

# Scaling properties of evolutionary paths in a biophysical model of protein adaptation

Michael Manhart<sup>1,‡</sup> and Alexandre V Morozov<sup>1,2</sup>

<sup>1</sup> Department of Physics and Astronomy, Rutgers University, Piscataway, NJ 08854, USA

<sup>2</sup> BioMaPS Institute for Quantitative Biology, Rutgers University, Piscataway, NJ 08854, USA

E-mail: morozov@physics.rutgers.edu

**Abstract.** The enormous size and complexity of genotypic sequence space frequently requires consideration of coarse-grained sequences in empirical models. We develop scaling relations to quantify the effect of this coarse-graining on properties of fitness landscapes and evolutionary paths. We first consider evolution on a simple Mount Fuji fitness landscape, focusing on how the length and predictability of evolutionary paths scale with the coarse-grained sequence length and number of alleles. We obtain simple scaling relations for both the weak- and strong-selection limits, with a non-trivial crossover regime at intermediate selection strengths. We apply these results to evolution on a biophysical fitness landscape designed to describe how proteins evolve new binding interactions while maintaining their folding stability. We combine numerical calculations for coarse-grained protein sequences with the scaling relations to obtain quantitative properties of the model for realistic binding interfaces and a full amino acid alphabet.

<sup>‡</sup> Current address: Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA 02138, USA

## 1. Introduction

The enormous size and dimensionality of genotypic sequence space are among the most salient features of molecular evolution. These features not only present technical challenges for experiments and computation, but raise major conceptual questions as well: how can populations efficiently find high-fitness states in such a large space? John Maynard Smith famously tackled this issue [1], arguing that positive selection acting on individual mutations was key to efficiently evolving functional protein sequences. However, this argument depends crucially on the structure of the fitness landscape and the underlying evolutionary dynamics. One expects a large population to ascend a steep and perfectly-smooth landscape quickly, while substantial landscape ruggedness or genetic drift will slow down adaptation.

The effect of ruggedness (due to epistatic interactions among genetic loci) on evolutionary paths has been a major focus of previous work. These studies have investigated both simple models of fitness landscapes, especially the uncorrelated random landscape [2–5] (also known as the “House of Cards” [6]) and the rough Mount Fuji model [5, 7, 8], as well as landscapes empirically measured in specific organisms [9, 10]. Populations in these studies are generally assumed to be under strong selection, so that evolutionary paths proceed strictly upward in fitness; thus a major goal is to determine the number and length of the accessible paths for different landscape topographies. More recent work has begun to consider the effect of population dynamics (e.g., clonal interference) on evolutionary predictability [11], a topic of central importance in evolutionary biology [12, 13].

In most cases the computational and experimental cost of analyzing empirical models has required simplified sequence spaces, especially binary sequences (indicating only the presence or absence of a mutation at each site) [3, 5, 8, 9], genomes or proteins with reduced lengths [14–17], and reduced sets of amino acids [16, 18] or protein structural components [19]. However, it is not clear how properties of landscapes and evolutionary paths change under these implicit coarse-graining schemes, which is essential for extending these models to more realistic biological systems. Specifically, we must determine how properties of a model scale with both the coarse-grained sequence length  $L$  and the coarse-grained number  $k$  of alleles at each site, the latter being important when multiple mutations at a single site are likely.

We first carry out this approach in a simple model of monomorphic populations undergoing substitutions on a smooth Mount Fuji landscape, showing how the scaling properties of the model depend crucially on the strength of selection relative to genetic drift. We then

consider evolution on a fitness landscape based on the biophysics of protein folding and binding, describing how proteins evolve new binding interactions while maintaining folding stability [18]. Using scaling relations, we extend numerical calculations of the model for coarse-grained representations of proteins, obtaining quantitative properties of the model for realistic binding interface sizes and a full amino acid alphabet.

## 2. Evolutionary paths on a smooth Mount Fuji landscape

We first consider a simple fitness landscape model, a smooth “Mount Fuji” (i.e., single-peaked) landscape [20]. Consider genotypic sequences of length  $L$  with  $k$  possible alleles  $\{A_1, \dots, A_k\}$  at each site, resulting in  $n_{\text{seq}} = k^L$  possible genotypes. We assume the alleles  $\{A_1, \dots, A_k\}$  are in increasing order of fitness rank. The sites could be residues in a protein, nucleotides in a DNA sequence, or larger genomic loci such as whole genes. In general we will interpret the sequences in the model as coarse-grained versions of actual biological sequences. For example, a 12-residue binding interface on a protein with 20 possible amino acids at each site could be coarse-grained into  $L = 6$  pairs of sites with  $k = 5$  alleles at each site, where each allele represents a class of amino acids grouped by physico-chemical properties (e.g., negative, positive, polar, hydrophobic, and other). This is analogous to block spin renormalization in Ising models [21].

Let the occupation number  $n_j(\sigma)$  of a sequence  $\sigma$  be the number of  $A_j$  alleles in the sequence, so that  $\sum_{j=1}^k n_j(\sigma) = L$ . We define the fitness of a sequence  $\sigma$  to be

$$\mathcal{F}(\sigma) = f^{\sum_{j=1}^k (j-1)n_j(\sigma)}, \quad (1)$$

where  $f$  is the multiplicative fitness change from each mutation: a mutation  $A_i \rightarrow A_j$  changes fitness by a factor of  $f^{j-i}$ . If  $f = 1$ , the fitness landscape is flat and evolution is neutral, while if  $f > 1$ , the landscape has a minimum point at  $A_1 A_1 \dots A_1$  ( $\mathcal{F} = 1$ ) and a maximum point at  $A_k A_k \dots A_k$  ( $\mathcal{F} = f^{L(k-1)}$ ). The model is non-epistatic since the fitness function factorizes over sites; thus all mutations have the same fitness effect regardless of the genetic background on which they occur. A more general Mount Fuji model could allow mutations at different sites and between different alleles to have different fitness benefits, although this will not affect the scaling properties of the model that are of primary interest here.

