

FILE S6: AN ALTERNATIVE APPROACH TO NEUTRAL DIVERSITY

Instead of calculating the distribution of neutral heterozygosity by first computing the distribution of real times, we could alternatively incorporate neutral mutations directly into the sum of ancestral paths framework. This completely bypasses the correspondence with real coalescence times. To do this, we characterize ancestral paths not only by the ordering of deleterious mutation and coalescence events, but also by the ordering of neutral mutations. This means that if we sample two individuals  $A$  and  $B$ , there are five types of events that can happen in their ancestral paths: a deleterious mutation (DM) in  $A$  or in  $B$ , a neutral mutation (NM) in either  $A$  or in  $B$ , and or a coalescence (C) event (if  $A$  and  $B$  are currently in the same class).

We now imagine that we sample two individuals from classes  $k$  and  $k'$ , and that they coalesce in class  $k - \ell$ . Our goal is to calculate the probability distribution of  $\pi_n$  given  $k$ ,  $k'$ , and  $\ell$ ,  $\rho(\pi_n|k, k', \ell)$ . We will find it helpful to divide the five types of events that can occur into two classes: neutral mutations on the one hand, and deleterious mutations or coalescence (which we call “steps”) on the other. We begin by computing the probability that a given number of NMs occur before the next DM or C events (i.e. the number of neutral mutations that occur at this “step”). We have

$$P(\text{a NMs, then DM in } k' \text{ or } k'|k', k) = \left( \frac{\frac{2U_n}{s}}{k' + k + \frac{2U_n}{s}} \right)^a \frac{k + k'}{k' + k + \frac{2U_n}{s}}, \quad (\text{S.74})$$

where we have made our usual assumption that  $Nh_k s k \gg 1$ , allowing us to neglect the rates of coalescence events (when  $k = k'$ ) in writing this expression.

This probability only depends on the sum of the current classes the individuals are in. At each subsequent step, regardless of the path taken, this sum of the classes will decrease by one. Therefore, the probability that  $a_i$  neutral mutations occur at step  $i$  is independent of the path taken. This observation allows us to calculate the probability that a given total number of neutral mutations have occurred since coalescence. We first calculate the probability that a given number of neutral mutations have occurred since the first deleterious mutation out of the  $k - \ell$  class. We will add in the additional neutral mutations once in the  $k - \ell$  class at the end.

In order for  $\pi_n$  neutral mutations to have occurred since the first deleterious mutation out of class  $k - \ell$ , we require that  $a_0$  mutations occurred at the first step,  $a_1$  mutations occurred at the second step, and so on, such that  $a_0 + a_1 + \dots + a_{k'-k+2\ell-1} = \pi_n$ . This gives

$$\rho(\pi_n = X|k', k, \ell) = \frac{\frac{(k'+k)!}{(2k-2\ell)!}}{\frac{(2U_n/s+k'+k)!}{(2U_n/s+2k-2\ell)!}} \sum_{|\vec{a}|=X} \left( \frac{2U_n/s}{2U_n/s+k+k'} \right)^{a_0} \cdots \left( \frac{2U_n/s}{2U_n/s+2k-2\ell+1} \right)^{a_{k'-k+2\ell-1}}. \quad (\text{S.75})$$

We can define  $x \equiv 2U_n/s + k + k'$ , recognize  $\pi_d = k' - k + 2\ell$ , and relabel the  $a_i$  as

$$a_0 \rightarrow X - b_0, \quad a_1 \rightarrow b_0 - b_1, \quad \dots \quad a_{\pi_d-2} \rightarrow b_{\pi_d-3} - b_{\pi_d-2}, \quad a_{\pi_d-1} \rightarrow b_{\pi_d-2}. \quad (\text{S.76})$$

This gives

$$\begin{aligned} \rho(\pi_n = X | k', k, \ell) &= \frac{\binom{k'+k}{\pi_d}}{\binom{2U_n+k'+k}{\frac{2U_n}{s} \pi_d}} \left(\frac{2U_n}{s}\right)^X \left(\frac{1}{x}\right)^X \sum_{b_0=0}^X \left(\frac{x}{x-1}\right)^{b_0} \\ &\quad \sum_{b_1=0}^{b_0} \left(\frac{x-1}{x-2}\right)^{b_1} \dots \sum_{b_{\pi_d-2}=0}^{b_{\pi_d-3}} \left(\frac{x-\pi_d+2}{x-\pi_d+1}\right)^{b_{\pi_d-2}}. \end{aligned} \quad (\text{S.77})$$

To simplify this expression, it is helpful to define a function  $\mathbf{f}$  such that:

$$\begin{aligned} \mathbf{f}(A, B) &\equiv \left(\frac{1}{x}\right)^X \sum_{b_0=0}^X \left(\frac{x}{x-1}\right)^{b_0} \\ &\quad \sum_{b_1=0}^{b_0} \left(\frac{x-1}{x-2}\right)^{b_1} \dots \sum_{b_{A-1}=0}^X \left(\frac{x-A+1}{x-A}\right)^{b_0} \sum_{b_A=0}^{b_{A-1}} \left(\frac{x-A}{x-B}\right)^{b_A} \end{aligned} \quad (\text{S.78})$$

