_b}^{x_{cm} - ks'_b - \zeta} \frac{e^{-\frac{(x-s_m)^2}{2v}}}{\Gamma(1 + \frac{|x|}{s'_b})} \frac{x + (k+1)s'_b}{|x + ks'_b|} \frac{dx}{x_{cm}} \right. \\
&\quad \left. + \frac{e^{-\frac{ks'_b \Delta}{v} - \frac{\Delta^2}{2v} + \frac{(2k+1)s'^2_b - x^2_{cm}}{2v}}}{(k+1)!} \frac{x_{cm} \sqrt{v}}{s'^2_b \sqrt{2\pi}} \int_{x_{cm} - (k+2)s'_b}^{x_{cm} - (k+1)s'_b} e^{\frac{x(s'_b - \Delta)}{v}} \Phi\left(\frac{x + (k+1)s'_b}{\sqrt{v}}\right) \frac{dx}{x_{cm}} \right\} \\
&\approx p_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), 0) \left(\frac{U'_b}{s'_b}\right)^k \left\{ \frac{\frac{v}{s'_b \sqrt{2\pi v}} v e^{-\frac{(\Delta - s'_b + x_{cm})^2}{2v}}}{(\Delta - s'_b + x_{cm})(s'_b - x_{cm}) \Gamma(k+2 - \frac{x_{cm}}{s'_b})} \right. \\
&\quad \left. + \frac{\frac{v}{s'_b \sqrt{2\pi v}} v e^{-\frac{(\Delta - s'_b + x_{cm})^2}{2v}}}{s'_b (s'_b - \Delta) \Gamma(k+2)} \right\} \tag{S143}
\end{aligned}$$

which is valid for $s'_b - x_{cm} + \mathcal{O}(\sqrt{v}) \lesssim \Delta \lesssim s'_b - \mathcal{O}(s'_b)$. The leading order solution for the fixation probability can therefore be summarized as

$$\log \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \begin{cases} \frac{x_c^2}{2v} - \frac{x_{cm}^2}{2v} + \frac{x_{cm}s_m}{v} & \text{if } |s_m| \ll x_c, s_m > 0, \\ \frac{x_c^2}{2v} - \frac{x_{cm}^2}{2v} - \log\left(1 - \frac{|s_m|}{s'_b}\right) + \log\left(\frac{v}{x_{cm}|s_m|}\right) & \text{if } \sqrt{v} \ll |s_m| \ll s'_b - x_{cm}, s_m < 0, \\ \frac{x_c^2}{2v} - \frac{x_{cm}^2}{2v} - \frac{|s_m|}{s'_b} \log\left(\frac{U'_b}{|s_m|}\right) & \text{if } |s_m| \gg s'_b, s_m < 0. \end{cases} \tag{S144}$$

5 Extension to continuous distributions of fitness effects

So far, we have focused on a simplified model of the mutation spectrum, where all new mutations confer the same characteristic fitness benefit. In reality, beneficial mutations can have a range of different fitness benefits, so a general evolvability modifier will involve an arbitrary perturbation to the distribution of fitness effects,

$$\delta\mu(s) = \mu_m(s) - \mu(s), \tag{S145}$$

which could include the addition or subtraction of mutations with a range of different fitness benefits. This forces us to return to the full solution in Eq. (S44) of SI Section 3,

$$\begin{aligned}
w_m(x) &= 2x_{cm} e^{\frac{x^2}{2v} - \frac{x^2_{cm}}{2v}} \\
&\times \int_0^\infty ds \int_{-\infty}^{x+s} \frac{\mu(s) + \delta\mu(s)}{v} \left(\frac{e^{\frac{ys}{v} - \frac{s^2}{2v}} \theta(y - x_{\min})}{1 + \frac{x_{cm}}{y} e^{\frac{y^2}{2v} - \frac{x^2_{cm}}{2v}}} + \frac{w_m(y) \theta(x_{\min} - y)}{2x_{cm} e^{\frac{(y-s)^2}{2v} - \frac{x^2_{cm}}{2v}}} \right) dy, \tag{S146}
\end{aligned}$$

along with the corresponding condition for x_{cm} ,

$$1 \approx \int_0^\infty ds \int_{-\infty}^{x_{cm}+s} \frac{\mu(s) + \delta\mu(s)}{v} \left(\frac{e^{\frac{ys}{v} - \frac{s^2}{2v}} \theta(y - x_{\min})}{1 + \frac{x_{cm}}{y} e^{\frac{y^2}{2v} - \frac{x^2_{cm}}{2v}}} + \frac{w_m(y) \theta(x_{\min} - y)}{2x_{cm} e^{\frac{(y-s)^2}{2v} - \frac{x^2_{cm}}{2v}}} \right) dy, \tag{S147}$$

both of which now include an integral over the fitness benefit of the next mutation (s) in addition to its landing fitness (y). The solutions to this equation are more complicated than the single-effect model in SI Section 4, since the y values that contribute the most to each integral will depend on the corresponding value of s (and vice versa).

However, previous work (42, 43) has shown that for a broad class of distributions with exponentially bounded tails, the integrals over s are dominated by a characteristic fitness effect s'_b when $x \geq x_c - s'_b$. This is a crucial simplification, as it implies that for this range of fitnesses, the shape of $w_m(x)$ and the location of x_{cm} can be well approximated by the single- s solutions in SI Section 4 for some appropriately chosen values of s'_b and U'_b (42, 43, 51, 52). These effective parameters will depend on the underlying values of $\mu_m(s)$, x_{cm} , and v , and must be determined self-consistently using the procedures described below.

In this section, we ask whether this same idea can be extended to the full range of x values in Eq. (S146), which is necessary for predicting the modifier fixation probability, $p_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m)$. We will see that the single- s approximation holds for certain classes of modifiers, but not others, and we will develop methods for approximating $p_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m)$ in both cases.

To carry out this analysis, we note that the integral formulation in Eqs. (S146) and (S147) naturally decomposes the fixation probability into contributions from different fitness benefits. By splitting $\mu_m(s)$ into its components, $\mu_m(s) = \mu(s) + \delta\mu(s)$, it also allows us to infer the relative contributions of $\mu(s)$ and $\delta\mu(s)$. We can make this connection even more explicit by rewriting Eq. (S147) in the form,

$$1 = \frac{\sqrt{2\pi v}}{2x_{cm}v e^{-\frac{x_{cm}^2}{2v}}} \int \mu(s) p_{\text{fix}}(s|x_{cm}, v) ds + \frac{\sqrt{2\pi v}}{2x_{cm}v e^{-\frac{x_{cm}^2}{2v}}} \int \delta\mu(s) p_{\text{fix}}(s|x_{cm}, v) ds, \quad (\text{S148})$$

where

$$p_{\text{fix}}(s|x_{cm}, v) = \int_{-\infty}^{x_{cm}} \frac{e^{-\frac{x^2}{2v}}}{\sqrt{2\pi v}} w_m(x+s) dx \quad (\text{S149})$$

is the fixation probability of a first-order mutation in the “dual” population defined by $v(\mu_m(s), N^*) = v(\mu(s), N)$ (see Eq. S38 in SI Section 3).

The decomposition in Eq. (S148) allows us to distinguish between two broad regimes depending on the relative contributions of $\mu(s)$ and $\delta\mu(s)$. When the fixation probability of the modifier is dominated by mutations from $\mu(s)$, then s'_b and x_{cm} will remain close to their wildtype values. This allows us to calculate x_{cm} and $p_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m)$ perturbatively, by treating the $\delta\mu(s)$ term as a small correction. We will refer to this case as the *perturbative regime*, since it will only apply for modifiers with weak-to-moderate fixation probabilities. For stronger second-order selection pressures, the fixation probability of the modifier will be dominated by mutations from the $\delta\mu(s)$ portion of $\mu_m(s)$. In this *modifier-dominated regime*, we will need to turn to the full solution of $w_m(x)$, with $\mu_m(s) \approx \delta\mu(s)$. We will discuss each of these cases in turn, after quickly reviewing the wildtype dynamics when $\delta\mu(s) = 0$.

5.1 Location of the interference threshold for the wildtype population

When $\delta\mu(s) = 0$, Eq. (S147) determines the location of interference threshold for the wildtype population,

$$1 = \int_0^\infty ds \int_{-\infty}^{x_c+s} \frac{\mu(s)}{v} \left(\frac{e^{\frac{ys}{v}}}{1 + \frac{x_c}{y} e^{\frac{y^2}{2v} - \frac{x_c^2}{2v}}} + \frac{w_m(y)\theta(x_{\min} - y)}{2x_c e^{\frac{(y-s)^2}{2v} - \frac{x_c^2}{2v}}} \right) dy. \quad (\text{S150})$$

The solution to this equation was previously described by Ref. (42). We reproduce it here for completeness, since we will build on this result in the following sections.

Ref. (42) showed that in our parameter regime of interest, the s integral in Eq. (S150) will be strongly peaked at a characteristic value, s_b , which will be determined self-consistently below. This observation allows us to evaluate the inner integral over y for s close to s_b . When $x_c - s_b \gg \sqrt{v}$, the contributions from $y < x_c$ will be dominated by y within $O(v/s)$ of x_c , while the contributions from $y > x_c$ will be dominated by fitnesses within $O(v/(x_c - s))$ of x_c . Equation (S150) then reduces to a single integral over s ,

$$1 \approx \int_0^\infty \mu(s) e^{\frac{x_c s}{v} - \frac{s^2}{2v}} \left[\frac{1}{s} + \frac{1}{x_c - s} \right] ds, \quad (\text{S151})$$

where the rapidly increasing exponential term is eventually counteracted by the rapidly decreasing $\mu(s)$ term. When $x_c - s_b \gg \sqrt{v}$, the competition between these two factors produces a sharp peak around $s \approx s_b \pm \Delta s_b$, where s_b is the value of s that maximizes the rapidly varying portion of the integrand,

$$x_c - s_b + v \frac{\partial \log \mu(s_b)}{\partial s} = 0, \quad (\text{S152})$$

and Δs_b is the characteristic width,

$$\Delta s_b = \left[\frac{1}{v} - \frac{\partial^2 \log \mu(s_b)}{\partial s^2} \right]^{-1/2}. \quad (\text{S153})$$

When $s_b \gg \Delta s_b$, we can evaluate Eq. (S151) with a Gaussian Laplace approximation to obtain

$$1 = \frac{U_b}{s_b} \left[1 - \frac{s_b}{x_c} \right]^{-1} e^{\frac{x_c s_b}{v} - \frac{s_b^2}{2v}}, \quad (\text{S154})$$

where the effective mutation rate,

$$U_b = \sqrt{2\pi \Delta s_b} \mu(s_b), \quad (\text{S155})$$

denotes the total rate of producing ‘‘driver’’ mutations within Δs_b of s_b . This matches the corresponding expression for the single-effect model in SI Section 4 if the effective selection strength and mutation rate are defined by Eqs. (S152) and (S155). Unlike SI Section 4, these effective parameters will now depend on the underlying values of $\mu(s)$ and N , and can shift if one or the other is varied (42, 43, 51, 52). This behavior will be crucial for understanding the fates of modifier mutations below.