We assume that the population is monomorphic: all organisms have the same genotype at any given time. This approximation holds when  $u \ll (LN \log N)^{-1}$ , where  $u$  is the per-site probability of

mutation per generation and  $N$  is the population size [22]. In this regime the population evolves through a series of substitutions, in which single mutants arise and fix one at a time. A substitution from genotype  $\sigma$  to  $\sigma'$  occurs at the rate [23]

$$W(\sigma'|\sigma) = Nu \phi(s), \quad (2)$$

where  $s = \mathcal{F}(\sigma')/\mathcal{F}(\sigma) - 1$  is the selection coefficient between  $\sigma$  and  $\sigma'$ , and  $\phi(s)$  is the fixation probability of a single mutant with selection coefficient  $s$ . We use the diffusion approximation to the Wright-Fisher model for the fixation probability [24]:

$$\phi(s) = \frac{1 - e^{-2s}}{1 - e^{-2Ns}}. \quad (3)$$

Note that when  $N|s| > 1$  this can be approximated by

$$\phi(s) \approx \begin{cases} 1 - e^{-2s} & \text{if } s > 0, \\ 0 & \text{if } s < 0. \end{cases} \quad (4)$$

That is, when selection is much stronger than genetic drift, deleterious mutations never fix, while beneficial mutations fix with a probability commensurate with their selective advantage. This is often referred to as the “strong-selection weak-mutation” (SSWM) limit [25].

### 2.1. The ensemble of evolutionary paths

For concreteness we consider the following evolutionary process: the population begins at the least fit genotype,  $A_1 A_1 \cdots A_1$ , and evolves according to (2) until it reaches the most fit genotype,  $A_k A_k \cdots A_k$ , for the first time. Define an evolutionary path  $\varphi$  as the ordered sequence of states  $\varphi = (\sigma_0, \sigma_1, \dots, \sigma_\ell)$  traversed by the population during this process, where  $\sigma_0 = A_1 A_1 \cdots A_1$  and  $\sigma_\ell = A_k A_k \cdots A_k$ . The probability of making a single substitution  $\sigma \rightarrow \sigma'$ , given a substitution out of  $\sigma$  occurs, is

$$Q(\sigma'|\sigma) = W(\sigma'|\sigma) \theta(\sigma), \quad (5)$$

where  $\theta(\sigma) = (\sum_{\sigma'} W(\sigma'|\sigma))^{-1}$  is the mean waiting time in  $\sigma$  before a substitution occurs. Thus the probability of taking a path  $\varphi$  is

$$\Pi[\varphi] = \prod_{i=0}^{\ell-1} Q(\sigma_{i+1}|\sigma_i). \quad (6)$$

Since the population is guaranteed to reach the final state eventually,  $\sum_{\varphi} \Pi[\varphi] = 1$ , where the sum is over all first-passage paths  $\varphi$  between the initial and final states.

We are interested in statistical properties of the evolutionary path ensemble. We can calculate many such properties using an exact numerical algorithm described in Appendix A [26, 27]. Here we are especially interested in the distribution of path lengths

$\ell$ , i.e., the number of substitutions experienced by the population before it first reaches the fitness maximum. The path length distribution  $\rho(\ell)$  is defined as

$$\rho(\ell) = \sum_{\varphi} \delta_{\ell, \mathcal{L}[\varphi]} \Pi[\varphi], \quad (7)$$

where  $\mathcal{L}[\varphi]$  is the length of path  $\varphi$ , and  $\delta$  is the Kronecker delta. We can similarly express the mean  $\bar{\ell}$  and variance  $\ell_{\text{var}}$  of path length. We also consider the path entropy  $S_{\text{path}}$ , defined as

$$S_{\text{path}} = - \sum_{\varphi} \Pi[\varphi] \log \Pi[\varphi]. \quad (8)$$

This quantity measures the predictability of evolution in sequence space: if only a single path is accessible, then  $S_{\text{path}} = 0$ , in which case evolution is perfectly predictable. Larger values of  $S_{\text{path}}$ , on the other hand, indicate a more diverse ensemble of accessible pathways, and thus less predictable evolution.

### 2.2. Neutral limit

We first consider properties of the evolutionary path ensemble in the case of neutral evolution ( $f = 1$  in (1)). For simple random walks on finite discrete spaces, previous work has shown that the mean path length scales with the total number of states [28, 29], while the distribution of path lengths will be approximately exponential [28]. Thus for neutral evolution,

$$\bar{\ell} \sim n_{\text{seq}} = k^L, \quad \ell_{\text{var}} \sim \bar{\ell}^2 \sim k^{2L}. \quad (9)$$

Conceptually, this means the population on average must explore the entire sequence space before reaching a particular point for the first time, and thus the average number of substitutions grows exponentially with the length of the sequence. Moreover, since the standard deviation is of the same order as the mean, paths much longer than the mean are likely.

Let  $\gamma$  be the average connectivity, defined as the average number of single substitutions accessible from each sequence; in neutral evolution all substitutions are accessible, so  $\gamma = L(k-1)$ . Since all substitutions are equally likely,  $Q(\sigma'|\sigma) = \gamma^{-1}$  for  $\sigma$  and  $\sigma'$  separated by a single mutation. The entropy of the neutral path ensemble is therefore [27]

$$\begin{aligned} S_{\text{path}} &= - \sum_{\varphi} \Pi[\varphi] \log \gamma^{-\mathcal{L}[\varphi]}, \\ &= \bar{\ell} \log \gamma \\ &\sim k^L \log L(k-1). \end{aligned} \quad (10)$$

The path entropy consists of two distinct components: the average path length and the average connectivity. We can consider the factor of  $\log \gamma$  as the average entropy contribution from each jump in the path.

