

Survival of the frequent at finite population size and mutation rate: filling the gap between quasispecies and monomorphic regimes

Bhavin S. Khatri^{1,2}

¹*The Francis Crick Institute, 1 Midland Road, London, NW1 1AT, U.K.*

²*Department of Life Sciences, Imperial College London, Silwood Park, Ascot, SL5 7PY*

(Dated: 23 July 2018)

In recent years, there has been increased attention on the non-trivial role that genotype-phenotype maps play in the course of evolution, where natural selection acts on phenotypes, but variation arises at the level of mutations. Understanding such mappings is arguably the next missing piece in a fully predictive theory of evolution. Although there are theoretical descriptions of such mappings for the monomorphic ($N\mu \ll 1$) and deterministic or very strong mutation ($N\mu \gg 1$) limit, given by developments of Iwasa's free fitness and quasispecies theories, respectively, there is no general description for the intermediate regime where $N\mu \sim 1$. In this paper, we address this by transforming Wright's well-known stationary distribution of genotypes under selection and mutation to give the probability distribution of phenotypes, assuming a general genotype-phenotype map. The resultant distribution shows that the degeneracies of each phenotype appear by weighting the mutation term; this gives rise to a bias towards phenotypes of larger degeneracy analogous to quasispecies theory, but at finite population size. On the other hand we show that as population size is decreased, again phenotypes of higher degeneracy are favoured, which is a finite mutation description of the effect of sequence entropy in the monomorphic limit. We also for the first time (to the author's knowledge) provide an explicit derivation of Wright's stationary distribution of the frequencies of multiple alleles.

INTRODUCTION

The past 100 years has seen our understanding of the mechanisms of evolution develop, from its initial mathematical foundations due initially to Wright and Fisher, and later Kimura, which encompass a description of the interplay between selection, mutation and drift, to the current day with descriptions of multi-locus evolution with recombination, linkage and epistasis¹. However, as powerful as these studies are, they lack a crucial missing ingredient in our understanding of evolution, which is the role of genotype-phenotype maps, where selection acts on phenotypes, but underlying variation arises at the genetic level. In general, such mappings will be very complex, but a common theme is that because these mappings will often be many-to-one, some phenotypes will have more genotypes associated with them than others. In the weak mutation or monomorphic limit, where the population-scaled mutation rate is small ($N\mu \ll 1$), there is a complete theory that predicts the equilibrium distribution of phenotypes in the monomorphic limit², which naturally leads from Iwasa's definition of *free fitness*³ (subsequently rediscovered by Sella & Hirsh⁴). A general prediction of these theories is that at small population sizes, as genetic drift dominates favouring phenotypes with larger sequence entropy, which is the log degeneracy of the phenotype. In the opposite regime, with infinitely large population sizes there is Eigen's quasispecies theory, which are deterministic sets of equations describing the growth and mutation of many genotypes⁵; one of its predictions is that at sufficiently large mutation rates regions of locally high robustness in genotype space are favoured^{6,7}. Translating to phenotype space, it is trivial to see that phenotypes that have higher (average) local robustness will be favoured over those with lower robustness. In both the monomorphic and

quasispecies regimes, we see analogous effects related to non-optimal degenerate (robust) phenotypes being favoured at small population sizes (large mutation rates).

In this paper, we present a theory that straddles both these regimes to calculate the equilibrium distribution of phenotype frequencies at arbitrary and finite population sizes and mutation rates. This is done by transforming Wright's equilibrium (stationary) distribution of multiple alleles, cast with the mutational structure of a genotype space, to the space of phenotypes assuming a simple many-to-one genotype to phenotype map. The result is a distribution of the same form as the distribution of the frequencies of genotypes, but where the mutation term for each phenotype is weighted by the degeneracy of that phenotype; this shows that phenotypes of high degeneracy will tend to be favoured as on average there are more mutational paths into them. We show explicitly in the case of two phenotypes a phase-transition in frequency to the more degenerate phenotype as population size is reduced and/or mutation rate increased. This theory represent a finite population sized description of the analogous phenomenon to survival of the flattest of quasispecies theory and a strong mutation description of survival of the frequent found in the monomorphic weak mutation regime.