Application to stretched exponential distributions. Following previous work (42, 43, 51), we can gain some intuition for these expressions by considering the class of stretched exponential distributions,

$$\mu(s) = \frac{U_0}{s_0} \cdot \frac{e^{-(s/s_0)^\beta}}{\Gamma(1 + \beta^{-1})}, \quad (\text{S156})$$

as a function of the tail parameter $\beta \geq 1$. Ref. (42) showed that the effective parameters for the exponential case ($\beta = 1$) are given by

$$s_b = x_c - \frac{v}{s_0} \approx x_c, \quad \Delta s_b = \sqrt{v}, \quad (\text{S157a})$$

leading to the approximate scaling,

$$x_c \sim s_0 \log(NU_0), \quad v \sim \frac{s_0^2 \log^2(NU_0)}{2 \log(s_0/U_0)}. \quad (\text{S157b})$$

Similarly, the solution for the $\beta \gg 1$ case is given by

$$s_b \approx \left(\frac{x_c s_0^\beta}{\beta v} \right)^{\frac{1}{\beta-1}} \left[1 - \frac{1}{1 + \beta(\beta-1) \frac{v}{s_0^2} \left(\frac{x_c}{\beta v} \right)^{\frac{\beta-2}{\beta-1}}} \right] \approx \left(\frac{x_c s_0^\beta}{\beta v} \right)^{\frac{1}{\beta-1}}, \quad (\text{S158a})$$

$$\Delta s_b = \left[\frac{1}{v} + \frac{\beta(\beta-1)}{s_0^2} \left(\frac{s'_b}{s_0} \right)^{\beta-2} \right]^{-1/2} \approx \frac{s_0}{\sqrt{\beta(\beta-1) \left(\frac{x_c s_0}{\beta v} \right)^{\frac{\beta-2}{\beta-1}}}}, \quad (\text{S158b})$$

with

$$x_c \sim \frac{2s_0 \log(Ns_0)}{\log(\sqrt{\beta}s_0/U_0)}, \quad v \sim \frac{2s_0^2 \log(Ns_0)}{\log^2(\sqrt{\beta}s_0/U_0)}. \quad (\text{S158c})$$

In our analysis below, we will find that the $\beta = 2$ case (corresponding to a half-Gaussian distribution) will serve as an important boundary case, so it will be useful to derive an explicit expression for this case as well. Equations (S152) and (S153) imply that

$$s_b = \frac{x_c}{1 + 2 \left(\frac{v}{s_0^2} \right)}, \quad \Delta s_b = \sqrt{v \left[\frac{1}{1 + 2 \left(\frac{v}{s_0^2} \right)} \right]}, \quad (\text{S159a})$$

which yields the asymptotic approximations

$$v \sim \frac{s_0^2 \log(NU_0)}{2 \log(s_0/U_0)}, \quad x_c \sim \sqrt{\frac{\log(Ns_0) \log(NU_0)}{\log(s_0/U_0)}}. \quad (\text{S159b})$$

When $v/s_0^2 \ll 1$, this solution behaves like the exponential case above (with $s_b \approx x_c$). This regime will be valid provided that

$$1 \ll \log(NU_0) \ll \log(s_0/U_0). \quad (\text{S160})$$

On the other hand, when $v/s_0^2 \gg 1$, the solution will behave more like the $\beta \gg 1$ case ($s_b \ll x_c$). This regime will be valid provided that

$$1 \ll \log(s_0/U_0) \ll \log(Ns_b) \ll \log^2(s_0/U_0). \quad (\text{S161})$$

which mirrors the conditions of validity for the single-effect model in SI Section 4.

5.2 Perturbative regime

5.2.1 Location of the interference threshold

In the perturbative regime, we assume that the dominant contribution to Eq. (S147) still comes from the $\mu(s)$ term. This suggests that x_{cm} will remain close to the wildtype interference threshold,

$$x_{cm} = x_c + \delta x_c, \quad (\text{S162})$$

where δx_c is a small correction whose magnitude will be determined self-consistently below.

If $\delta x_c \ll v/s_b$, then the integral over s will remain strongly peaked around the same value as the wildtype integral in Eq. (S151). This implies that

$$s_b(N^*) \approx s_b(N), \quad U_b(N^*) = U_b(N), \quad (\text{S163})$$

where s_b and U_b are defined as in Eqs. (S152) and (S155), and that

$$\int \mu(s) e^{\frac{(x_c + \delta x_c)s}{v} - \frac{s^2}{2v}} ds \approx U_b \tilde{p}_{\text{fix}}(s_b | x_c) \left(1 + \frac{s_b \delta x_c}{v} \right) = U_b \tilde{p}_{\text{fix}}(s_b | x_c + \delta x_c). \quad (\text{S164})$$

From the self-consistency condition for N in Eq. (S72), we also have

$$\frac{2N^*(x_c + \delta x_c) e^{-\frac{(x_c + \delta x_c)^2}{2v}}}{\sqrt{2\pi v}} = \frac{1}{s_b} = \frac{2Nx_c e^{-\frac{x_c^2}{2v}}}{\sqrt{2\pi v}}, \quad (\text{S165})$$

and

$$\frac{N^*}{N} \approx e^{\frac{x_c \delta x_c}{v}}. \quad (\text{S166})$$

Substituting these results into Eq. (S148), we obtain

$$1 = \frac{s_b}{v} \cdot U_b \tilde{p}_{\text{fix}}(s | x_c) \left(1 + \frac{s_b \delta x_c}{v} \right) [1 + I(\delta\mu, \delta x_c)], \quad (\text{S167})$$

where

$$I(\delta\mu, \delta x_c) \equiv \int \frac{\delta\mu(s)}{U_b} e^{\log \tilde{p}_{\text{fix}}(s | x_c + \delta x_c) - \log \tilde{p}_{\text{fix}}(s_b | x_c + \delta x_c)} ds \quad (\text{S168})$$

is the integral from Eq. (8) in the main text. Equation (S167) shows that $I(\delta\mu, \delta x_c)$ can be interpreted as the relative contribution of mutations from $\delta\mu(s)$ vs $\mu(s)$ when $\delta x_c \ll v/s_b$. Substituting for the wildtype x_c condition [$U_b \tilde{p}_{\text{fix}}(s_b) = v/s_b$], we can rearrange Eq. (S167) to obtain

$$\delta x_c = -\frac{v}{s_b} I(\delta\mu, \delta x_c), \quad (\text{S169})$$

which is valid when $\delta x_c \ll v/s_b$. This shows that $\delta x_c \ll v/s_b$ is equivalent to the condition that $I(\delta\mu, \delta x_c) \ll 1$. If $\delta x_c \ll v/s$ for the fitness effects that dominate Eq. (S168), then we can perform the same Taylor expansion obtain an analytical expression for δx_c ,

$$\delta x_c = -\frac{v}{s_b} I(\delta\mu, 0) = -\frac{v}{s_b} \int \frac{\delta\mu(s)}{U_b} e^{\log \tilde{p}_{\text{fix}}(s | x_c) - \log \tilde{p}_{\text{fix}}(s_b | x_c)} ds, \quad (\text{S170})$$

which can be evaluated in terms of the wildtype interference threshold x_c . For stronger fitness benefits, δx_c must be obtained by solving the implicit relation in Eq. (S169).

5.2.2 Fixation probabilities of modifiers

When the above assumptions hold, the fixation probability of the modifier can be obtained from our previous calculations in Eqs. (S83) and (S93). Since $\delta x_c \ll v/s_b \ll \sqrt{v} \ll x_c$, the fixation probability of a modifier without direct costs or benefits will be dominated by exceptionally fit genetic backgrounds ($x_c - s_b \lesssim x \lesssim x_c$), so that Eq. (S83) reduces to,

$$\tilde{p}_{\text{fix}}(\delta\mu) \approx \frac{x_{cm}}{x_c} e^{\frac{x_c^2}{2v} - \frac{x_{cm}^2}{2v}} \approx \exp\left[-\frac{x_c \delta x_c}{v}\right]. \quad (\text{S171})$$

Substituting our expression for δx_c in Eq. (S169) then yields

$$\tilde{p}_{\text{fix}}(\delta\mu) \approx \exp\left[\frac{x_c}{s_b} \cdot I(\delta\mu, \delta x_c)\right]. \quad (\text{S172})$$

This shows that the conditions of validity for the perturbative regime,

$$\delta x_c \ll v/s_b \iff I(\delta\mu, \delta x_c) \ll 1, \quad (\text{S173})$$

are equivalent to the assumption that

$$|\log \tilde{p}_{\text{fix}}(\delta\mu)| \ll x_c/s_b. \quad (\text{S174})$$

Since x_c can be asymptotically larger than s_b , this implies that a small change to the distribution of fitness effects can be strongly selected ($|\log \tilde{p}_{\text{fix}}(\delta\mu)| \gg 1$) even when it leads to a negligible change in s_b and U_b . The perturbative regime is therefore an example of a modifier mutation that cannot be reduced to the single-effect model in SI Section 4.

It is straightforward to extend this calculation to consider modifiers that also include direct costs or benefits. Since the second-order selection pressures are bounded by Eq. (S174), most of the interesting behavior will occur when $|s_m| \sim \mathcal{O}(x_c/s_b) \ll s_b$. The fixation probability of these modifiers will continue to be dominated by exceptionally fit genetic backgrounds ($x_c - s_b \lesssim x \lesssim x_c$), so that Eq. (S93) reduces to

$$\log \tilde{p}_{\text{fix}}(\delta\mu, s_m) \approx \log \tilde{p}_{\text{fix}}(\delta\mu) + \log \tilde{p}_{\text{fix}}(s_m), \quad (\text{S175})$$

where $\tilde{p}_{\text{fix}}(s_m)$ is the fixation probability of a first-order mutation from SI Section 4.1.3.

5.3 Modifier-dominated regime (multiple mutations)

For stronger second-order selection pressures ($|\log \tilde{p}_{\text{fix}}(\delta\mu)| \gtrsim x_c/s_b$), the fixation probability in Eq. (S146) will start to depend more sensitively on the mutations that are added (or removed) by $\delta\mu(s)$. For a strongly beneficial modifier, the fixation probability will be dominated by the new mutations that are added by $\delta\mu(s)$. In contrast, a strongly deleterious modifier will be dominated by the mutations that remain once $\delta\mu(s)$ is subtracted from $\mu(s)$. In both cases, we will need to develop a general solution for $w_m(x)$ for arbitrary distributions $\mu_m(s) \neq \mu(s)$.

Following our analysis in SI Section 4, it will be useful to distinguish between two characteristic regimes, depending on whether the dominant fitness benefits are smaller or larger than x_{cm} . We will continue to refer to these as the *multiple mutations* and *quasi-sweeps* regimes, respectively, since they will qualitatively resemble their analogues in SI Section 4. We will start by analyzing the multiple mutations regime, while the quasi-sweeps case will be considered in SI Section 5.4 below.