It is worth noting that mean path length (and the distribution of path lengths in general) does *not* have explicit dependence on connectivity: it only depends on the size of the space. So it is the enormous size, not the dimensionality, of sequence space that causes neutral evolution to require so many steps to reach a particular point. In contrast, path entropy, and thus evolutionary predictability, depends on *both* the size and dimensionality of sequence space.

### 2.3. Strong-selection limit

We now consider evolutionary paths in the strong-selection limit. Here all beneficial mutations are selected so strongly ( $f \gg 1$  in (1)) that their fixation probabilities are all approximately 1, while deleterious mutations never occur. Thus evolutionary paths proceed strictly upward on the fitness landscape. This is sometimes called the “adaptive walk” scenario [2], and it is identical to zero-temperature Monte Carlo with energy replaced by negative fitness [3]. Since the fitness landscape is non-epistatic and reverse mutations are impossible, each site can be considered to evolve independently. In particular, we can decompose the total path length into a sum of independent path lengths for individual sites, so that the path length cumulants for the whole sequence are simply sums of the cumulants for individual sites. (Note that the restriction to first-passage paths effectively couples all the sites because they must all reach their final states simultaneously, and so site independence is only valid when reverse mutations are prohibited.)

In Appendix B we show that  $\bar{\ell} = H_{k-1}$  (the  $(k-1)$ th harmonic number) for a single site in the strong-selection limit, and hence the mean length for  $L$  sites is  $LH_{k-1}$ . Since  $H_{k-1} = \log k + b + \mathcal{O}(k^{-1})$ , the mean length scales as

$$\bar{\ell} \sim L(\log k + b). \quad (11)$$

We explicitly include the  $\mathcal{O}(1)$  constant  $b$  here since it may be comparable to  $\log k$  if  $k$  is not too large. For the harmonic numbers,  $b$  is exactly the Euler-Mascheroni constant  $\gamma_{\text{EM}} \approx 0.57721$ , but we use generic notation here as this same scaling form will be fit to an empirical model in the next section. Equation (11) implies that  $\bar{\ell}$  scales approximately logarithmically with the size  $n_{\text{seq}}$  of sequence space, compared with the linear scaling seen in the neutral case. Moreover, Appendix C shows that  $\rho(\ell)$  is approximately Poisson, and thus the variance  $\ell_{\text{var}}$  should obey the same scaling as  $\bar{\ell}$ .

The average connectivity of sequence space is reduced compared to the neutral case, since only beneficial substitutions are allowed. The connectivity averaged over all sequences is  $L(k-1)/2$  (Appendix D); the reduction by a factor of 2 is intuitively explained by the fact that every allowed beneficial substitution

has a prohibited deleterious substitution. For the path entropy under strong selection, we take as an ansatz the same dependence on  $\bar{\ell}$  and  $\gamma$  as in (10), albeit with different  $L, k$  scaling:

$$S_{\text{path}} \sim \bar{\ell} \log \gamma \sim L(\log k + b) \log \frac{1}{2} L(k-1). \quad (12)$$

We numerically verify this ansatz in the next section (figure 1).

### 2.4. Coarse-graining and landscape-dependence of scaling relations

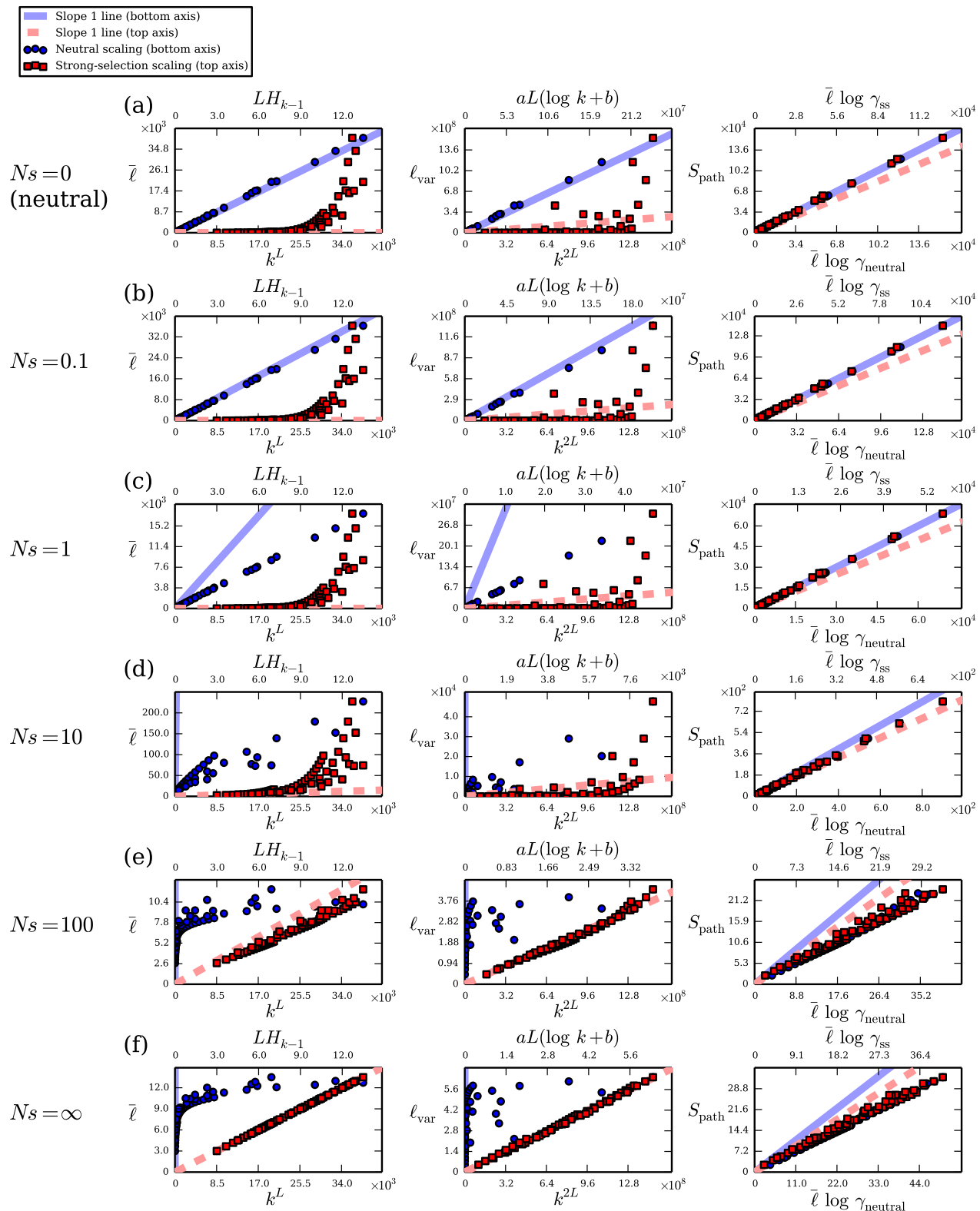
The path scaling relations depend qualitatively on whether the fitness landscape is flat (neutral evolution) or infinitely-steep (strong selection). How does the transition between these two limits occur at intermediate selection strengths, where selection and stochastic fluctuations (genetic drift) compete more equally? We now implement a concrete scheme for coarse-graining sequence space on the fitness landscape of (1). Let  $s = f^{L(k-1)} - 1$  be the total selection coefficient between the maximum and minimum fitness points on the landscape. As we vary the coarse-grained sequence parameters  $L$  and  $k$ , we will rescale the permutation fitness benefit  $f$  in (1) such that the total selection coefficient  $s$  is held fixed:

$$f = (1 + s)^{1/(L(k-1))}. \quad (13)$$

Thus the fitness benefit of each individual mutation decreases as we increase  $L$  and  $k$ . For each choice of  $s$ , we numerically calculate path statistics for a range of  $L$  and  $k$  using the method of Appendix A.

In figure 1 we show the scaling of  $\bar{\ell}$ ,  $\ell_{\text{var}}$ , and  $S_{\text{path}}$  calculated in this manner for several values of relative selection strength  $Ns$ . For  $Ns = 0$ , we not only confirm the neutral scaling relations (9) but also observe that any proportionality factors and additive constants are so negligible that the scaling relations are actually approximate equalities (figure 1a). The predicted relation for the path entropy (10) also holds exactly. Moreover, weak selection appears to preserve these scaling relations: they still hold even at  $Ns = 0.1$  (figure 1b). When selection becomes comparable to genetic drift ( $Ns = 1$ , figure 1c), the neutral scaling relations still hold qualitatively, although the slopes of  $\bar{\ell} \sim k^L$  and  $\ell_{\text{var}} \sim k^{2L}$  are no longer close to 1, indicating different proportionality factors.

At the other extreme ( $Ns = \infty$ , figure 1f), the predicted scaling relations (11) for path length hold as expected. We also verify that  $S_{\text{path}} \sim \bar{\ell} \log \gamma$  even for strong selection, albeit with a proportionality factor less than 1. This scaling maintains at finite but large selection strengths of  $Ns = 100$  (figure 1e). At intermediate selection strengths ( $Ns = 10$ , figure 1d), however, neither set of scaling relations for  $\bar{\ell}$  and  $\ell_{\text{var}}$



**Figure 1.** Scaling properties of evolutionary paths on the Mount Fuji landscape (1) for different values of  $Ns$ , where  $N = 1000$  and  $s$  is the total selection coefficient from the least fit to the most fit sequence on the landscape: (a)  $Ns = 0$  (neutral evolution), (b)  $Ns = 0.1$ , (c)  $Ns = 1$ , (d)  $Ns = 10$ , (e)  $Ns = 100$ , and (f)  $Ns = \infty$ . The left column shows mean path length (number of substitutions)  $\bar{\ell}$ , the middle column shows path length variance  $\ell_{\text{var}}$ , and the right column shows path entropy  $S_{\text{path}}$ . Each panel plots numerical data against both neutral scaling parameters on the bottom axes (blue circles and the solid blue line of slope 1;  $\gamma_{\text{neutral}} = L(k-1)$ ), as well as strong-selection scaling parameters on the top axes (red squares and the dashed red line of slope 1;  $\gamma_{\text{ss}} = L(k-1)/2$ ). Numerical values of the variance  $\ell_{\text{var}}$  are fitted to a function of the form  $aL(\log k + b)$  for each value of  $Ns$  separately. We show all  $L > 1$  and  $k > 2$  such that  $kL < 4 \times 10^4$ .



holds, indicating that they are no longer a simple function of sequence space size  $k^L$ .

### 3. Evolutionary paths in a biophysical model of protein adaptation

Simple model landscapes defined in genotype space, such as (1), have produced many theoretical results and guided analysis of some data [2–5, 8, 10]. However, their purely phenomenological nature allows for little interpretation of their parameters and includes no basis in the underlying molecular processes — interactions among proteins, DNA, RNA, and other biomolecules — that govern cells. Thus a promising alternative is to develop models of fitness that explicitly account for these molecular properties [14, 16, 17, 30, 31]. We now consider the scaling properties of evolutionary paths in such a model based on the biophysics of protein folding and binding [18, 26, 27].

#### 3.1. Protein energetics and coarse-graining

Consider a protein with two-state folding kinetics [32]. In the folded state, the protein has an interface that binds a target molecule. Because the protein can bind *only* when it is folded, the protein has three possible structural states: folded and bound, folded and unbound, and unfolded and unbound. Let the free energy of folding be  $E_f$  (sometimes known as  $\Delta G$ ), so that an intrinsically-stable protein has  $E_f < 0$ . Let the free energy of binding, relative to the chemical potential of the target molecule, be  $E_b$ , so that  $E_b < 0$  indicates a favorable binding interaction. Note that  $E_b$  becomes more favorable as the chemical potential of the target molecule is increased.