TRANSFORMING DISTRIBUTION FROM GENOTYPE TO PHENOTYPE

We assume a genotype space denoted by a vector \mathbf{g} , where $g_i = \{\sigma_k\}$, where σ_k represents possible symbols at each site i of the genome, from an alphabet of size \mathcal{A} . We assume each genotype has fitness $f_{\mathbf{g}}$ and that the mutation rate from \mathbf{g} to \mathbf{g}' is $\mu_{\mathbf{g} \rightarrow \mathbf{g}'} = \mu_0 \delta(\rho(\mathbf{g}, \mathbf{g}') - 1)$, where

$\rho(\mathbf{g}, \mathbf{g}')$ is the Hamming distance between sequences \mathbf{g} and \mathbf{g}' . The rate of mutations into any given state is then the same as any other state, $\mu_{\mathbf{g}} = \mu = \sum_{\mathbf{g}'} \mu_{\mathbf{g}\mathbf{g}'} = L(\mathcal{A}-1)\mu_0$. In this case, the equilibrium distribution of the frequency of genotypes, $x_{\mathbf{g}}$:

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{\mathbf{g}} e^{2Nf_{\mathbf{g}}x_{\mathbf{g}}} x_{\mathbf{g}}^{2N\mu-1} \delta(1 - \sum_{\mathbf{g}} x_{\mathbf{g}}) \quad (1)$$

which we derive in the Appendix, as to the author's knowledge this has not been done explicitly in the literature⁸. Here N is the effective population size and the vector \mathbf{x} is

$$p(\mathbf{z}) = \frac{1}{Z} \int d\mathbf{x} \prod_{\mathbf{g}} e^{2Nf_{\mathbf{g}}x_{\mathbf{g}}} x_{\mathbf{g}}^{2N\mu-1} \delta(1 - \sum_{\mathbf{g}} x_{\mathbf{g}}) \prod_{\xi} \delta(z_{\xi} - \sum_{\mathbf{g} \in \xi} x_{\mathbf{g}}) \quad (2)$$

$\sum_{\mathbf{g}} x_{\mathbf{g}} = \sum_{\xi} z_{\xi}$ and so the first delta function constraint enforces $\sum_{\xi} z_{\xi} = 1$. The product over genotypes can be decomposed into a product over phenotypes and a product over genotypes which map to the same phenotype, where these genotypes have the same fitness, by definition, f_{ξ} :

$$p(\mathbf{z}) = \frac{1}{Z} \delta(1 - \sum_{\xi} z_{\xi}) \prod_{\xi} e^{2Nf_{\xi}z_{\xi}} \int \prod_{\mathbf{g} \in \xi} dx_{\mathbf{g}} x_{\mathbf{g}}^{2N\mu-1} \delta(z_{\xi} - \sum_{\mathbf{g} \in \xi} x_{\mathbf{g}}) \quad (3)$$

$$= \frac{1}{Z} \delta(1 - \sum_{\xi} z_{\xi}) \prod_{\xi} e^{2Nf_{\xi}z_{\xi}} \int \prod_{\mathbf{g} \in \xi} dx_{\mathbf{g}} x_{\mathbf{g}}^{2N\mu-1} \left(z_{\xi} - \sum_{\mathbf{g} \in \xi} x_{\mathbf{g}} \right)^{2N\mu-1} \quad (4)$$

Integrating over the frequency of a single genotype of each phenotype, we are left to evaluate a multidimensional integral over the remaining genotypes, which are coupled due to the constraint that the sum over their frequencies should be z_{ξ} . To perform the integral we modify the transformation in⁹, which transforms the unit simplex to the unit cube; here we will transform the simplex over all genotypes that belong to a given phenotype constrained to sum to frequency z_{ξ} to the unit cube over transformed genotype frequencies u_i (switching to a linear index i corresponding to the i^{th} genotype \mathbf{g}_i):

$$u_i = \frac{x_i}{z_{\xi} - \sum_{j < i} x_j} \quad (5)$$

The inverse of this transformation is

$$x_i = z_{\xi} u_i \prod_{j < i} (1 - u_j). \quad (6)$$

Making the change of variable in the integral, and using the fact that

$$z_{\xi} - \sum_i x_i = \frac{x_{\Omega_{\xi}}}{u_{\Omega_{\xi}}} = z_{\xi} \prod_{j < \Omega_{\xi}} (1 - u_j)$$

we have:

$$p(\mathbf{z}) = \frac{1}{Z} \delta(1 - \sum_{\xi} z_{\xi}) \prod_{\xi} e^{2Nf_{\xi}z_{\xi}} \int \left| \frac{\partial \mathbf{x}}{\partial \mathbf{u}} \right| \prod_{\mathbf{g}_i \in \xi} du_i \left(z_{\xi} u_i \prod_{j < i} (1 - u_j) \right)^{2N\mu-1} \left(z_{\xi} \prod_{j < \Omega_{\xi}} (1 - u_j) \right)^{2N\mu-1} \quad (7)$$