5.3.1 Location of the interference threshold

When the relevant fitness benefits in $\mu_m(s)$ are small compared to x_{cm} , we can repeat the calculation in SI Section 5.1 to obtain an analogous condition for x_{cm} ,

$$1 \approx \int_0^\infty \mu_m(s) e^{\frac{x_{cm}s}{v} - \frac{s^2}{2v}} \left[\frac{1}{s} + \frac{1}{x_{cm} - s} \right] ds. \quad (\text{S176})$$

This integral will be peaked around a new value s'_b that satisfies

$$x_{cm} - s'_b + v \frac{\partial \log \mu_m(s'_b)}{\partial s} = 0, \quad (\text{S177})$$

along with a new characteristic width,

$$\Delta s'_b = \left[\frac{1}{v} - \frac{\partial^2 \log \mu_m(s'_b)}{\partial s^2} \right]^{-1/2}, \quad (\text{S178})$$

which will generally differ from the wildtype values of s_b and Δs_b above. When $s'_b \gg \Delta s'_b$ an analogous Gaussian Laplace approximation yields a simple expression for x_{cm} ,

$$1 = \frac{U'_b}{s'_b} \left[1 - \frac{s'_b}{x_{cm}} \right]^{-1} e^{\frac{x_{cm}s'_b}{v} - \frac{s'^2_b}{2v}}, \quad (\text{S179})$$

where we have defined the effective mutation rate

$$U'_b = \sqrt{2\pi \Delta s'^2_b \mu_m(s'_b)}. \quad (\text{S180})$$

Together, Eqs. (S177-S179) allow us to solve for s'_b , U'_b , and x_{cm} as a function of $\mu_m(s)$ and v . These expressions will be self-consistently valid when $x_{cm} - s'_b \gg \sqrt{v}$.

5.3.2 Extending the shoulder solution to lower fitness values.

With the location of x_{cm} fixed by Eq. (S179), we can repeat our calculation in SI Section 4.1.2 to extend the shoulder solution for $w_m(x)$ to progressively lower fitness values. Provided that $x + s'_b \gtrsim x_{cm} + \mathcal{O}(\Delta s'_b)$, the dominant contribution to Eq. (S146) will also be contained within the region of integration. This demonstrates that the shoulder solution will continue to be valid for fitnesses as low as

$$x_{\min}(\mu_m(s)) = x_{cm} - s'_b + \Delta s'_b. \quad (\text{S181})$$

In particular, it shows that for $x > x_{\min}(\mu_m(s))$, the shape of $w_m(x)$ can be approximated by a single-effect model with s'_b and U'_b defined by Eqs. (S177) and (S180) above.

Hopping vs leapfrogging. When $x \lesssim x_{\min}(\mu_m(s))$, the upper limit in the landing fitness integral in Eq. (S146) will again become important. However, the crucial difference from our previous analysis in SI Section 4.1 is that this upper limit is not a fixed parameter, but can vary as a function of s . In other words, it is theoretically possible for the mutant offspring to land at *any* relative fitness value if it acquires a sufficiently large beneficial mutation from $\mu_m(s)$. This is a new effect that is only present when we allow for a continuous distribution of fitness effects. This will lead to two distinct classes of solutions for $w_m(x)$ depending on the tail of $\mu_m(s)$.

If the tail of $\mu_m(s)$ decays sufficiently rapidly, then a modifier lineage that arises below $x_{cm} - s'_b$ will typically fix by “hopping” to higher relative fitnesses through a sequence of smaller mutations, similar to the single-effect models in SI Section 4. However, if the tail of $\mu_m(s)$ decays more slowly, then a modifier arising below $x_{cm} - s'_b$ will typically fix by generating a *single* large driver mutation that bypasses these intermediate fitness values, and lands at or above interference threshold at x_{cm} . This “leapfrogging” regime is potentially empirically relevant, since it will apply when $\mu_m(s)$ decays like an exponential distribution.

By splitting the integral in Eq. (S146) into contributions from above and below the interference threshold, we can obtain a corresponding expression for the fixation probability,

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \underbrace{\int_{x_{cm}-x}^{\infty} \frac{\mu_m(s) 2(x+s) ds}{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} |x|}}_{\text{landing above } x_{cm}} + 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \underbrace{\int_{x_{\min}-x}^{x_{cm}-x} \mu_m(s) e^{\frac{xs}{v} + \frac{s^2}{2v}} \frac{ds}{s}}_{\text{landing below } x_{cm}}, \quad (\text{S182})$$

which will be valid for $x \gtrsim x_{cm} - 2s'_b + \mathcal{O}(v/s'_b + \Delta s_b)$. It will be helpful to evaluate these expressions by defining a new variable,

$$\Delta x \equiv x_{cm} - s'_b - x, \quad (\text{S183})$$

such that $\Delta x \in (0, s'_b)$.

When $\Delta x \gtrsim \mathcal{O}(\Delta s_b)$, then the lower limit of the integral in the first term in Eq. (S182) will lie in the high fitness tail of $\mu_m(s)$. An exponential Laplace approximation around $s = x_{cm} - x = s'_b + \Delta x$ then yields

$$2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \underbrace{\int_{x_{cm}-x}^{\infty} \frac{\mu_m(s) 2(x+s) ds}{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} |x|}}_{\text{landing above } x_{cm}} \approx \frac{2x_{cm}}{|x|} \frac{\mu_m(s'_b + \Delta x)}{-\partial_s \log \mu_m(s'_b + \Delta x)}. \quad (\text{S184})$$

The integrand in the second term in Eq. (S182) has a critical point at $s = s^*(\Delta x)$, which satisfies,

$$x_{cm} - s'_b - \Delta x + s^* + v \frac{\partial \log \mu_m(s^*)}{\partial s} = 0. \quad (\text{S185})$$

This critical point will correspond to a local maximum if the second derivative is negative,

$$1 + v \frac{\partial^2 \log \mu_m(s^*)}{\partial s^2} < 0. \quad (\text{S186})$$

The value of this integral will therefore crucially depend on the tail of $\mu_m(s)$.

If the curvature of $\mu_m(s)$ is sufficiently small [$\partial_s^2 \log \mu_m(s^*) > -1/v$], then s^* will correspond to a local minimum of the integrand, and the integral will therefore be dominated by the upper limit of integration. An exponential Laplace approximation then yields,

$$w_m(x) \approx \frac{2x_{cm}}{|x|} \cdot \frac{\mu_m(s'_b + \Delta x)}{-\partial_s \log \mu_m(s'_b + \Delta x)} + \frac{2x_{cm}}{s'_b + \Delta x} \cdot \frac{\mu_m(s'_b + \Delta x)}{\frac{x_{cm}}{v} + \partial_s \log \mu_m(s'_b + \Delta x)}. \quad (\text{S187})$$

Intuitively speaking, this says that the fixation probability is dominated by mutations of size $s \approx s'_b + \Delta x$ that result in a landing fitness very close to the interference threshold x_{cm} . Since Δx can be as large as s'_b (which is $\gg \Delta s'_b$), this constitutes an example of the “leapfrogging” behavior above, since it will involve a jump much larger than s'_b . From the form of the condition in Eq. (S186), we can see that this leapfrogging behavior

will always occur for the exponential distribution in Eq. (S156), as well as the Gaussian case ($\beta = 2$) when $v/s_0^2 \ll 1$. In these cases, the single- s equivalence will break down for initial fitness below $x_{cm} - s'_b$, and we must use Eq. (S187) instead.

On the other hand, if the curvature of $\mu_m(s)$ is large and negative [$-\partial_s^2 \log \mu_m(s^*) \gg 1/v$], then s^* will be a local maximum, and a Gaussian Laplace approximation yields

$$\int_{x_{\min}-x}^{x_{cm}-x} \mu_m(s) e^{\frac{xs}{v} + \frac{s^2}{2v}} \frac{ds}{s} \approx \sqrt{\frac{2\pi}{-\partial_s^2 \log \mu_m(s^*)}} \frac{\mu_m(s^*)}{s^*} e^{\frac{xs^*}{v} - \frac{s^{*2}}{2v}}. \quad (\text{S188})$$

When $-\partial_s^2 \log \mu_m(s^*) \gg 1/v$, we can also solve Eq. (S185) perturbatively to show that $s^*(\Delta x)$ is close to s'_b :

$$s^* \approx s'_b \left(1 + \frac{1 - \frac{\Delta x}{s'_b}}{-v \partial_s^2 \log \mu_m(s'_b)} \right) \approx s'_b, \quad (\text{S189})$$

so that Eq. (S182) reduces to

$$w_m(x) \approx \underbrace{\frac{2x_{cm}}{|x|} \cdot \frac{\mu_m(x_{cm} - x)}{-\partial_s \log \mu_m(s'_b + \Delta x)}}_{\text{landing above } x_{cm} \text{ "leapfrogging"}} + \underbrace{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v} + \frac{xs'_b}{v} + \frac{s'^2_b}{2v}} \frac{\sqrt{2\pi \Delta s'^2_b} \mu_m(s'_b)}{s'_b}}_{\text{landing below } x_{cm} \text{ "hopping"}}. \quad (\text{S190})$$

The balance between these two terms will determine whether the typical successful lineage “hops” or “leapfrogs” to fixation. Comparing the magnitudes of these terms, we find that the hopping term in Eq. (S190) will dominate if

$$\log \left[\frac{\mu_m(s'_b)}{\mu_m(s'_b + \Delta x)} \right] - \frac{x_{cm} \Delta x}{v} + \frac{\Delta x^2}{2v} \gtrsim \mathcal{O}(1). \quad (\text{S191})$$

This condition is more stringent than the one defining s'_b in Eq. (S176), which has an additional $s'_b \Delta x/v$ term on the left hand side. This means that even if mutations above s'_b are negligible for $x \gtrsim x_{cm} - s'_b$, they may become relevant for larger direct costs. When $\beta \geq 2$, the decay of $\mu_m(s)$ will still be bounded by its Gaussian approximation,

$$-\log \mu_m(s'_b + \Delta x) \geq -\log \mu_m(s'_b) - \frac{(s'_b - x_{cm}) \Delta x}{v} + \frac{\Delta x^2}{2} \left(\frac{1}{\Delta s'^2_b} - \frac{1}{v} \right), \quad (\text{S192})$$

so Eq. (S191) reduces to

$$-\frac{s'_b \Delta x}{v} + \frac{\Delta x^2}{2 \Delta s'^2_b} \gtrsim \mathcal{O}(1). \quad (\text{S193})$$

The leapfrogging term will therefore be negligible for

$$\Delta x \gtrsim \mathcal{O}(s'_b \cdot \Delta s'^2_b / v), \quad (\text{S194})$$

which will constitute the majority of the interval $(0, s'_b)$ when $v \gg \Delta s'^2_b$ — this is the same condition used to evaluate the hopping term above. For more general choices of $\mu_m(s)$, the Gaussian approximation may not

longer bound the value of $\mu_m(s'_b + \Delta x)$, and full condition in Eq. (S191) must be used instead. When these conditions are met, the fixation probability in Eq. (S190) will be well-approximated by the hopping term,

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v} - \frac{s'_b(x_{cm} - s'_b - x)}{v}} \left(\frac{x_{cm} - s'_b}{x_{cm}} \right), \quad (\text{S195})$$

where we have used the auxilliary condition in Eq. (S179) to eliminate the explicit dependence on $\mu_m(s'_b)$. This expression is identical to the result in SI Section 4.1.2, which shows that the single- s mapping continues to apply for $x \gtrsim x_{cm} - 2s'_b$.

