

Random Field Model Reveals Structure of the Protein Recombinational Landscape - Text S1

Philip A. Romero¹, Frances H. Arnold^{1,*},

¹ Division of Chemistry and Chemical Engineering, California Institute of Technology, Pasadena, CA, USA

* E-mail: frances@cheme.caltech.edu

Expected value of the library variance

Consider a recombination library L , generated by recombining b sequence fragments from p parental sequences. From the definition of the library mean M_L (main text equation 8), the expected value of the library mean within the random field is

$$E[M_L] = \frac{1}{p^b} \sum_{\mathbf{c} \in L} E[\mathcal{E}_{\mathbf{c}}] \quad (1)$$

and the variance of the library mean is

$$\text{Var}[M_L] = \frac{1}{p^{2b}} \sum_{\mathbf{c1} \in L} \sum_{\mathbf{c2} \in L} \text{Cov}[\mathcal{E}_{\mathbf{c1}}, \mathcal{E}_{\mathbf{c2}}], \quad (2)$$

where the expected value of the random field $E[\mathcal{E}_{\mathbf{c}}]$ is defined in main text equation 6, and the covariance within the random field $\text{Cov}[\mathcal{E}_{\mathbf{c1}}, \mathcal{E}_{\mathbf{c2}}]$ is defined in main text equation 7.

Similarly, the expected value of the library variance V_L (main text equation 9) is given by

$$E[V_L] = \frac{1}{p^b} \sum_{\mathbf{c} \in L} E[(\mathcal{E}_{\mathbf{c}} - M_L)^2] \quad (3)$$

which can be expanded to

$$E[V_L] = \frac{1}{p^b} \sum_{\mathbf{c} \in L} [(E[\mathcal{E}_{\mathbf{c}}] - E[M_L])^2 + \text{Var}[\mathcal{E}_{\mathbf{c}}] + \text{Var}[M_L] - 2 \text{Cov}[\mathcal{E}_{\mathbf{c}}, M_L]] \quad (4)$$

where $\text{Var}[\mathcal{E}_{\mathbf{c}}] = \text{Cov}[\mathcal{E}_{\mathbf{c}}, \mathcal{E}_{\mathbf{c}}]$ and

$$\text{Cov}[\mathcal{E}_{\mathbf{c}}, M_L] = -E[\mathcal{E}_{\mathbf{c}}]E[M_L] + \frac{1}{p^b} \sum_{\mathbf{c2} \in L} (E[\mathcal{E}_{\mathbf{c}}]E[\mathcal{E}_{\mathbf{c2}}] + \text{Cov}[\mathcal{E}_{\mathbf{c}}, \mathcal{E}_{\mathbf{c2}}]). \quad (5)$$

From this, we can substitute equations 1, 2, and 5 into equation 4 to get an expression for the expected value of the library variance.

Additive component of a chimera's energy

An additive energy function can be defined by considering how individual mutations contribute to variation in the library. Depending on its structural context, a mutation's effect may be constant or varied throughout the library. A chimera's additive energy, which accounts for purely additive and averaged epistatic effects, is given by

$$E_{A,\mathbf{c}} = \sum_i b_{\mathbf{c},P}^i \varepsilon_P^i + \sum_i b_{\mathbf{c},N}^i \varepsilon_N^i, \quad (6)$$

where $b_{\mathbf{c},P}^i$ and $b_{\mathbf{c},N}^i$ specify how the energy terms ε_P^i and ε_N^i contribute to additive energy of chimera \mathbf{c} . The b 's are analogous to the a 's from main text equation 1. However their values are not binary but rather are determined by the average contribution that an interaction makes to the library. For an interaction i between positions p_1^i and p_2^i , $b_{\mathbf{c},P}^i$ (and equivalently $b_{\mathbf{c},N}^i$) is given by

$$b_{\mathbf{c},P}^i = \begin{cases} a_{\mathbf{c},P}^i & \text{if } p_1^i \text{ and } p_2^i \text{ intra-fragment,} \\ f(i|\mathbf{c}, p_1) + f(i|\mathbf{c}, p_2) - f(i) & \text{if } p_1^i \text{ and } p_2^i \text{ inter-fragment,} \end{cases} \quad (7)$$

where $f(i)$ is the frequency of interaction i within the entire library L , $f(i|\mathbf{c}, p_1)$ is the frequency of interaction i in the subset of the library that has the same residue at position p_1^i as chimera \mathbf{c} , and $f(i|\mathbf{c}, p_2)$ is the frequency of interaction i in the subset of the library that has the same residue at position p_2^i as chimera \mathbf{c} . This additive energy can be used to calculate a library's additive variance V_A .