The determinant of the Jacobian can be simply evaluated since $|\partial x_i / \partial u_j| = 0$ for $j > i$ and so the determinant is the product of the diagonal elements:

$$\left| \frac{\partial \mathbf{x}}{\partial \mathbf{u}} \right| = \prod_i^{\Omega_\xi - 1} \frac{\partial x_i}{\partial u_i} = \prod_i^{\Omega_\xi - 1} \left[z_\xi \prod_{j < i} (1 - u_j) \right] = z_\xi^{\Omega_\xi - 1} \prod_i^{\Omega_\xi - 1} \left[\prod_{j < i} (1 - u_j) \right] \quad (8)$$

Plugging this into above we get:

$$p(\mathbf{z}) = \frac{1}{Z} \delta\left(1 - \sum_\xi z_\xi\right) \prod_\xi^n e^{2N f_\xi z_\xi} z_\xi^{\Omega_\xi - 1} \left(z_\xi^{2N\mu - 1}\right)^{\Omega_\xi - 1} z_\xi^{2N\mu - 1} \int \prod_{\mathbf{g}_i \in \xi}^{\Omega_\xi - 1} \left[du_i [u_i(1 - u_i)]^{2N\mu - 1} \prod_{j < i} (1 - u_j)^{2N\mu} \right] \quad (9)$$

$$= \frac{1}{Z'} \delta\left(1 - \sum_\xi z_\xi\right) \prod_\xi^n e^{2N f_\xi z_\xi} z_\xi^{2N\mu\Omega_\xi - 1} \quad (10)$$

$$(11)$$

where Z' is the normalisation factor. The key result here is that the degeneracy of phenotypes enhances the effective mutation rate into that phenotype giving a bias to increase the frequency of that phenotype.

EXAMPLE: TWO PHENOTYPES

Let's assume there are two phenotypes with log-fitness f_1 and f_2 , with selection coefficient $s = f_1 - f_2$, degeneracies Ω_1 , and Ω_2 and a base-pair mutation rate μ . If the frequency of phenotype 1 is denoted z , and phenotype 2 $1 - z$, the probability density is given by Eqn.9:

$$p(z) = \frac{1}{Z} e^{2Nsz} z^{2N\mu\Omega_1 - 1} (1 - z)^{2N\mu\Omega_2 - 1} \quad (12)$$

where $Z = \frac{\Gamma(2N\mu\Omega_1)\Gamma(2N\mu\Omega_2)}{\Gamma(2N\mu(\Omega_1 + \Omega_2))} {}_1F_1(2N\mu\Omega_1; 2N\mu(\Omega_1 + \Omega_2); 2Ns)$. As shown in Fig.1, if we allow phenotype 1 to be advantageous with $s > 0$, there is a shift from phenotype 1 being at high frequency, when the degeneracies are equal ($\Omega_1 = \Omega_2$) to favouring phenotype 2 (reduction in frequency z), when there is a strong bias in the genotype-phenotype map towards phenotype 2 ($\Omega_2 \gg \Omega_1$). For large population sizes (a) we are in the strong mutation regime, where $2N\mu\Omega_1 \gg 1$ and $2N\mu\Omega_2 \gg 1$, and we see there is a mutation-selection balance in favour of the advantageous phenotype, when degeneracies are equal, and this equilibrium moves to smaller frequencies as the ratio of the degeneracies Ω_2/Ω_1 increases; this is an analogous finite N description of population delocalisation as found for infinite N deterministic quasispecies modelling of populations^{6,7}, except here we capture the broad fluctuations around the mutation-selection balance equilibrium. On the other hand when the effect of mutations is weak, $2N\mu\Omega_1 \ll 1$ and $2N\mu\Omega_2 \ll 1$, as shown in (b), we see distributions characteristic of the mostly monomorphic composition of the population at any given time, where distributions are condensed at $z = 0$ and $z = 1$; nonetheless we see that when the ratio of the degeneracies (Ω_2/Ω_1) becomes large, the distributions shift from a larger density near $z = 1$ to one

where a larger density at $z = 0$. This is the correct polymorphic extension of the populations in the monomorphic weak mutation regime, which has recently seen attention using such concepts as free fitness and sequence entropy^{2-4,10-12}.

We can also examine the effect of changing the mutation rate or population size in a two phenotype system when there is a large bias in degeneracy. This is most effectively probed by calculating the mean frequency of the phenotype $\langle z \rangle = \int dz zp(z)$:

$$\langle z \rangle = \frac{1}{Z} \int_0^1 dz z e^{2Nsz} z^{2N\mu\Omega_1 - 1} (1 - z)^{2N\mu\Omega_2 - 1} \quad (13)$$

$$= \frac{1}{2N\bar{Z}} \frac{d}{ds} \int_0^1 e^{2Nsz} z^{2N\mu\Omega_1 - 1} (1 - z)^{2N\mu\Omega_2 - 1} \quad (14)$$