We can continue these calculations to extend $w_m(x)$ to lower fitness values. At the next rung $[x_{cm} - 3s'_b \lesssim x \lesssim x_{cm} - 2s'_b]$, we can once again write x as

$$x = x_{cm} - 2s'_b - \Delta x, \quad (\text{S196})$$

for some $\Delta x \in (0, s'_b)$, so that Eq. (S146) reduces to

$$\begin{aligned} w_m(x) \approx & 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \underbrace{\int_{2s'_b + \Delta x}^{\infty} \frac{\mu_m(s) 2(x+s)}{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} |x|} ds}_{\text{landing above } x_{cm}} \\ & + 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \underbrace{\int_{s'_b + \Delta x}^{2s'_b + \Delta x} \mu_m(s) e^{\frac{(x_{cm} - 2s'_b - \Delta x)s}{v} + \frac{s^2}{2v}} \frac{ds}{s}}_{\text{landing between } x_{cm} - s'_b \text{ and } x_{cm}} \\ & + 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v} - \frac{s'_b(x_{cm} - s'_b - x)}{v}} \underbrace{\int_{\Delta x}^{s'_b + \Delta x} \mu_m(s) e^{\frac{s^2}{2v} + \frac{s(x_{cm} - s'_b - \Delta x)}{v}} \frac{ds}{s + s'_b}}_{\text{landing below } x_{cm} - s'_b}. \end{aligned} \quad (\text{S197})$$

The first term in Eq. (S197) can be evaluated with a similar Laplace approximation as Eq. (S184), yielding

$$2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} \underbrace{\int_{2s'_b + \Delta x}^{\infty} \frac{\mu_m(s) 2(x+s)}{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v}} |x|} ds}_{\text{landing above } x_{cm}} \approx \frac{2x_{cm}}{|x|} \frac{\mu_m(2s'_b + \Delta x)}{-\partial_s \log \mu_m(2s'_b + \Delta x)}. \quad (\text{S198})$$

The integral in the third term in Eq. (S197) is similar to the hopping term in Eq. (S182) (up to a prefactor). This implies that it will have a critical point at $s = s^*(\Delta x)$ with the same location and curvature as Eq. (S185). Finally, the second term in Eq. (S197) corresponds to intermediate leapfrogging events of size $s \in (s'_b, 2s'_b)$. The integrand in this term will have a critical point at a new value, s_2^* , defined by

$$x_{cm} - 2s'_b - \Delta x + s_2^* + v \frac{\partial \log \mu_m(s_2^*)}{\partial s} = 0, \quad (\text{S199})$$

which will generally be less than s^* . If the curvature of $\mu_m(s)$ is small $[-\partial_s^2 \log \mu_m(s'_b) < 1/v]$, then the critical points in the second and third terms will once again correspond to local minima, and the fixation probability will be dominated by the first term,

$$w_m(x) \approx \frac{2x_{cm}}{|x|} \frac{\mu_m(2s'_b + \Delta x)}{-\partial_s \log \mu_m(2s'_b + \Delta x)}, \quad (\text{S200})$$

similar to Eq. (S187) above. On the other hand, if $-\partial_s^2 \log \mu_m(s'_b) \gg 1/v$, then $s^*(\Delta x) \approx s_2^*(\Delta x) \approx s'_b$ will both correspond to local maxima, and the fixation probability will be dominated by the first and third terms:

$$w_m(x) \approx \underbrace{\frac{2x_{cm}}{|x|} \frac{\mu_m(2s'_b + \Delta x)}{-\partial_s \log \mu_m(2s'_b + \Delta x)}}_{\text{leapfrogging}} + \underbrace{2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v} - \frac{s'_b(x_{cm} - \frac{5}{2}s'_b - 2x)}{v}} \left(\frac{x_{cm} - s'_b}{x_{cm}} \right) \frac{\mu_m(s'_b) \sqrt{2\pi \Delta s'_b}}{2s'_b}}_{\text{hopping}}. \quad (\text{S201})$$

By comparing the magnitudes of these terms, we can see that the hopping term will dominate if

$$\log \left[\frac{\mu_m(s'_b)}{\mu_m(2s'_b + \Delta x)} \right] - \left[\frac{x_{cm}(s'_b + \Delta x)}{v} - \frac{s'^2_b + \Delta x^2}{2v} \right] \gtrsim O(1), \quad (\text{S202})$$

and the corresponding fixation probability will be given by

$$w_m(x) \approx \frac{1}{2} \cdot 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v} - \frac{(2x_{cm} - 3s'_b - 2x)s'_b}{v}} \left(\frac{x_{cm} - s'_b}{x_{cm}} \right)^2, \quad (\text{S203})$$

which matches the single-effect result from SI Section 4.1. For the class of stretched exponential distributions in Eq. (S156) with $\beta \geq 2$, the mutation spectrum term will again be bounded by its Gaussian approximation,

$$-\log \mu_m(2s'_b + \Delta x) \geq -\log \mu_m(s'_b) - \frac{(s'_b - x_{cm})(s'_b + \Delta x)}{v} + \frac{(s'_b + \Delta x)^2}{2} \left(\frac{1}{\Delta s'^2_b} - \frac{1}{v} \right), \quad (\text{S204})$$

so that Eq. (S202) reduces to

$$-\frac{s'_b(s'_b + \Delta x)}{v} + \frac{s'^2_b + \Delta x^2}{2\Delta s'^2_b} + \frac{s\Delta x}{\Delta s'^2_b} \left(1 - \frac{\Delta s'^2_b}{v} \right) \gtrsim O(1). \quad (\text{S205})$$

This shows that the hopping approximation will be valid whenever $v \gg \Delta s'^2_b$, which is the same condition that we obtained for the previous interval above. One can continue this argument to show that for relative fitnesses of the form,

$$x = x_{cm} - ns'_b - \Delta x, \quad (\text{S206})$$

with $\Delta x \in (0, s'_b)$, the hopping term will be given by

$$w_m(x) \approx \frac{1}{n!} \cdot 2x_{cm} e^{\frac{x^2 - x_{cm}^2}{2v} - \frac{ns'_b}{v} (x_{cm} - \frac{(n+1)}{2}s'_b - x)} \left(\frac{x_{cm} - s'_b}{x_{cm}} \right)^n, \quad (\text{S207})$$

and will constitute the dominant contribution whenever

$$\log \left[\frac{\mu_m(s'_b)}{\mu_m(ns'_b + \Delta x)} \right] - \left[\frac{x_{cm}(n-1)s'_b + x_{cm}\Delta x}{v} - \frac{(n-1)^2 s'^2_b + \Delta x^2}{2v} \right] \gtrsim O(1). \quad (\text{S208})$$

This shows that for the stretched exponential distributions in Eq. (S156) with $\beta \geq 2$, the single- s approximation will continue to apply for arbitrarily large fitness costs as long as $v \gg \Delta s'^2_b$.

5.3.3 Fixation probabilities of modifiers

We can now use our solution for $w_m(x)$ to evaluate the fixation probability of a modifier mutation. When the dynamics are dominated by hopping, we have seen that the functional form of $w_m(x)$ is well-approximated by the single- s model in SI Section 4.1. This implies that the fixation probability of the modifier can be predicted using our earlier expressions in SI Sections 4.1.4 and 4.1.5

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) = \tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow U'_b, s'_b), s_m), \quad (\text{S209})$$

if (U_b, s_b) and (U'_b, s'_b) are chosen to coincide with the effective parameters defined in SI Sections 5.1 and 5.3.1 above. This approximation will be valid when $\log \tilde{p}_{\text{fix}}(\mu_m) \gtrsim x_c/s_b$ and $x_{cm} - s'_b \gtrsim \sqrt{v}$. We used this expression to obtain the theoretical predictions for the exponential distribution in Fig. 4C,D in the main text.

This single- s approximation will also be valid for $\mu_m(s)$ that allow for leapfrogging, as long as the direct or indirect costs of the modifier are not too strong ($s_m + x_c - x_{cm} \gtrsim -s'_b$). For larger costs ($s_m + x_c - x_{cm} \lesssim -s'_b$), we will need to modify Eq. (S209) to include the additional contributions from the leapfrogging term. Substituting the solution for $w_m(x)$ in Eq. (S187) into Eq. (S93), we obtain

$$\begin{aligned} \tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) &\approx \frac{x_{cm}}{x_c s_b} e^{\frac{x_c^2}{2v}} \int_{-\infty}^{x_c + s_m} \frac{\mu_m(x_{cm} - x)}{-\partial_s \log \mu_m(x_{cm} - x)} e^{-\frac{(x-s_m)^2}{2v}} \frac{dx}{|x|}, \\ &\approx \frac{(x_{cm} - s'_b) s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \int_{x_{cm} - (x_c + s_m)}^{\infty} \frac{\mu_m(u)}{\mu_m(s'_b)} \cdot \frac{e^{\frac{(x_{cm} - s'_b)^2}{2v} - \frac{(u - (x_{cm} - s_m))^2}{2v}}}{-|x_{cm} - u| \cdot \partial_s \log \mu_m(u)} \frac{du}{\sqrt{2\pi \Delta s_b'^2}}, \end{aligned} \quad (\text{S210})$$

where we have substituted the auxiliary condition for x_{cm} . This integral will have a local maximum at u^* , where u^* satisfies

$$x_{cm} - s_m - u^* + v \partial_s \log \mu_m(u^*) = 0, \quad (\text{S211})$$

and a characteristic width Δ_u ,

$$\Delta_u = \left[\frac{1}{v} - \frac{\partial^2 \log \mu_m(u^*)}{\partial s^2} \right]^{-1/2}. \quad (\text{S212})$$

We can therefore evaluate it with a Gaussian Laplace approximation to obtain,

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{(x_{cm} - s'_b) s'_b}{x_c s_b} e^{\frac{x_c^2 - x_{cm}^2}{2v}} \frac{\mu_m(u^*)}{\mu_m(s'_b)} \cdot \frac{e^{\frac{(x_{cm} - s'_b)^2}{2v} - \frac{(u^* - (x_{cm} - s_m))^2}{2v}}}{-|x_{cm} - u^*| \cdot \partial_s \log \mu_m(u^*)} \frac{\sqrt{2\pi \Delta_u^2}}{\sqrt{2\pi \Delta s_b'^2}}. \quad (\text{S213})$$

In the case of an exponential distribution, $u^* = s'_b - s_m$ and $\Delta_u = \sqrt{v}$, so Eq. (S213) reduces to,

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{v s'_b}{x_c s_b |s_m|} e^{\frac{x_c^2 - x_{cm}^2}{2v} + \frac{s_m}{s'_0}}. \quad (\text{S214})$$

For a half-Gaussian distribution, $u^* = s'_b (1 - s_m/x_{cm})$ and $\Delta_u = \sqrt{\frac{1}{v} + \frac{2}{s_0^2}}$, so that

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \frac{v x_{cm}^2}{x_c s_b s'_b s_m^2} e^{\frac{x_c^2 - x_{cm}^2}{2v} + \frac{2 s_m x_{cm} - \delta^2}{s_0'^2 + 2v}} \approx \frac{v x_{cm}^2}{x_c s_b s'_b s_m^2} e^{\frac{x_c^2 - x_{cm}^2}{2v} - \frac{s_m}{s_0'^2}} \quad (\text{S215})$$

where the last approximation follows when $s'_0 \gg \sqrt{v}$. These results imply that the leading order scaling for large direct costs when leapfrogging is valid is given by

$$\log \tilde{p}_{\text{fix}} \approx \frac{x_c^2}{2v} - \frac{x_{cm}^2}{2v} + \log(\mu_m(s_m)s'_0/U'_b). \quad (\text{S216})$$

5.4 Modifier-dominated regime (quasi-sweeps)

The results in the previous section apply when $x_{cm} - s'_b \gg \sqrt{v}$. As the gap between x_{cm} and s'_b shrinks, we will eventually reach a point where a single established mutation will be sufficient to drive the modifier lineage to fixation. To solve for the shape of $w_m(x)$ in this “quasi-sweeps” regime, we must revisit the equations for $w_m(x)$ and x_{cm} in Eqs. (S146) and (S147).