The folding and binding energies depend on the protein’s genotype (amino acid sequence)  $\sigma$ . We assume that adaptation only affects “hotspot” residues at the binding interface [33, 34]. We focus on  $L$  binding hotspot residues which, to a first approximation, make additive contributions to the total folding and binding free energies [35]:

$$\begin{aligned} E_f(\sigma) &= E_f^{\text{ref}} + \sum_{i=1}^L \epsilon_f(i, \sigma^i), \\ E_b(\sigma) &= E_b^{\text{min}} + \sum_{i=1}^L \epsilon_b(i, \sigma^i), \end{aligned} \quad (14)$$

where  $\epsilon_f(i, \sigma^i)$  and  $\epsilon_b(i, \sigma^i)$  are matrix entries capturing the energetic contributions of amino acid  $\sigma^i$  at position  $i$ . Folding and binding energetics are probed experimentally and computationally by measuring the changes (often denoted by  $\Delta\Delta G$ ) in  $E_f$  or  $E_b$  resulting from single-point mutations. For folding energies,

these mutational effects are universally distributed over many proteins [36]; in accordance with this observation, we sample entries of  $\epsilon_f$  from a Gaussian distribution with mean 1.25 kcal/mol and standard deviation 1.6 kcal/mol. The reference energy  $E_f^{\text{ref}}$  captures the fixed contribution to the folding energy from all other residues in the protein, and is defined as the energy of a reference sequence  $\sigma_{\text{ref}}$  (so that  $\epsilon_f(i, \sigma_{\text{ref}}^i) = 0$  for all  $i \in \{1, \dots, L\}$ ).

The parameter  $E_b^{\text{min}}$  is defined as the binding energy of the best-binding genotype  $\sigma_{\text{bb}}$  with the minimum  $E_b$ :  $\epsilon_b(i, \sigma_{\text{bb}}^i) = 0$  for all  $i \in \{1, \dots, L\}$ . Since binding hotspot residues typically have a 1–3 kcal/mol penalty for mutations away from the wild-type amino acid [33, 34], we sample the other entries of  $\epsilon_b$  from an exponential distribution defined in the range of  $(1, \infty)$  kcal/mol, with mean 2 kcal/mol. This distribution is consistent with alanine-scanning experiments which probe energetics of amino acids at the binding interface [37]. The exact shapes of the distributions for  $\epsilon_f$  and  $\epsilon_b$  are unimportant for large enough  $L$  due to the central limit theorem.

For large  $L$  and a full amino acid alphabet ( $k = 20$ ), numerical calculations over all  $k^L$  sequences are not possible. However, we can consider coarse-grained versions of the model by grouping sites and amino acids into classes. If we then determine how properties of the model scale with “effective”  $L$  and  $k$  under such a coarse-graining procedure, we can extrapolate these properties to “physical” values of  $L$  and  $k$ . As we vary  $L$  and  $k$ , we want to hold the global distribution of sequence energies fixed, similar to our coarse-graining scheme in the previous section. The additive energy model in (14) implies that energy scales linearly with length  $L$  and does not depend on  $k$ . Thus we can obtain effective  $\epsilon_f$  and  $\epsilon_b$  matrices for the coarse-grained model:

$$\epsilon_{f,\text{eff}} = \frac{L_{\text{phys}}}{L_{\text{eff}}} \epsilon_{f,\text{phys}}, \quad \epsilon_{b,\text{eff}} = \frac{L_{\text{phys}}}{L_{\text{eff}}} \epsilon_{b,\text{phys}}. \quad (15)$$

For simplicity we will drop the “eff” labels and hereafter interpret  $L$ ,  $k$ ,  $\epsilon_f$ , and  $\epsilon_b$  as these effective, coarse-grained parameters unless otherwise indicated.

#### 3.2. Evolutionary model

Without loss of generality, we assume the protein contributes fitness 1 to the organism when it is both folded and bound. Let  $f_{\text{ub}}, f_{\text{uf}} \in [0, 1]$  be the multiplicative fitness penalties for being unbound and unfolded, respectively: the fitness is  $f_{\text{ub}}$  if the protein is unbound but folded, and  $f_{\text{ub}}f_{\text{uf}}$  if the protein is both unbound and unfolded. Then the fitness of the protein averaged over all three possible structural states is

given by [18]

$$\mathcal{F}(E_f, E_b) = \frac{e^{-\beta(E_f+E_b)} + f_{ub}e^{-\beta E_f} + f_{ub}f_{uf}}{e^{-\beta(E_f+E_b)} + e^{-\beta E_f} + 1}, \quad (16)$$

where  $\beta = 1.7 \text{ (kcal/mol)}^{-1}$  is inverse room temperature and the structural states are assumed to be in thermodynamic equilibrium.

We assume that the population begins as perfectly adapted to binding a target molecule characterized by energy matrix  $\epsilon_{b1}$  with minimum binding energy  $E_{b1}^{\min}$  (defining a fitness landscape  $\mathcal{F}_1$ ). The population is then subjected to a selection pressure which favors binding a new target, with energy matrix  $\epsilon_{b2}$  and minimum binding energy  $E_{b2}^{\min}$  (fitness landscape  $\mathcal{F}_2$ ). The population evolves in the monomorphic limit with the SSWM dynamics in (2) and (4). Thus the evolutionary paths are first-passage paths leading from the genotype corresponding to the global maximum on  $\mathcal{F}_1$  to a local or global maximum on  $\mathcal{F}_2$ , with fitness increasing monotonically along each path.