$$= \frac{1}{2N\bar{Z}} \frac{dZ}{ds} \quad (15)$$

$$= \frac{\Omega_1}{\Omega_1 + \Omega_2} \frac{{}_1F_1(2N\mu\Omega_1 + 1; 2N\mu(\Omega_1 + \Omega_2) + 1; 2Ns)}{{}_1F_1(2N\mu\Omega_1; 2N\mu(\Omega_1 + \Omega_2); 2Ns)} \quad (16)$$

where we have used the fact that $\frac{d}{dx} {}_1F_1(a; b; x) = \frac{a}{b} {}_1F_1(a + 1; b + 1; x)$. We expect that in the limit that $2N\mu \rightarrow 0$ to recover the monomorphic limit^{2-4,10}, where the probability of each phenotype is given by $p_1 = \frac{e^{2Ns\Omega_1}}{e^{2Ns\Omega_1} + e^{2Ns\Omega_2}}$ and $p_2 = \frac{e^{2Ns\Omega_2}}{e^{2Ns\Omega_1} + e^{2Ns\Omega_2}}$, such that $\langle z \rangle = p_1 \times 1 + p_2 \times 0 = \frac{e^{2Ns\Omega_1}}{e^{2Ns\Omega_1} + e^{2Ns\Omega_2}}$. Using the series definition of the confluent hypergeometric function it is straightforward to show that

$$\lim_{2N\mu \rightarrow 0} \{{}_1F_1(2N\mu\Omega_1 + 1; 2N\mu(\Omega_1 + \Omega_2) + 1; 2Ns)\} = e^{2Ns} \quad (17)$$

$$\lim_{2N\mu \rightarrow 0} \{{}_1F_1(2N\mu\Omega_1; 2N\mu(\Omega_1 + \Omega_2); 2Ns)\} = 1 + \frac{\Omega_1(e^{2Ns} - 1)}{\Omega_1 + \Omega_2}, \quad (18)$$

which leads to the following expression for the limit of the mean frequency:

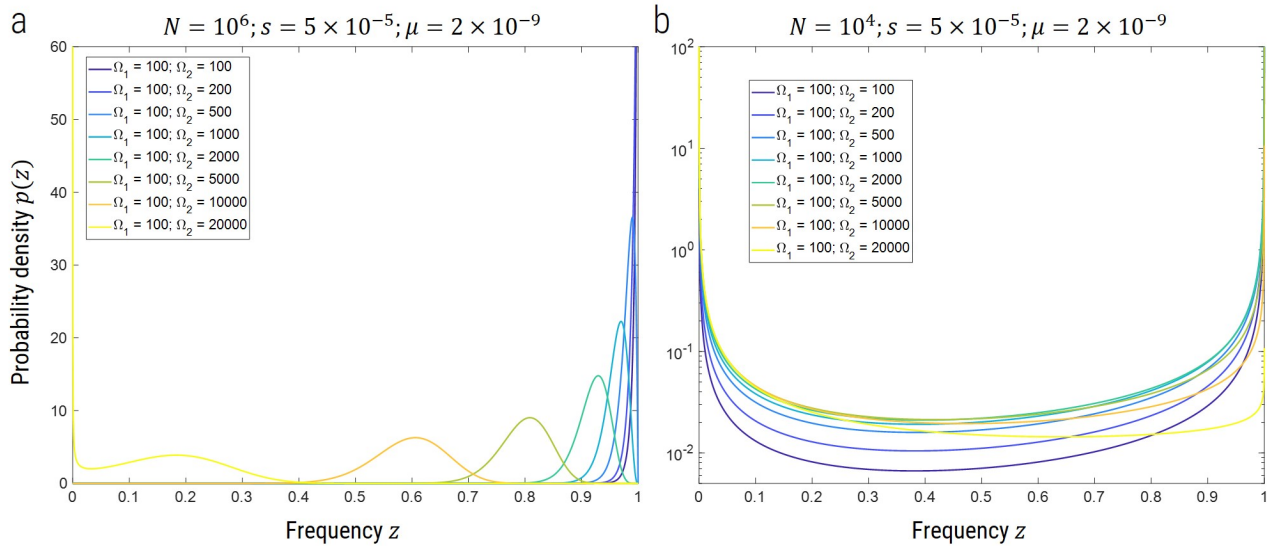


FIG. 1. Plot of the phenotypic frequency distribution for different values for $\mu = 2 \times 10^{-9}$, $s = 5 \times 10^{-5}$ and a) $N = 10^6$ and b) $N = 10^4$, for various values of the degeneracy of each phenotype as shown in the legend. N.B. the probability density is plotted on a linear scale in a) and log scale in b) for clarity.

$$\lim_{2N\mu \rightarrow 0} = \frac{\Omega_1 e^{2Ns}}{\Omega_1 e^{2Ns} + \Omega_2}, \quad (19)$$

which agrees with the monomorphic expectation.