5.4.1 Location of the interference threshold

Motivated by our solution for the quasi-sweeps regime in SI Section 4.2, we anticipate that the dominant contributions to the y-integral in Eq. (S147) come from the Haldane region of the shoulder solution ($y > x_{cm}$). When $s \gtrsim x_{cm} + \mathcal{O}(\sqrt{v})$, the integrand will be peaked around a characteristic value $y^* = s \pm \mathcal{O}(\sqrt{v})$. Performing a Gaussian Laplace approximation around this maximum yields,

$$1 \approx \frac{\sqrt{2\pi}}{x_{cm}\sqrt{v}} e^{\frac{x_{cm}^2}{2v}} \int_{x_{cm}}^{\infty} s\mu_m(s) ds. \quad (\text{S217})$$

This is similar to the single-effect model in SI Section 4.2.1 if we take

$$U'_b = \int_{x_{cm}}^{\infty} \mu_m(s) ds, \quad s'_b = \frac{1}{U'_b} \int_{x_{cm}}^{\infty} s\mu_m(s) ds. \quad (\text{S218})$$

The decomposition in Eqs. (S146) and (S147) shows that this s'_b integral coincides with the distribution of fixed mutations in the modifier lineage. In contrast to the multiple mutations regime in SI Section 5.3, the distribution in Eq. (S218) may have a typical scale but will generally not be strongly peaked.

This quasi-sweeps regime will be self-consistently valid if Eq. (S217) is much larger than the remaining contributions from $y \leq x_{cm}$:

$$\int^{x_{cm}} \frac{\mu_m(s)}{v} e^{\frac{x_{cm}s}{v} - \frac{s^2}{2v}} \ll \frac{\sqrt{2\pi}}{x_{cm}\sqrt{v}} e^{\frac{x_{cm}^2}{2v}} \int_{x_{cm}}^{\infty} s\mu_m(s) ds. \quad (\text{S219})$$

For the stretched exponential distributions in Eq. (S156), this will only be true if $x_{cm} \lesssim s'_0$, so that $U'_b \sim \mathcal{O}(U'_0)$ and $s'_b \sim \mathcal{O}(s'_0)$. This shows that the quasi-sweeps regime will be valid when

$$s'_0 \gg x_{cm} \sim \sqrt{2v \log\left(\frac{v}{U'_0 s'_0}\right)}. \quad (\text{S220})$$

5.4.2 Extending the shoulder solution to lower fitness values

We can continue this line of reasoning to compute the shape of $w_m(x)$ for fitness values below x_{cm} . For mutations with fitness effects $s \gtrsim x_{cm} - x + \mathcal{O}(\sqrt{v})$, the y-integral in Eq. (S146) will continue to have a

large contribution from the Haldane region of the shoulder solution. For these values of s , we can repeat the Gaussian integration in Eq. (S128) to obtain

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2}{2v}} \int_{x_{cm}-x+O(\sqrt{v})}^{\infty} \mu_m(s) \left\{ \frac{s\sqrt{2\pi}}{x_{cm}\sqrt{v}} \cdot \left[\Phi\left(\frac{x}{\sqrt{v}}\right) - \Phi\left(\frac{x_{cm}-s}{\sqrt{v}}\right) \right] - \frac{e^{-\frac{x^2}{2v}} - e^{-\frac{(x_{cm}-s)^2}{2v}}}{x_{cm}} \right\} ds, \quad (\text{S221})$$

where $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-u^2/2} du$ is the Gaussian cumulative distribution function. When $x \gtrsim O(\sqrt{v})$, the Gaussian cumulative function is close to one, and Eq. (S221) reduces to

$$w_m(x) \approx 2x_{cm} e^{\frac{x^2-x_{cm}^2}{2v}} \frac{\sqrt{2\pi}}{x_{cm}\sqrt{v}} \int_{x_{cm}}^{\infty} \mu_m(s) s ds = 2x_{cm} e^{\frac{x^2-x_{cm}^2}{2v}}, \quad (\text{S222})$$

where the last line follows from substituting the auxiliary condition for x_{cm} in Eq. (S217). This shows that the shoulder solution will be valid down to $x_{\min} = O(\sqrt{v})$, similar to the single-effect model in SI Section 4.2.2.

When $x \lesssim O(\sqrt{v})$, the shape of $\Phi(x)$ will start to become important. At this point, the boundary terms at $x_{cm} - s$ will become negligible, so that Eq. (S221) reduces to

$$\begin{aligned} w_m(x) &\approx 2x_{cm} e^{\frac{x^2-x_{cm}^2}{2v}} \int_{x_{cm}-x}^{\infty} \mu_m(s) \left[\frac{s\sqrt{2\pi} v e^{\frac{x_{cm}^2}{2v}}}{v x_{cm}} \cdot \Phi\left(\frac{x}{\sqrt{v}}\right) - \frac{e^{\frac{x_{cm}^2-x^2}{2v}}}{x_{cm}} \right] ds, \\ &\approx 2x_{cm} e^{\frac{x^2-x_{cm}^2}{2v}} \left[\Phi\left(\frac{x}{\sqrt{v}}\right) \left(1 - \frac{\int_{x_{cm}}^{x_{cm}-x} s \mu_m(s) ds}{\int_{x_{cm}}^{\infty} s \mu_m(s) ds} \right) - \frac{\sqrt{v}}{s'_b \sqrt{2\pi}} e^{-\frac{x^2}{2v}} \left(1 - \frac{\int_{x_{cm}}^{x_{cm}-x} \mu_m(s) ds}{\int_{x_{cm}}^{\infty} \mu_m(s) ds} \right) \right], \end{aligned} \quad (\text{S223})$$

where the second line follows from substituting the auxiliary condition for x_{cm} in Eq. (S217). This will coincide with the analogous expression for the single- s model in 4.1.2 if

$$\int_{x_{cm}}^{x_{cm}-x} \mu_m(s) ds \ll \int_{x_{cm}}^{\infty} \mu_m(s) ds. \quad (\text{S224})$$

This condition will be satisfied if $x_{cm} - x \lesssim s'_b$, or alternatively, if $x \gtrsim x_{cm} - s'_b$.

For lower initial fitnesses, the contributions from mutations that land below the interference threshold ($x + s < x_{cm}$) can begin to become important. We can estimate the onset of these effects by considering initial fitnesses in the range $-O(\sqrt{v}) \gtrsim x \gtrsim x_{cm} - s'_b$. In this limit, Eq. (S223) reduces to

$$w_m(x) \approx \frac{2U'_b(s'_b - |x|)}{|x|} \approx \frac{2x_{cm}\sqrt{v}}{\sqrt{2\pi}|x|} \left(1 - \frac{|x|}{s'_b} \right) e^{-\frac{x_{cm}^2}{2v}}, \quad (\text{S225})$$

where the second line follows from substituting the auxiliary condition for x_{cm} in Eq. (S217). The corresponding contributions from mutations with $x + s < x_{cm}$ in Eq. (S146) are dominated by the upper limit of the y -integral, which yields an additional correction of order

$$\delta w_m(x) \approx 2x_{cm} e^{\frac{x^2-x_{cm}^2}{2v}} \int_{O(\sqrt{v})-x}^{x_{cm}-x} \mu_m(s) e^{\frac{xs}{v} + \frac{s^2}{2v}} \frac{ds}{s}. \quad (\text{S226})$$

If the curvature of $\mu_m(s)$ is not too strong [$\partial_s \log \mu_m(s'_b) \ll s'_b/v \ll s'_b$], this s -integral will be dominated by the upper limit of integration, so that

$$\delta w_m(x) \approx \frac{2v \cdot \mu(x_{cm} - x)}{x_{cm} - x}. \quad (\text{S227})$$

This will be small compared to Eq. (S225) as long as $s'_b \gtrsim x_{cm} \gg \sqrt{v}$, which shows that Eq. (S223) will be a good approximation for $w_m(x)$ as long as $x \gtrsim x_{cm} - s'_b$.

For lower initial fitnesses ($x \lesssim x_{cm} - s'_b$), these additional contributions will start to become important. In principle, we can generalize the recursive approach in SI Section 4.2.2 to extend our solution for $w_m(x)$ below this point. However, these corrections will no longer be universal, and will depend on other features the mutation spectrum beyond the first two moments in Eq. (S218). Since our results in the main text will only require the portion of the solution for $x \gtrsim x - s'_b$, we will leave these calculations for future work.

6 Incorporating deleterious mutations

While the previous sections have focused on changes to the beneficial mutation spectrum, the vast majority of new mutations are neutral or deleterious (49). The most general evolvability modifiers could alter the rates and fitness costs of these deleterious mutations as well.

Previous work (50, 112, 114) has shown in our parameter regime of interest, deleterious mutations can be divided into two broad categories: (i) quasi-neutral mutations ($|s| \ll v/x_c$), which can frequently hitchhike to fixation with beneficial mutations, and (ii) “purgeable” deleterious mutations ($|s| \gg v/x_c$), which rarely reach high frequencies, but can collectively reduce the effective population size if their mutation rate is sufficiently high. Previous estimates suggest that the overall mutation rates in many natural and laboratory microbial populations are sufficiently low ($U \lesssim 10^{-3}$ (66)) that both categories have a negligible impact on the overall the rate of adaptation (50, 112). However, differences in the deleterious distribution of fitness effects can still have a strong influence on the fixation probability of a modifier mutation (26).