### 3.3. Case 1: selection for binding strength

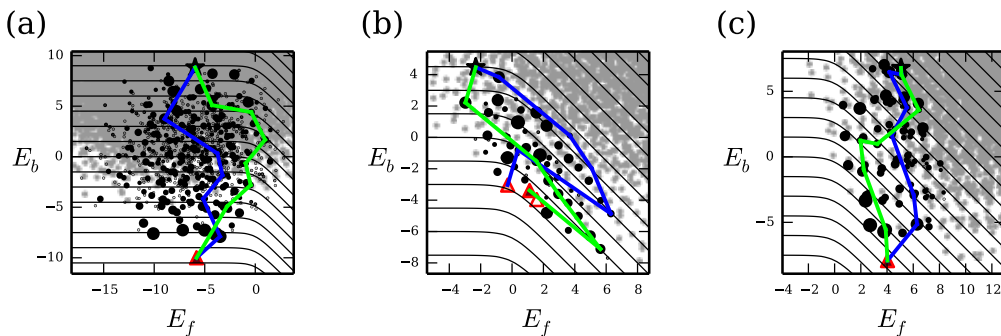
There are three qualitatively distinct cases of the fitness landscape in (16), depending on the values of the parameters  $f_{ub}$  and  $f_{uf}$ . These cases correspond to different biological scenarios for the selection pressures on binding and folding. In the simplest scenario (“case 1”), proteins are selected for their binding function ( $f_{ub} < 1$ ), but misfolding carries no additional fitness penalty (e.g. due to toxicity of misfolded proteins) beyond loss of function ( $f_{uf} = 1$ ). Thus we say there is direct selection for binding only [18]. Three examples of adaptation in this regime are shown in figure 2; the main determinant of the qualitative nature of adaptation is the overall folding stability  $E_f$ .

Although the model is non-epistatic at the level of the energy traits (since (14) is additive), there can be epistasis at the level of fitness (16) due to its nonlinear dependence on energy. Indeed, there is widespread magnitude epistasis, which occurs when the fitness effect of a mutation has different magnitude on different genetic backgrounds, although it is always beneficial or always deleterious. Sign epistasis, which occurs when a mutation can be beneficial on one background but deleterious on another, manifests itself as curvature in the fitness contours in energy space, as shown in figure 2. However, we see that the landscape is largely free of sign epistasis except near  $E_f = 0$ , where there is a higher probability of multiple local fitness maxima (figure 2b). Overall, this suggests that the scaling relations from the non-epistatic Mount Fuji model may provide a reasonable approximation for this model of protein adaptation; the approximately additive nature of protein traits as in (14) has led

to applications of the Mount Fuji model to proteins previously [7, 38, 39].

In figure 3 we show scaling properties of the fitness landscape for the three stability regimes of the model for case 1 (corresponding to the examples in figure 2). The minimum path length  $\ell_{\min}$  is the Hamming distance between the initial and final states for adaptation; for a randomly-chosen initial sequence,  $\ell_{\min} = L(1 - 1/k)$  on average. Indeed, this relation accurately describes the stable protein regime (figure 3a). For stable proteins, there is no selection pressure to improve stability further, so the global fitness maximum is almost always the best-binding sequence. Since the binding energetics for the old and new targets are uncorrelated, the initial and final states are uncorrelated as well. For marginally-stable and unstable proteins,  $\ell_{\min}$  still scales with  $L(1 - 1/k)$ , but with a reduced slope. This is due to the fact the initial and final states become correlated in these two cases. We can think of this effect as a reduction in the effective length  $L$ , since more beneficial mutations are already present in the initial state. We see similar behavior in the average connectivity  $\gamma$  and accessible size  $n_{\text{seq}}$  of sequence space (figure 3b,c). Note that a random initial state reduces the average connectivity of the accessible sequence space by an additional factor of 2, yielding  $\gamma = L(k - 1)/4$  (see Appendix D).

Whereas stable and unstable proteins almost always have a single fitness maximum, marginally-stable proteins have a sizable probability of multiple maxima owing to greater sign epistasis (figure 2b). In a purely random, uncorrelated fitness landscape, the average number  $m$  of local maxima is  $k^L/(L(k - 1) + 1)$  [2]. This has the form  $n_{\text{seq}}/(\gamma + 1)$ : the number of maxima increases with the total size of the space and decreases with the connectivity. We empirically test this scaling for the average number of maxima for a marginally-stable protein, and we find good agreement (figure 3d). By fitting numerically-calculated values of  $m$  as a power law of  $n_{\text{seq}}/(\gamma + 1)$ , we obtain an anomalous scaling exponent of  $\approx 0.18$ ; the fact this is less than 1 reflects the correlated nature of our fitness landscape. The fitted scaling relation allows us to accurately determine the average number of local maxima for binding hotspots and amino acid alphabets much larger than we can directly calculate. By also fitting  $\gamma$  as a linear function of  $L(k - 1)/4$  (figure 3b) and  $n_{\text{seq}}$  as a power law of  $((k + 1)/2)^L$  (figure 3c), for a marginally-stable protein with  $L_{\text{phys}} = 12$  hotspot residues and an amino acid alphabet of size  $k_{\text{phys}} = 20$ , we estimate the number of local maxima to be  $\approx 84.7$ . This is a large number of maxima in absolute terms, although it is far smaller than the expected number on an uncorrelated random landscape of the same size ( $k^L/(L(k - 1) + 1) \approx 1.8 \times 10^{13}$ ).



**Figure 2.** Example landscapes of protein adaptation with direct selection for binding only (case 1), zoomed into the region of energy space accessible to evolutionary paths in our model. (a) Stable protein with  $E_f^{\text{ref}} = -20$  kcal/mol, (b) marginally-stable protein with  $E_f^{\text{ref}} = -5$  kcal/mol, and (c) intrinsically-unstable protein with  $E_f^{\text{ref}} = 5$  kcal/mol. In all panels  $f_{\text{ub}} = 0$ ,  $f_{\text{uf}} = 1$ , and  $E_{\text{b1}}^{\text{min}} = E_{\text{b2}}^{\text{min}} = -10$  kcal/mol. The coarse-grained sequence parameters are  $L = 6$  and  $k = 5$ , with energies rescaled according to (15) with  $L_{\text{phys}} = 12$ . The black star indicates the initial state for adaptation (global maximum on  $\mathcal{F}_1$ ); red triangles indicate local fitness maxima on  $\mathcal{F}_2$ , shaded according to their commitment probabilities (probability of reaching that final state starting from the initial state); black circles indicate intermediate states along paths, sized proportional to their path density (total probability of paths passing through them); small gray circles are genotypes inaccessible to adaptation. The black contours indicate constant fitness  $\mathcal{F}_2$  (the fitness difference between adjacent contours is non-uniform so that they are equidistant in energy space), while example paths are shown in blue and green.