We can also take the limit $2Ns \rightarrow 0$, keeping $2N\mu$ finite, which is simply evaluated as $\lim_{x \rightarrow 0} \{ {}_1F_1(a; b; x) \} = 1$, giving

$$\lim_{2Ns \rightarrow 0} = \frac{\Omega_1}{\Omega_1 + \Omega_2}. \quad (20)$$

So whether in the monomorphic ($2N\mu \ll 1$) or polymorphic limit ($2N\mu > 1$), as selection becomes very weak, we get the purely neutral result that the average frequency is determined solely by the relative degeneracies of each phenotype.

In Fig.2, we investigate the mean frequency $\langle z \rangle$ for $\Omega_1 = 100$ and $\Omega_2 = 10000$ and $s = 5 \times 10^{-5}$, for various values of N and μ . We see there is a strong delocalisation transition for both μ and N ; as has been essentially described previously for deterministic quasispecies in what is known as the “survival of the flattest”^{6,7}, for an increasing mutation rate, we find the less advantageous, but more genotypically degenerate phenotype is favoured. However, here this calculation also shows that concurrently decreasing the effective population size increases progress towards this delocalisation transition, again favouring the more degenerate phenotype.

In Fig.3, we have an equivalent plot for the case when the degeneracies of each phenotype is the same. We see that we have a transition to the neutral state for increasing μ or

decreasing N , giving an equal likelihood of each phenotype, signified by $\langle z \rangle = 1/2$.

DISCUSSION & CONCLUSIONS

In this paper, we have derived the equilibrium distribution of the frequency of phenotypes, assuming a general many-to-one genotype phenotype map and that the mutation rates between nearest neighbour genotypes is uniform. The result shows that the equilibrium distribution is of the same form as genotypes but the mutation terms are weighted by the degeneracy of each phenotype. This gives rise to a bias towards phenotypes of higher degeneracy as the mutation rate is increased and/or the population size decreased, as we show explicitly for the two phenotype case. This calculation generalises the equilibrium distribution of phenotypes in the monomorphic regime to the polymorphic, describing the analogous effect of the increasing effect of sequence entropy (log degeneracy of phenotypes), as population size is decreased. However, here we note that there is no obvious way to express the equilibrium frequency distribution in terms of an analogous quantity such as sequence entropy. Iwasa³ and Barton & Coe¹⁰ describe a free fitness of the frequencies of genotypes, in terms the mean fitness, a genetic diversity and an entropy term. It is possible here to construct an analogous free fitness of the frequency of phenotypes, by separating out the terms that depend only on fitness, mutation and drift, however, it is not clear that this is meaningful; here the mutation term is weighted by the degeneracy, giving an effective mutation driven delocalisation to more frequent phenotypes, yet this same degeneracy in the small population size monomorphic limit gives rise to