We can incorporate purgeable deleterious mutations into our framework using the argument outlined in Ref. (26), which we reproduce here for completeness. Ref. (26) showed that it is useful to rewrite the mutation term in Eq. (S34) as separate beneficial and deleterious contributions,

$$\begin{aligned} 0 = & \underbrace{x \cdot w_m(x)}_{\text{selection}} - \underbrace{v \cdot \partial_x w_m(x)}_{\text{competition w/ wildtype}} - \underbrace{\frac{1}{2} \cdot w_m(x)^2}_{\text{drift while rare}} + \underbrace{\int_0^\infty \mu_m(s) [w_m(x+s) - w_m(x)] ds}_{\text{beneficial mutations}} \\ & + \underbrace{\int_0^\infty \mu_m(-\delta) [w_m(x-\delta) - w_m(x)] d\delta}_{\text{deleterious mutations}} \end{aligned} \quad (\text{S228})$$

where we are indexing the deleterious mutations by their overall magnitude $\delta \equiv -|s|$. The definition of a purgeable mutation is that they rarely reach high frequencies, which implies that the fixation probability after acquiring a purgeable mutation must be close to zero. If all of the deleterious mutations in μ_m fall in this purgeable category, we can therefore neglect the $w_m(x - \delta)$ term in Eq. (S228). This allows us to rewrite Eq. (S228) in the simpler form,

$$0 = (x - U'_d)w_m(x) - v \cdot \partial_x w_m(x) - \frac{1}{2} \cdot w_m(x)^2 + \int_0^\infty \mu_m(s) [w_m(x+s) - w_m(x)] ds \quad (\text{S229a})$$

where U'_d denotes the total mutation rate for purgeable mutations:

$$U'_d = \int_0^\infty \mu_m(-|\delta|) d\delta. \quad (\text{S229b})$$

This equation has the same form as the purely beneficial case analyzed above, but with the x coordinate shifted by $-U'_d$. The solutions can therefore be expressed in the simple form

$$w_m(x) = w_m^b(x - U'_d), \quad (\text{S230})$$

where $w_m^b(x)$ denotes our existing results for the purely beneficial case. Likewise, previous work (26, 50) has shown that the background fitness distribution satisfies an analogous formula,

$$f(x) = f^b(x - U_d), \quad (\text{S231a})$$

where U_d is the rate of producing purgeable mutations in the wildtype background,

$$U_d = \int_0^\infty \mu(-|\delta|) d\delta. \quad (\text{S231b})$$

Substituting these expressions into Eq. (2), we find that the fixation probability of the modifier can be written as

$$p_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) = \int f^b(x - U_d) \cdot w_m^b(x - U'_d + s_m) dx = \int f^b(y) \cdot w_m^b(y + s_m - U'_d + U_d) dy, \quad (\text{S232})$$

which implies that differences in the purgeable mutation rate behave like effective direct cost,

$$s_m^{\text{eff}} = U_d - U'_d. \quad (\text{S233})$$

We validated this approximation using Wright-Fisher simulations (Fig. 3), and found that it encompasses a broad range of deleterious fitness effects, including those much smaller than a typical beneficial driver mutation (Fig. 3B). Deviations eventually occur for more weakly selected mutations, which have a lower impact on the fixation probability than Eq. (S233) would predict. This suggests that purgeable mutations represent an upper bound on the strength of second-order selection.

Interestingly, the mapping in Eq. (S233) implies that the same modulation effect that occurs between first- and second-order selection will also apply for selection on simultaneous changes to robustness ($U_d \rightarrow U'_d$) and evolvability ($\mu(s) \rightarrow \mu_m(s)$). In particular, it implies that larger populations are more likely to trade reduced robustness for long-term gains in evolvability. Since empirical deleterious mutation rates are often comparatively low ($|s_m^{\text{eff}}| \lesssim s_b$), even the maximum possible reduction in robustness will be overpowered by marginal increases in evolvability (Fig. 3C). The opposite scenario – greedily selecting for increased robustness despite long-term reductions in evolvability – is also theoretically possible, but we expect that it will be less relevant in practice due to the low deleterious mutation rates in many organisms.

Finally, we note that robustness and evolvability do not have to trade off with one other – several recent studies have shown that they can sometimes interact synergistically as well (63, 65). Our analytical framework also provides new predictions for these synergistic mutations, allowing us to determine how each phenotype contributes to the lineage's long-term evolutionary fate (Fig. 3E). Our results show that larger populations will tend to weigh proportional enhancements in evolvability ($\Delta v/v$) more heavily than comparable increases in robustness ($\Delta U_d/U_d$). This is qualitatively different from the successive sweeps picture (SI Section 2), which predicts that the opposite ordering should occur. These examples illustrate how clonal interference can reshape our intuition about second-order selection for evolvability.

7 Extensions to more general fitness landscapes

Our analysis above focused on the simplest possible model of an evolvability modifier, where the local distribution of fitness effects could be approximated by a pair of fixed distributions,

$$\mu(s|\vec{g}) \approx \mu(s), \quad \mu_m(s|\vec{g}) \approx \mu_m(s), \quad (\text{S234})$$

for all genotypes \vec{g} (SI Section 3.1). In this section, we consider several extensions of our model that relax this assumption in different ways. In each case, we find that our existing theory provides a useful baseline for incorporating these new effects.

7.1 Weak Macroscopic Epistasis

It is clear that some deviations from Eq. (S234) will be too small to be relevant to natural selection. We can formalize this idea by considering subtle deviations of the form,

$$\mu(s|\vec{g}) = \mu(s) + \epsilon(s, \vec{g}), \quad (\text{S235})$$

where the function $\epsilon(s, \vec{g})$ represents a small perturbation. In this case, our existing results in SI Section 5 allow us to precisely define what we mean by “small”. In particular, if we let $\delta\mu(s) = \epsilon(s, \vec{g})$, then we can conclude that any perturbation for which the integral Eq. (8) is small [$I(\delta\mu) \lesssim s_b/x_c$] will be essentially invisible to natural selection. Any deviations from Eq. (S234) that fall below this minimum resolution will therefore not affect our main results.

We note that this space of “negligible” perturbations can be quite large from the perspective of the original DFE. Fig. 4 shows that even large fluctuations in $\mu(s)$ can be tolerated for fitness effects $\lesssim s_b(\mu, N)$, even if they cause in large shifts in the overall mean and height of the DFE. Conversely, much smaller perturbations at fitness effects $\gtrsim s_b(\mu, N)$ can lead to important deviations from our original model, even if they nominally appear to satisfy Eq. (S234).

We also note that this argument constitutes an upper bound on the impact of a given deviation from Eq. (S234), since our derivation of Eq. (8) assumed that the evolvability differences were permanent. This means that some values of $\epsilon(s, \vec{g})$ that exceed the resolution limit above could still have a small influence on the results because their evolvability differences are only transient. We consider such cases in more detail below.

7.2 Transient differences in evolvability

In many cases of interest, the evolvability benefits of a modifier will not persist indefinitely (as in Eq. S234) but will only apply within some local region of genotype space. For example, in the stability-activity landscape in Extended Data Fig. 2C, the stability enhancing modifier creates an opportunity for K additional mutations to accumulate. At a formal level, this constitutes a large deviation from Eq. (S234).

However, it is also clear that such permanent shifts are not truly necessary. Second-order selection can only take place while the modifier is competing with the wildtype; once one of these lineages fixes, future changes in $\mu(s|\vec{g})$ can no longer contribute to the fixation probability. Our results above allow us to estimate the size of this critical window. In particular, our heuristic analysis in SI Section 4.1.4 shows that the benefits of an evolvability modifier accumulate over $K^* \sim \max\{x_{cm}/s'_b, 1\}$ mutational steps. This suggests that changes to the DFE that occur outside this horizon will have a negligible impact on our results.

We tested this prediction by considering a simple model where the modifier reverts to the wildtype DFE after K mutational steps:

$$\mu_m(s|\vec{g}) = \begin{cases} \mu_m(s) & \text{if } |\vec{g} - \vec{g}_m| < K, \\ \mu(s) & \text{else.} \end{cases} \quad (\text{S236})$$

Simulations of this model show that our existing theory continues to provide a good approximation even when K is as low as 2 or 3 (Fig. 4A), consistent with the moderate values of $q \equiv x_c/s_b$ that are attained for many empirically relevant parameter values (42, 51). Similar results are also observed if the modifier reverts to a dead-end rather than the wildtype (Fig. 4B), with differences arising only for the smallest values of K . In both cases, we find that the evolvability benefits still confer exponentially large advantages over the classical SSWM expectation (Fig. 4A), even when they last for just a single mutation ($K = 1$).

7.3 Global diminishing returns epistasis

Many evolving populations exhibit a form of diminishing returns epistasis, where the fitness benefits of new mutations systematically decline with a lineage's absolute fitness (9, 10, 44, 58, 69, 98). To understand how this phenomenon might impact our results, we considered a simple model of global epistasis inspired by the long-term evolution experiment in Fig. 1A (58), where the typical fitness benefits of new mutations decline exponentially with the total fitness,

$$\mu(s|\vec{g}) = U_b \cdot \delta(s - \tilde{s}_b e^{-X(\vec{g})/\theta}), \quad (\text{S237})$$

where θ controls the strength of the diminishing returns effect. In the simplest scenario, we can assume that the modifier distribution exhibits a similar decline,

$$\mu_m(s|\vec{g}) = U'_b \cdot \delta(s - \tilde{s}'_b e^{-X(\vec{g})/\theta}), \quad (\text{S238})$$

but with different values of U'_b and \tilde{s}'_b . This generalizes the simple toy model in Fig. 2 to allow for steadily declining fitness effects.

Adiabatic approximation. The simplest behavior occurs when θ is sufficiently large. Suppose that the modifier arises at time t_0 when the mean fitness of the population is $\bar{X}(t_0)$. If we could neglect the additional deceleration that occurs over the lifetime of a single mutation, we could simply apply our existing results with $s_b = \tilde{s}_b e^{-\bar{X}(t_0)/\theta}$ and $s'_b = \tilde{s}'_b e^{-\bar{X}(t_0)/\theta}$. We can use our heuristic picture in SI Section 4.1.4 to determine when this ‘‘adiabatic approximation’’ will be valid. The fate of a mutation is determined over $T_c \sim x_c/v$ generations, during which time the total fitness of the population increases by $\Delta\bar{X} \sim x_c$. If $\theta \gg x_c$, then diminishing returns epistasis will have a negligible effect over the lifetime of a single mutation, even if it had a large effect in setting the initial scale of selection ($\bar{X}(t_0) \gtrsim \theta$). This provides a self-consistent justification for our adiabatic approximation above.

This simple picture leads to novel predictions when combined with our existing results in Eqs. (3) and (4). In particular, it suggests that a selection-strength modifier will become less strongly selected over time as the population becomes better adapted, while the benefits of a mutator allele will remain roughly constant over the same time window. Since x_c is roughly proportional to s_b (Eq. S50), we expect that θ/x_c will grow increasingly large at long times as s_b becomes progressively smaller (Fig. 5). This suggests that the adiabatic approximation will often be useful for understanding the long-term dynamics of a population (69).

Beyond the adiabatic approximation. The situation becomes more complicated when x_c is comparable to θ , since the fitness benefits of new mutations will now decline within the scale of the population fitness distribution. The dynamics in this regime are poorly understood even in the absence of second-order selection (69). However, we can still account for these effects in a crude way by leveraging our heuristic picture in SI Section 4.1.4, and focusing on the leading-order corrections when x_c/θ is small but finite. We briefly review the results for our original model below, and then show how they can be extended to calculate the leading-order corrections from diminishing returns epistasis.