In other words,  $\mathbf{f}(A, B)$  is a set of  $A$  nested sums, each of the same form, except for the final sum, which can have a different denominator. Using this definition, we have

$$P(\pi_n = X | k', k, \ell) = \frac{\binom{k'+k}{\pi_d}}{\binom{2U_n+k'+k}{\frac{2U_n}{s} \pi_d}} \left(\frac{2U_n}{s}\right)^X \mathbf{f}(\pi_d - 2, \pi_d - 1). \quad (\text{S.79})$$

The virtue of this definition is that this sum can be solved recursively. We have

$$\sum_{b_A=0}^{b_{A-1}} \left(\frac{x-A}{x-B}\right)^{b_A} = \frac{x-B}{A-B} - \frac{x-A}{A-B} \left(\frac{x-A}{x-B}\right)^{b_{A-1}}. \quad (\text{S.80})$$

Therefore we have

$$\mathbf{f}(A, B) = \frac{x-A}{B-A} \mathbf{f}(A-1, B) - \frac{x-B}{B-A} \mathbf{f}(A-1, A). \quad (\text{S.81})$$

Repeatedly inserting this result yields:

$$\begin{aligned} \mathbf{f}(A, A+1) &\rightarrow \frac{(x-A)(x-A-1)}{1} \left( \frac{\mathbf{f}(A-1, A+1)}{x-A-1} - \frac{\mathbf{f}(A-1, A)}{x-A} \right) \\ \mathbf{f}(A, A+1) &\rightarrow \frac{(x-A+1)(x-A)(x-A-1)}{2} \left[ \frac{\mathbf{f}(A-2, A+1)}{x-A-1} - \frac{2\mathbf{f}(A-2, A)}{x-A} + \frac{\mathbf{f}(A-2, A-1)}{x-A+1} \right] \\ &\vdots \\ \mathbf{f}(A, A+1) &\rightarrow (m+1) \binom{x-A-1+m}{m+1} \sum_{i=0}^m \frac{(-1)^{i+m}}{x-A-1+i} \binom{m}{i} \mathbf{f}(A-m, A+1-i). \end{aligned} \quad (\text{S.82})$$

Note that  $\mathbf{f}(-1, B) = 1/B^X$ , since there are no more sums to compute. Thus, for  $m = A + 1$  we have

$$\mathbf{f}(A, A + 1) = (A + 2) \binom{x}{A + 2} \sum_{i=0}^{A+1} \frac{(-1)^{i+A+1}}{(x - A - 1 + i)^{X+1}} \binom{A + 1}{i}. \quad (\text{S.83})$$

Relabeling the sum and taking  $A = \pi_d - 2$ , we have

$$\mathbf{f}(\pi_d - 2, \pi_d - 1) = \pi_d \binom{x}{\pi_d} \sum_{i=0}^{\pi_d-1} \frac{(-1)^i}{(x - i)^{X+1}} \binom{\pi_d - 1}{i}. \quad (\text{S.84})$$

We can now substitute these results into our expression for  $\pi_n$ , to find

$$\rho_1(\pi_n = X | k', k, \ell) = \pi_d \binom{k' + k}{\pi_d} \left( \frac{2U_n}{s} \right)^X \sum_{i=0}^{\pi_d-1} \frac{(-1)^i}{(2U_n/s + k + k' - i)^{X+1}} \binom{\pi_d - 1}{i} \quad (\text{S.85})$$

Note, however, that this is only the distribution of neutral mutations since the first deleterious mutation out of class  $k - l$ . It is also possible for neutral mutations to occur prior to the coalescence event. Adding in this factor, we find

$$\begin{aligned} \rho(\pi_n = X | k', k, \ell) &= \pi_d \binom{k' + k}{\pi_d} \sum_{i=0}^{\pi_d-1} (-1)^i \binom{\pi_d - 1}{i} \\ &\times \sum_{X=0}^{\pi_n} \frac{(2U_n/s)^X}{(2U_n/s + k + k' - i)^{X+1}} \left( \frac{2N_{k-l}U_n}{1 + 2N_{k-l}U_n + 2N_{k-l}s(k-l)} \right)^{\pi_n - X}. \end{aligned} \quad (\text{S.86})$$

Rearranging this expression gives

$$\rho(\pi_n | k', k, \ell) = \sum_{i=0}^{\pi_d-1} \pi_d (-1)^{\pi_d - i - 1} \binom{\pi_d - 1}{i} \binom{k' + k}{\pi_d} \frac{B}{A - B} \left( \frac{(\frac{2U_n}{s})^{\pi_n}}{(\frac{2U_n}{s} + B)^{\pi_n + 1}} - \frac{(\frac{2U_n}{s})^{\pi_n}}{(\frac{2U_n}{s} + A)^{\pi_n + 1}} \right), \quad (\text{S.87})$$

where we have defined

$$A = k' + k - i, \quad B = 2(k - \ell) + \frac{1}{Nsh_{k-l}}, \quad (\text{S.88})$$

identical to our earlier result.