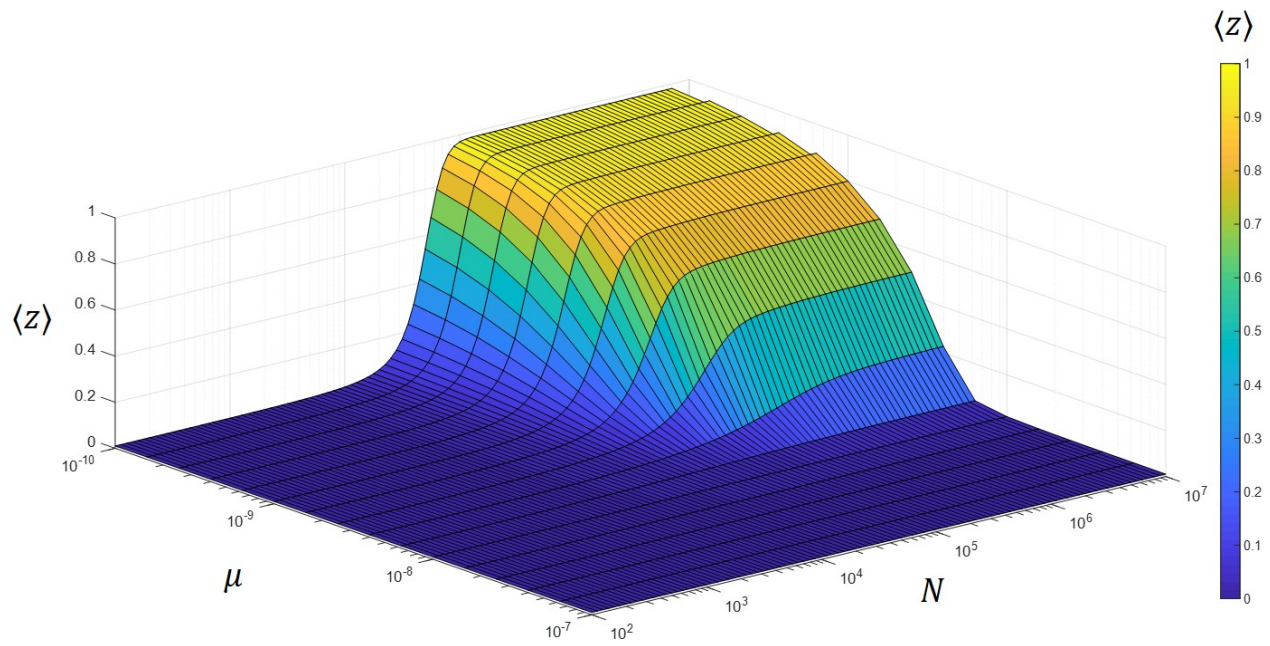


FIG. 2. Plot of the mean frequency of phenotype 1, $\langle z \rangle$, as function of effective population size N and mutation rate μ for a selection coefficient $s = 5 \times 10^{-5}$ (favouring phenotype 1) and degeneracies $\Omega_1 = 100$ and $\Omega_2 = 10000$.

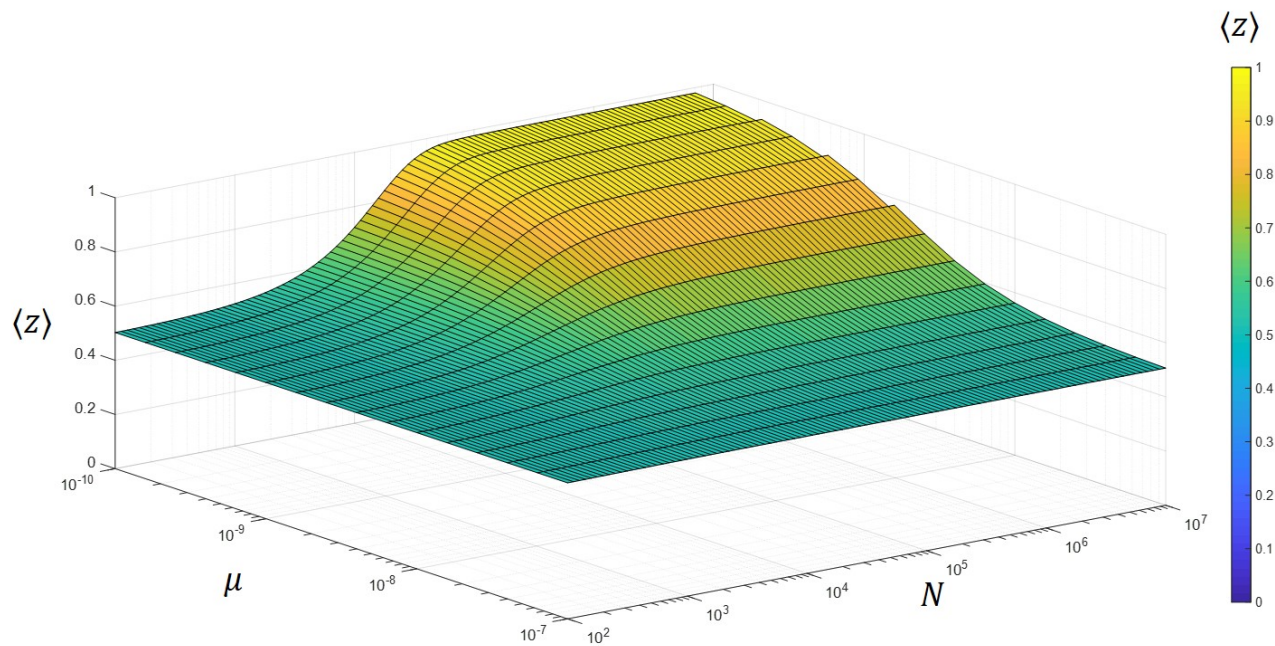


FIG. 3. Plot of the mean frequency of phenotype 1, $\langle z \rangle$, as function of effective population size N and mutation rate μ for a selection coefficient $s = 5 \times 10^{-5}$ (favouring phenotype 1) and degeneracies $\Omega_1 = 10000$ and $\Omega_2 = 10000$.

the (Boltzmann) sequence entropy of phenotypes.