Recall that for small changes in the selection coefficient ($s'_b - s_b \ll s'_b$), successful modifiers will typically arise in the high-fitness nose of the population ($x \approx x_c$) and acquire x_c/s_b mutations before reaching appreciable frequencies (SI Section 4.1.4). In each of these steps ($j = 1, \dots, x_c/s_b$), a selection-strength modifier will grow for $\tau_{\text{est}} \approx s_b/v$ generations while the next nose establishes. During this establishment time, the modifier will produce $\sim \exp\left[\frac{s_b}{v} \cdot j(s'_b - s_b)\right]$ more mutations than a wildtype individual with the same fitness, leading to the scaling

$$\log \tilde{p}_{\text{fix}} \sim \frac{s_b}{2v} \cdot \frac{x_c}{s_b} \left(\frac{x_c}{s_b} - 1 \right) (s'_b - s_b). \quad (\text{S239})$$

This picture becomes more complicated with diminishing returns epistasis since the wildtype and modifier selection coefficients will decline with each subsequent mutation,

$$s'_{b,j} \equiv s'_b(t_0) e^{-\sum_{i=1}^{j-1} s'_{b,i}/\theta}, \quad (\text{S240a})$$

$$s_{b,j} \equiv s_b(t_0) e^{-\sum_{i=1}^{j-1} s_{b,i}/\theta}, \quad (\text{S240b})$$

where $s'_b(t_0)$ and $s_b(t_0)$ denote the effective selection coefficients at the time that the modifier arises. As a result, the modifier's growth advantage, and the time it spends growing, will also vary with each mutation. We can straightforwardly extend our heuristic prediction to account for these differences,

$$\log \tilde{p}_{\text{fix}} \sim \sum_{i=1}^{\frac{x_c}{s_b} - 1} \sum_{j=1}^i \tau_{\text{est},j} \cdot s_{b,j} \cdot \left(\frac{s'_{b,j}}{s_{b,j}} - 1 \right), \quad (\text{S241})$$

where $\tau_{\text{est},j}$ denotes the establishment time of the j^{th} mutational step. These establishment times will emerge from the non-equilibrium dynamics imposed by diminishing returns epistasis, including the declining rate of fitness increase and width of the fitness distribution. While these dynamics remain poorly characterized even in the absence of second order selection, we find that we can obtain accurate predictions by assuming that the product $\tau_{\text{est},j} \cdot s_{b,j}$ remains approximately constant during the lifetime of the modifier:

$$\tau_{\text{est},j} \cdot s_{b,j} \approx \tau_{\text{est}}(t_0) \cdot s_b(t_0). \quad (\text{S242})$$

This assumption is motivated by the logarithmic dependence of $\tau_{\text{est}} \cdot s_b$ on s_b in the absence of diminishing returns epistasis (51; Eq. S50), which suggests that treating it as constant over fixation timescales will often be a good approximation. With this assumption, Eq. (S241) reduces to the simpler form

$$\log \tilde{p}_{\text{fix}} \sim \sum_{i=1}^{\frac{x_c}{s_b} - 1} \sum_{j=1}^i \tau_{\text{est}}(t_0) \cdot s_b(t_0) \cdot \left(\frac{s'_{b,j}}{s_{b,j}} - 1 \right). \quad (\text{S243})$$

Under the model in Eq. (S240), the diminishing returns effects on $s'_{b,j}$ and $s_{b,j}$ exactly cancel each other, so that we end up with the same adiabatic approximation above. This suggests that the next-order corrections remain small when the modifier and wildtype follow similar diminishing returns schedules.

We can also consider modifiers that alter the form of the diminishing returns epistasis itself. The simplest example is one that changes the diminishing returns parameter θ to a new value θ_m once the modifier mutation arises:

$$\mu_m(s|\vec{g}) = U'_b \cdot \delta \left(s - \tilde{s}'_b e^{-X(\vec{g}_m)/\theta - [X(\vec{g}) - X(\vec{g}_m)]/\theta_m} \right), \quad (\text{S244})$$

where \vec{g}_m denotes the founding genotype of the modifier lineage. The difference between θ and θ_m is invisible in the adiabatic approximation, but the next order corrections can be obtained from Eq. (S243), by substituting

$$s'_{b,j} \equiv s'_b(t_0) e^{-\sum_{i=1}^{j-1} s'_{b,i}/\theta_m}, \quad (\text{S245a})$$

$$s_{b,j} \equiv s_b(t_0) e^{-\sum_{i=1}^{j-1} s_{b,i}/\theta}. \quad (\text{S245b})$$

The diminishing returns contributions no longer cancel, leading to modest corrections to our existing theory (blue lines, Fig. 5D). In particular, when $\theta_m \gg \theta$, these diminishing returns modifiers can be positively selected even when their initial DFEs are the same [$s'_b(t_0) = s_b(t_0)$], since the modifier experiences fewer diminishing returns effects at later times. These differences disappear once $s'_b(t_0)$ is sufficiently large that we enter the quasi-sweeps regime. In this case, we have seen that the success of the modifier is driven by its ability to acquire its first additional mutation (SI Section 4.2.3), where the new value of θ_m has not yet taken effect. These extensions of our heuristic analysis suggest that the simple model we have studied in this work may be a powerful tool to understand selection for evolvability on more general fitness landscapes.

8 Numerical Methods

8.1 Theoretical predictions

The theory curves in each of the figures were generated using the following procedures.

Figure 2. To generate the theoretical predictions in Fig. 2, we first solved for x_{cm} numerically as a function of s'_b , U'_b , and v . We found that numerical accuracy was improved by using a modified version of Eqs. (S57) and (S119),

$$1 = \frac{U'_b}{s'_b} \left[e^{\frac{x_{cm}s'_b}{v} - \frac{s'^2_b}{2v}} + \frac{\sqrt{2\pi}s'^2_b}{x_{cm}\sqrt{v}} e^{\frac{x^2_{cm}}{2v}} \Phi \left(\frac{s'_b - x_{cm}}{\sqrt{v}} \right) \right], \quad (\text{S246})$$

which has the same asymptotic behavior, but smoothly captures the transition between the multiple-mutations ($x_{cm} - s'_b \gg \sqrt{v}$) and quasi-sweeps regimes ($s'_b - x_{cm} \gg \sqrt{v}$). Numerical solutions to Eq. (S246) were obtained using the `fsolve` function from the SciPy library (115), with the measured value of v obtained from our forward-time simulations. We used the same procedure to solve for the wildtype interference threshold x_c using the analogous version of Eq. (S246). If $x_{cm} > s'_b$, we estimated the fixation probability of the modifier using the multiple-mutations expression in Eq. (S85), while the quasi-sweeps prediction in Eq. (S136) was used for $x_{cm} < s'_b$. We used this procedure to generate all of the theory lines in panels A-C.

The boundary for the gray region in Fig. 2D was obtained by identifying the values of U'_b and s'_b that minimize the true fixation probability $\tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b))$ at a fixed value of $x = \tilde{p}_{\text{fix}}(s_b \rightarrow s'_b) \cdot \tilde{p}_{\text{fix}}(s_b \rightarrow s'_b)$. Using the leading-order approximations in Eqs. (S89), (S90), and S91, the fixation probability of the compound modifier can be approximated as

$$\tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) \approx x^{s'_b/s_b} \cdot \tilde{p}_{\text{fix}}(s_b \rightarrow s'_b)^{1-s'_b/s_b}, \quad (\text{S247})$$

which eliminates the explicit dependence on U'_b . Minimizing this function with respect to s'_b yields an analytical approximation for the boundary,

$$y(x) = \min_{s'_b|x} \left\{ \log \tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b)) \right\} \approx \log(x) - \frac{1}{4} \log^2(x) \left[\frac{x_c^2}{2v} - \frac{x_c s_b}{2v} \right], \quad (\text{S248})$$

which was used to define the gray region in Fig. 2D.

Figure 3. The theoretical predictions for Fig. 3 were obtained using a similar procedure as in Fig. 2. For a modifier with a direct cost ($s_m < 0$; Fig. 3A,C), we used the multiple-mutations expression in Eq. (S100) if $x_{cm} > s'_b$, and the quasi-sweeps expression in Eq. (S140) if $x_{cm} < s'_b$; the integrals in Eq. (S140) were computed numerically using the quad function from the SciPy library (115). For a modifier with a direct benefit ($s_m > 0$; Fig. 3B,C), we used the corresponding predictions in Eq. (S98) when $x_{cm} < x^*$, or the dead-end predictions in Eq. (S116) otherwise, with x^* defined by Eq. (S112).

The transition line in Fig. 3C was obtained by solving for the critical value of s_m that satisfies

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) = 1. \quad (\text{S249})$$

The theoretical line was obtained by numerically solving Eq. (S249) using the fixation probability predictions described above. The simulated values were calculated by linearly interpolating the observed fixation probabilities using the `polyfit` package in the Numpy library.

Figure 4. Predictions for the continuous distributions of fitness effects in Fig. 4 were obtained using the single-s mapping described in SI Section 5. We first calculated the value of x_c by numerically solving the wildtype version of Eq. (S246), with the effective parameters s_b and U_b defined by Eqs. (S152) and (S155). Using these estimates, we next checked whether the modifier was in the perturbative regime by numerically solving for δx_c using Eq. (S169). If $\delta x_c < v/s_b$, we used the corresponding predictions from the perturbative regime in Eq. (S175), with $s'_b \approx s_b$. If $\delta x_c > v/s_b$, we turned to the corresponding predictions from the modifier-dominated regime in SI Sections 5.3 and 5.4.

To do so, we numerically solved for x_{cm} using Eq. (S246) with s'_b and U'_b defined by Eqs. (S177) and (S180). For the bimodal distributions in Fig. 4, the solutions to Eqs. (S177) and (S180) can be expressed as a piecewise function,

$$s'_b, U'_b = \begin{cases} s_b^\mu(x_{cm}), U_b^\mu(x_{cm}) & \text{if } U_b^\mu \tilde{p}_{\text{fix}}(s_b^\mu | x_{cm}, v) > U_1 \tilde{p}_{\text{fix}}(s_1 | x_{cm}, v), \\ s_1, U_1 & \text{else,} \end{cases} \quad (\text{S250})$$

where $s_b^\mu(x_{cm})$ and $U_b^\mu(x_{cm})$ refer to the effective parameters of the wildtype distribution, but evaluated at x_{cm} rather than x_c . The fixation probability of the modifier was then estimated as

$$\tilde{p}_{\text{fix}}(\mu(s) \rightarrow \mu_m(s), s_m) \approx \tilde{p}_{\text{fix}}((U_b, s_b) \rightarrow (U'_b, s'_b), s_m), \quad (\text{S251})$$

where the right hand side was calculated using the methods described for Figs. 2 and 3 above.

To capture the transition region between the perturbative- and modifier-dominated regimes, we continued to use our perturbative approximation if the modifier-dominated estimate of $|x_{cm} - x_c|$ was less than the perturbative calculation of $|\delta x_c|$. This convention ensures that our theoretical predictions are continuous at the border between the two regimes.

Figure 5. The theoretical predictions for Fig. 5C were obtained using Eq. (S100). The theoretical predictions for Fig. 5D were obtained using the same procedure in Fig. 2 for a beneficial modifier with no direct cost or benefit.

Extended Data Fig. 3. The theoretical predictions for Fig. 3B,C were obtained using Eq. (S74). The theoretical predictions for Fig. 3D,E were obtained using the procedure used in Fig. 2 for a beneficial modifier without a short term cost.