In figure 4 we show the scaling of path statistics  $\bar{\ell}$ ,  $\ell_{\text{var}}$ , and  $S_{\text{path}}$ . We find that the strong-selection scaling relations describe these cases of the protein model very well, despite the complexities of the energy and fitness model relative to the simple Mount Fuji case. The main discrepancy is in the path length variance, indicating that the distributions  $\rho(\ell)$  are not as close to Poisson as in the Mount Fuji model. We expect this is mainly due to the epistasis present in the protein model. Nevertheless, the scaling is accurate enough to extend the model to larger binding interfaces and a full amino acid alphabet. For example, using the fitted coefficients  $a$  and  $b$  (figure 4a,b), we estimate  $\bar{\ell} \approx 25.5$  and  $\ell_{\text{var}} \approx 11.3$  for a marginally-stable protein with  $L_{\text{phys}} = 12$  hotspot residues and an amino acid alphabet of size  $k_{\text{phys}} = 20$ . Comparing these against the estimated  $\ell_{\text{min}} \approx 9.7$  (fitted as a linear function of  $L(1 - 1/k)$ ; figure 3a), we see that many more substitutions than the minimum are likely.

### 3.4. Cases 2 and 3: selection for folding stability

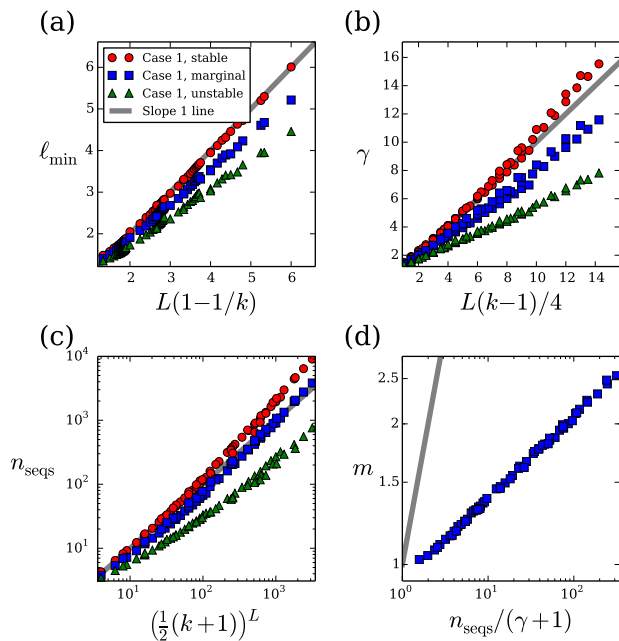
The fitness landscape changes qualitatively when there are additional selection pressures against misfolding beyond loss of function, e.g., for proteins that form toxic aggregates when misfolded [40–42]. The first possibility is that the protein has a non-functional binding interaction ( $f_{\text{ub}} = 1$ ) but is deleterious when misfolded ( $f_{\text{uf}} < 1$ ; “case 2”). Here the relative binding strengths of the old and new targets lead to different patterns of adaptation. In figure 5a, we show an example of adaptation when both the old and new targets have potentially strong (but non-functional)

binding affinity, while figure 5b shows an example when the old target has weak affinity while the new one has strong affinity. Figure 5c shows the case when the old target has strong affinity and the new target has little to no affinity.

Finally, the most general case is to have distinct selection pressures on both binding and folding ( $0 < f_{\text{ub}} < 1$  and  $f_{\text{uf}} < 1$ ; “case 3”). Adaptation in this scenario often resembles binding-only selection (figure 2), except when both binding and folding are of marginal strength (i.e.,  $E_f \simeq 0$  and  $E_b \simeq 0$ ). In this case, the distribution of genotypes in energy space straddles a straight diagonal fitness contour, leading to a distinct pattern of evolutionary paths that gain extra folding stability first, only to lose it later as binding improves (figure 5d).

We show the scaling properties of the evolutionary paths for cases 2 and 3 in figure 6. In general, the predicted scaling relations are less accurate compared to binding-only selection (case 1). This is due to the increased sign epistasis in these regimes (note significant curvature in the fitness contours in figure 5). Selection for both binding and folding (case 3) is particularly epistatic in the  $E_f \simeq 0$ ,  $E_b \simeq 0$  regime, leading to the largest deviations from the Mount Fuji scaling (figure 6). On the other hand, the degree of epistasis here is still far from the maximally-epistatic, uncorrelated random landscape [2,6]; in that model we should have  $\bar{\ell} \sim \log L$  [3], which is clearly not the case in our biophysical model.