This calculation is also a finite population size description

of the quasispecies limit, which are deterministic infinite population size models⁵. Quasispecies models predict the

phenomenon of *survival of the flattest*, which is a delocalisation transition due to higher robustness or smaller 'local curvature' being favoured at larger mutation rates^{6,7}; here we make no explicit statement about robustness, and we see an analogous phenomenon arises irrespective of how different genotypes of the same phenotype are connected; the difference between this and the quasispecies is that here we have an equilibrium calculation, which assumes a certain ergodicity or accessibility of a representative number of states for each phenotype. Whether this equilibrium is reached on relevant evolutionary timescales is an open question; as discussed in¹³, certain phenotypes may be more likely to arrive as they arise more frequently in the local neighbourhood. On the other hand there are broad reasons to expect the structure of genotype-phenotype maps to be ergodic¹⁴; in particular, simulations of the genotype-phenotype map for spatial patterning in development studied in¹¹ was found to be ergodic, despite simulation times which could not exhaustively search the whole genotype space.

We also present for the first time, to the author's knowledge, an explicit derivation of Wright's multi-allele frequency distribution⁸. The derivation makes clear the "curl-free" or "potential" assumption that gives rise to the stationary or equilibrium distribution of allele frequencies. In other words there cannot be any circulating fluxes of probability in the equilibrium state; this arises due to the uniform mutation assumption, and we expect that for arbitrary mutation structure the curl-free assumption will not be satis-

fied. It is an open question as to which forms of mutation structure can give rise to potential solutions of the form of Wright's stationary distribution. We believe the results of this paper will likely be robust to such considerations in the limit of large degeneracies of each phenotype, where the product of degeneracy and mutation rate into a phenotype becomes some effective average mutation rate over microscopic mutation rates.

To conclude these results provide a quantitative theory to calculate the equilibrium distribution of the frequency of phenotypes for a general genotype-phenotype map. As such it provides a bridge between recent results showing the importance of degeneracy of phenotypes in the weak mutation, monomorphic, regime and infinite population size, deterministic, quasispecies calculations that demonstrate delocalisation phenomena due to increasing mutation rates.

ACKNOWLEDGEMENTS

I thank Ard Louis for useful discussions and for initially posing the question.

APPENDIX

The Fokker-Planck equation describing the stochastic dynamics of the gene frequencies of multiple alleles is given by

$$\frac{\partial p(\mathbf{x}, t)}{\partial t} = - \sum_{i=1}^{n-1} \frac{\partial}{\partial x_i} (A_i(\mathbf{x})p(\mathbf{x}, t)) + \frac{1}{2} \sum_{i=1}^{n-1} \sum_{j=1}^{n-1} \frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} (B_{ij}(\mathbf{x})p(\mathbf{x}, t)) \quad (21)$$

where n is the number of alleles, and the mean change in allele frequency (convective force) on the i^{th} allele is

$$A_i(\mathbf{x}) = (f_i - \bar{f}(\mathbf{x}))x_i + \mu(1 - nx_i) \quad (22)$$

and the co-variance of the change in allele frequencies of the i^{th} and j^{th} allele (effective diffusion matrix) is

$$B_{ij}(\mathbf{x}) = \frac{1}{N}(\delta_{ij} - x_i)x_j \quad (23)$$

The equilibrium solution $p^*(\mathbf{x}, t)$ can be found by setting the flux $\mathbf{J}(\mathbf{x}) = 0$, where the flux is defined by the continuity equation to be $\partial_t p(\mathbf{x}, t) = -\nabla \cdot \mathbf{J}$:

$$J_i(\mathbf{x}) = A_i(\mathbf{x})p^* - \frac{1}{2} \sum_{j=1}^{n-1} \frac{\partial}{\partial x_j} (B_{ij}(\mathbf{x})p^*) \quad (24)$$

$$= \left(A_i(\mathbf{x}) - \frac{1}{2} \sum_{j=1}^{n-1} \frac{\partial B_{ij}}{\partial x_j} \right) p^* - \frac{1}{2} \sum_{j=1}^{n-1} B_{ij}(\mathbf{x}) \frac{\partial p^*}{\partial x_j}. \quad (25)$$

Multiplying through by the inverse B_{ij}^{-1} and summing it is simple to show that

$$\frac{\partial p^*}{\partial x_i} = 2 \sum_{j=1}^{n-1} B_{ij}^{-1} \left(A_i(\mathbf{x}) - \frac{1}{2} \sum_{j=1}^{n-1} \frac{\partial B_{ij}}{\partial x_j} \right) p^* \quad (26)$$

$$= \psi_i(\mathbf{x})p^* \quad (27)$$

Now if the vector field is *conservative*, i.e. it is the gradient of a scalar function, $\psi = \nabla \Psi(\mathbf{x})$, then the solution to