Extended Data Fig. 4. The theoretical predictions for Fig. 4 were obtained using the same procedure in Fig. 2 for a beneficial modifier with no direct cost or benefit.

Extended Data Fig. 5. The forward-time simulations with diminishing returns epistasis in Fig. 5 followed a similar procedure to that outlined in Methods. To enhance reproducibility, the simulations were allowed to “burn-in” for $\Delta t = 2 \cdot 10^4$ generations before initiating diminishing returns. After this burn-in period, diminishing returns was initiated and the mean fitness and typical selection coefficient, $s_b(\bar{X})$, were recorded every 50 generations. The simulation results in Fig. 5A,B are the mean and average of 10 simulations.

After the typical selection coefficient in the population had declined to the chosen value $s_b(\bar{X}) = 2.5 \cdot 10^{-2}$, the fitness and abundance of each lineage was saved. Modifier lineages were then introduced at a constant rate for $t_U = 1000$ generations. This value was capped to ensure a successful modifier arose in a population with the chosen selection coefficient. To generate the comparisons of the neutral modifier and background selection coefficients in Fig. 5C, $s_b(\bar{X} + x_c)$ was obtained using an estimate of x_c from the model without diminishing returns. This estimate was obtained using a similar procedure to that in Fig. 2 with v numerically evaluated in a population with a constant selection coefficient, $s_b = 2.5 \cdot 10^{-2}$. The results in Fig. 5C were then obtained using Eq. (S240).

After $t_U = 1000$, U_m was set to zero and the simulation continued until a modifier took over or all modifier lineages were purged. If a modifier did not take over, the population reverted back to the saved lineage distribution and this process was repeated n times until modifier took over. This allowed us to apply the same procedure in Methods to calculate the fixation probability of a modifier with $T = n \cdot t_U$. The theoretical predictions for the $\theta = \theta_m$ line were obtained using the same procedure in Fig. 2 for a selection strength modifier with $s_b = \tilde{s}_b$ and $s'_b = \tilde{s}'_b$. The theoretical predictions for the $\theta = \infty$ line were obtained using $\min\{\text{Eq. (S243), Eq. (S136)}\}$, where $s_{b,0} \cdot \tau_{\text{est},0}$ was calculated from the relation $s_{b,0} \cdot \tau_{\text{est},0} = s_{b,0}^2/v$. In this case, v was evaluated numerically in Wright Fisher simulations without diminishing returns for $s_b = s_b(\bar{X} + x_c)$.

8.2 Empirical example from Ref. (8)

The relative fitness estimates for the modifier example in Fig. 1A were obtained from Ref. (8). This study examined two strains of *E. coli* that were isolated from generation 500 of Lenski’s long-term evolution experiment (79). Relative fitnesses of the two strains at the first timepoint in Fig. 1A were obtained from head-to-head competitions under the same conditions as the original experiment. Ref. (8) reported these relative fitness values using the metric,

$$W \equiv \frac{\log\left(\frac{N_2(\Delta t) \cdot 2^{\Delta t}}{N_2(0)}\right)}{\log\left(\frac{N_1(\Delta t) \cdot 2^{\Delta t}}{N_1(0)}\right)}, \quad (\text{S252})$$

where $N_i(0)$ is the number of colonies of strain i observed at the beginning of the competition, $N_i(\Delta t)$ is the (adjusted) number of colonies observed at the end of the competition, and Δt is the length of the competition

in generations. To enable more direct comparisons with our theory, we converted these W estimates to the relative (log) fitness,

$$\Delta X \equiv \frac{1}{\Delta t} \log \left[\frac{f_2(\Delta t)}{f_1(\Delta t)} \cdot \frac{f_1(0)}{f_2(0)} \right], \quad (\text{S253})$$

where $f_i(t)$ is the relative frequency of strain i at generation t .

To perform this conversion, we assumed that the colony counts in Ref. (8) were adjusted so that a similar overall number of colonies were measured at both timepoints. This implies that $N_i(\Delta t)/N_i(0) = f_i(\Delta t)/f_i(0)$. If we also assume that the competition experiments were started at a 1:1 ratio, so that $f_1(0) \approx f_2(0) \approx 1/2$, we can obtain a relation between the ΔX and W metrics,

$$W = 1 + \frac{\Delta X \cdot \Delta t}{\log \left(\frac{2^{\Delta t+1}}{1+e^{\Delta X \cdot \Delta t}} \right)}, \quad (\text{S254})$$

which depends on the duration Δt . We used this expression to convert the average and 95% confidence intervals reported in Ref. (8) into the relative fitness values in Fig. 1A:

Initial timepoint. The fitness measurements for the initial timepoint were conducted over $\Delta t = 47$ generations, and yielded a fitness disadvantage of $W=0.937$ (0.934,0.941) [mean and 95% confidence intervals from replicate competition assays]. Numerical conversion using Eq. (S254) yielded a log fitness difference of $\Delta X = -4.4\%$ (-4.17%, -4.66%).

Second timepoint. The fitness measurements at the second timepoint were obtained after evolving each isolate for an additional 883 generations in 20 independent replay experiments. Pooled fitness measurements of the evolved populations were performed over $\Delta t = 6.6$ generations, and yielded a fitness benefit of $W=1.02$ (1.003, 1.039). Numerical conversion using Eq. (S254) yielded a log fitness difference of $\Delta X = 1.4\%$ (0.2%, 2.6%).

Supplementary References

81. Johnson, T. The approach to mutation–selection balance in an infinite asexual population, and the evolution of mutation rates. *Proc. R. Soc. Lond. B* **266**, 2389–2397 (1999).
82. Desai, M. M. & Fisher, D. S. The balance between mutators and nonmutators in asexual populations. *Genetics* **188**, 997–1014 (2011).
83. James, A. & Jain, K. Fixation probability of rare nonmutator and evolution of mutation rates. *Ecol. Evol.* **6**, 755–764 (2016).
84. Kondrashov, D. A. & Kondrashov, F. A. Topological features of rugged fitness landscapes in sequence space. *Trends Genet.* **31**, 24–33 (2015).
85. Leigh, E. G. The evolution of mutation rates. *Genetics* **73**, Suppl 73:1–18 (1973).
86. Giles, L. E. Natural selection and mutability. *Am. Nat.* **104**, 301–305 (1970).
87. Painter, P. R. Mutator genes and selection for the mutation rate in bacteria. *Genetics* **79**, 649–660 (1975).
88. Gillespie, J. H. Mutation modification in a random environment. *Evolution* **35**, 468–476 (1981).
89. Ishii, K., Matsuda, H., Iwasa, Y. & Sasaki, A. Evolutionarily stable mutation rate in a periodically changing environment. *Genetics* **121**, 163–174 (1989).
90. Kessler, D. A. & Levine, H. Mutator dynamics on a smooth evolutionary landscape. *Phys. Rev. Lett.* **80**, 2012–2015 (1998).
91. Johnson, T. Beneficial mutations, hitchhiking and the evolution of mutation rates in sexual populations. *Genetics* **151**, 1621–1631 (1999).
92. Tanaka, M. M., Bergstrom, C. T. & Levin, B. R. The evolution of mutator genes in bacterial populations: the roles of environmental change and timing. *Genetics* **164**, 843–854 (2003).
93. Andre, J.-B. & Godelle, B. The evolution of mutation rate in finite asexual populations. *Genetics* **172**, 611–626 (2006).
94. Wylie, C. S., Ghim, C.-M., Kessler, D. & Levine, H. The fixation probability of rare mutators in finite asexual populations. *Genetics* **181**, 1595–1612 (2009).
95. Gardiner, C. *Handbook of Stochastic Methods* (Springer, 1985).
96. Gerrish, P. J. & Lenski, R. E. The fate of competing beneficial mutations in an asexual population. *Genetica* **102**, 127–144 (1998).
97. Good, B. H. & Desai, M. M. Fluctuations in fitness distributions and the effects of weak linked selection on sequence evolution. *Theor. Popul. Biol.* **85**, 86–102 (2013).
98. Reddy, G. & Desai, M. M. Global epistasis emerges from a generic model of a complex trait. *Elife* **10**, e64740 (2021).

99. Božić, N. M., Guo, Y., Rycroft, C. H. & Amir, A. How microscopic epistasis and clonal interference shape the fitness trajectory in a spin glass model of microbial long-term evolution. *Elife* <https://doi.org/10.7554/eLife.87895.2> (2023).
100. Weissman, D. B., Desai, M. M., Fisher, D. S. & Feldman, M. W. The rate at which asexual populations cross fitness valleys. *Theor. Popul. Biol.* **75**, 286–300 (2009).
101. Ochs, I. E. & Desai, M. M. The competition between simple and complex evolutionary trajectories in asexual populations. *BMC Evol. Biol.* **15**, 55 (2015).
102. Merrikh, C. N. & Merrikh, H. Gene inversion potentiates bacterial evolvability and virulence. *Nat. Commun.* **9**, 4662 (2018).
103. Horton, J. S., Flanagan, L. M., Jackson, R. W., Priest, N. K. & Taylor, T. B. A mutational hotspot that determines highly repeatable evolution can be built and broken by silent genetic changes. *Nat. Commun.* **12**, 6092 (2021).
104. Bloom, J. D., Gong, L. I. & Baltimore, D. Permissive secondary mutations enable the evolution of influenza oseltamivir resistance. *Science* **328**, 1272–1275 (2010).
105. Javanmardi, K. et al. Antibody escape and cryptic cross-domain stabilization in the SARS-CoV-2 omicron spike protein. *Cell Host Microbe* **30**, 1242–1254 (2022).
106. Rodrigues, J. V. et al. Biophysical principles predict fitness landscapes of drug resistance. *Proc. Natl Acad. Sci. USA* **113**, E1470–E1478 (2016).
107. Li, Y. et al. Patterns of somatic structural variation in human cancer genomes. *Nature* **578**, 112–121 (2020).
108. Venkataram, S. et al. Development of a comprehensive genotype-to-fitness map of adaptation-driving mutations in yeast. *Cell* **166**, 1585–1596.e22 (2016).
109. Tsimring, L. S., Levine, H. & Kessler, D. A. RNA virus evolution via a fitness-space model. *Phys. Rev. Lett.* **76**, 4440 (1996).
110. Rouzine, I. M., Wakeley, J. & Coffin, J. M. The solitary wave of asexual evolution. *Proc. Natl Acad. Sci. USA* **100**, 587–592 (2003).
111. Hallatschek, O. The noisy edge of traveling waves. *Proc. Natl Acad. Sci. USA* **108**, 1783–1787 (2011).
112. Melissa, M. J., Good, B. H., Fisher, D. S. & Desai, M. M. Population genetics of polymorphism and divergence in rapidly evolving populations. *Genetics* **221**, iyac053 (2022).
113. Neher, R. A. & Shraiman, B. I. Genetic draft and quasi-neutrality in large facultatively sexual populations. *Genetics* **188**, 975–996 (2011).
114. Schiess, S., Szöllősi, G. J., Mustonen, V. & Lässig, M. Emergent neutrality in adaptive asexual evolution. *Genetics* **189**, 1361–1375 (2011).
115. Virtanen, P. et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* **17**, 261–272 (2020).