Eqn.26 can be found using the standard integrating factor method:

$$p^*(\mathbf{x}) = \frac{1}{Z} e^{\int_{\sigma} \psi(\mathbf{x}') \cdot d\mathbf{x}'} = \frac{1}{Z} e^{\Psi(\mathbf{x})}. \quad (28)$$

The vector field ψ is conservative if it is free from rotation, which in 3 dimensions or less means it is curl free. In higher dimensions an equivalent condition is:

$$\frac{\partial \psi_i}{\partial x_j} = \frac{\partial \psi_j}{\partial x_i} \quad (29)$$

which we loosely refer to as the “curl-free” condition. We can evaluate $\psi_i(\mathbf{x})$, by using the fact that the inverse of B_{ij} is given by¹⁵:

$$B_{ij}^{-1} = N \left(\frac{\delta_{ij}}{x_i} + \frac{1}{x_n} \right) \quad (30)$$

where $x_n = 1 - \sum_{i=1}^{n-1} x_i$ to give

$$\psi_i(\mathbf{x}) = 2N(f_i - f_n) + (2N\mu - 1) \left(\frac{1}{x_i} - \frac{1}{x_n} \right). \quad (31)$$

Evaluating the partial derivative wrt x_j we find:

$$\frac{\partial \psi_i}{\partial x_j} = -\frac{2N\mu}{x_n^2}. \quad (32)$$

As this does not depend on x_i or x_j then this satisfies Eqn.29 and ψ is a curl-free vector field. Note that this is essentially a restatement of the fact that in order to find an equilibrium solution we are assuming detailed balance is obeyed. To find $\Psi(\mathbf{x})$ we can integrate by inspection to get:

$$\Psi(\mathbf{x}) = 2N \sum_{i=1}^n f_i x_i + (2N\mu - 1) \sum_{i=1}^n \ln(x_i), \quad (33)$$

where note that the summations run to the index $i = n$. Plugging into Eqn.28, we find our final result

$$p^*(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^n e^{2N f_i x_i} x_i^{2N\mu-1}. \quad (34)$$

To the author’s knowledge this is the first time this derivation has appeared explicitly in the literature.

REFERENCES

- ¹R. Neher and B. Shraiman, “Statistical genetics and evolution of quantitative traits,” *Reviews of Modern Physics* (2011).
- ²B. S. Khatri and R. A. Goldstein, “A coarse-grained biophysical model of sequence evolution and the population size dependence of the speciation rate,” *Journal of theoretical biology* **378**, 56–64 (2015).
- ³Y. Iwasa, “Free fitness that always increases in evolution,” *Journal of Theoretical Biology* **135**, 265 – 281 (1988).
- ⁴G. Sella and A. E. Hirsh, “The application of statistical physics to evolutionary biology.” *Proc Natl Acad Sci U S A* **102**, 9541–9546 (2005).
- ⁵M. Eigen, “Selforganization of matter and the evolution of biological macromolecules.” *Naturwissenschaften* **58**, 465–523 (1971).
- ⁶M. Eigen, J. McCaskill, and P. Schuster, “Molecular quasi-species,” *The Journal of Physical Chemistry* **92**, 6881–6891 (1988).
- ⁷C. O. Wilke, “Quasispecies theory in the context of population genetics,” *BMC Evolutionary Biology* **5**, 1–8 (2005).
- ⁸S. Wright, “Genetics, Palaeontology and Evolution,” (Princeton University Press, 1949) Chap. Adaptation and selection, p. 383.
- ⁹G. J. Baxter, R. A. Blythe, and A. J. McKane, “Exact solution of the multi-allelic diffusion model.” *Math Biosci* **209**, 124–170 (2007).
- ¹⁰N. H. Barton and J. B. Coe, “On the application of statistical physics to evolutionary biology.” *J Theor Biol* **259**, 317–324 (2009).
- ¹¹B. S. Khatri, T. C. B. McLeish, and R. P. Sear, “Statistical mechanics of convergent evolution in spatial patterning.” *Proc Natl Acad Sci U S A* **106**, 9564–9569 (2009).
- ¹²B. S. Khatri and R. A. Goldstein, “Simple biophysical model predicts faster accumulation of hybrid incompatibilities in small populations under stabilizing selection,” *Genetics* **201**, 1525–1537 (2015).
- ¹³S. Schaper and A. A. Louis, “The arrival of the frequent: how bias in genotype-phenotype maps can steer populations to local optima,” *PloS one* **9**, e86635 (2014).
- ¹⁴T. C. B. McLeish, “Are there ergodic limits to evolution? ergodic exploration of genome space and convergence,” *Interface Focus* **5**, 20150041 (2015).
- ¹⁵P. Antonelli and C. Strobeck, “The geometry of random drift i. stochastic distance and diffusion,” *Advances in Applied Probability* **9**, 238–249 (1977).