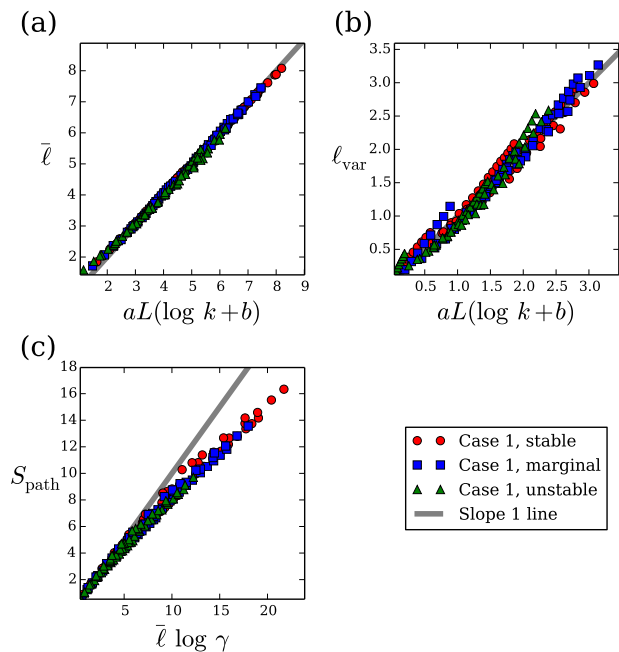




**Figure 3.** Scaling of landscape properties for three regimes of protein adaptation with direct selection for binding only (case 1). (a) Minimum path length  $\ell_{\min}$ , equal to the Hamming distance between the initial and final states, versus  $L(1-1/k)$ ; (b) average connectivity  $\gamma$  versus  $L(k-1)/4$ ; (c) average number  $n_{\text{seq}}$  of accessible sequences versus  $((k+1)/2)^L$ ; (d) average number  $m$  of local fitness maxima versus  $n_{\text{seq}}/(\gamma+1)$ . In all panels red circles are for stable proteins, blue squares are for marginally-stable proteins, and green triangles are for intrinsically-unstable proteins, with all energy and fitness parameters the same as in figure 2. Each point represents an average over  $10^4$  realizations of the folding and binding energy matrices; trivial realizations where the initial state is already a local maximum are excluded. We include all  $L > 1$  and  $k > 2$  such that  $k^L < 4 \times 10^4$ , coarse-grained according to (15) with  $L_{\text{phys}} = 12$ . Slope 1 lines from the origin are shown in gray to guide the eye.

#### 4. Discussion

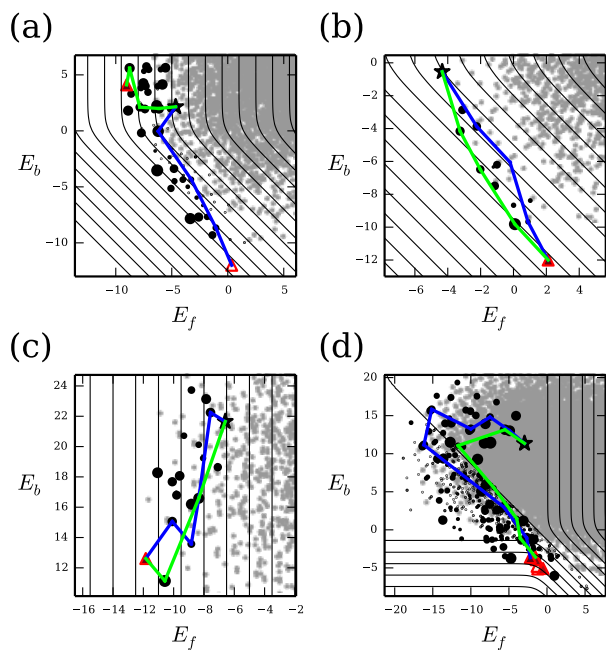
Developing models of fitness landscapes based on the physics of proteins and other biomolecules has emerged as a powerful approach for understanding molecular evolution [14, 16, 17, 30, 31]. However, the empirical nature of these models often makes explicit analytic treatments impossible, while the enormous size of sequence space often restricts numerical calculations or simulations to short sequences  $L$  or reduced alphabet sizes  $k$ . While analyses with small  $L$  and  $k$  may preserve qualitative properties of the models, quantitatively extending these results to more realistic parameter values is essential for comparison with experimental data. Here we have developed a scaling approach in which we empirically fit small  $L$  and  $k$  calculations to scaling relations in order to obtain precise quantitative properties of the model for arbitrarily large  $L$  and  $k$ . The scaling analysis moreover confirms that small  $L$  and  $k$  calculations



**Figure 4.** Scaling of path properties for three regimes of protein adaptation with direct selection for binding only (case 1). (a) Mean path length (average number of substitutions)  $\bar{\ell}$  and (b) path length variance  $\ell_{\text{var}}$  versus  $aL(\log k + b)$ , where the parameters  $a$  and  $b$  are fitted separately for  $\bar{\ell}$  and  $\ell_{\text{var}}$  and for stable, marginal, and unstable proteins. (c) Path entropy  $S_{\text{path}}$  versus  $\bar{\ell} \log \gamma$ . All symbols are the same as in figure 3.

largely preserve qualitative properties of the model expected for realistic sequence spaces. Although the scaling relations are derived for a much simpler, purely non-epistatic model, they are surprisingly robust to the widespread magnitude epistasis and limited sign epistasis observed in the biophysical fitness model.

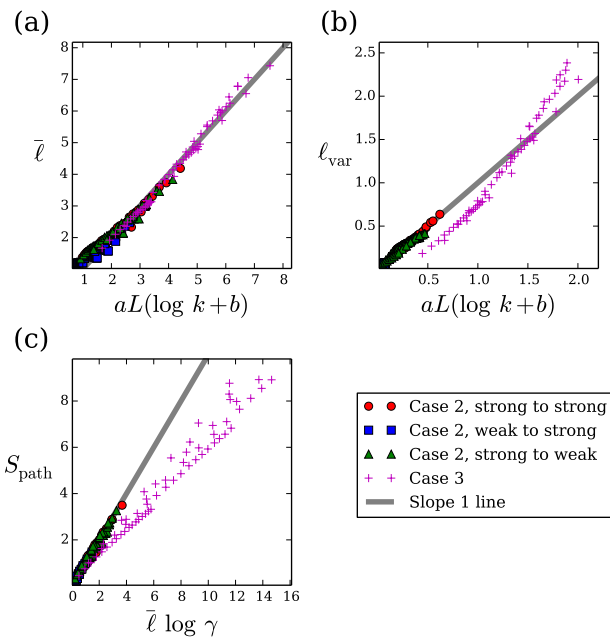
We also gain important conceptual insights from the scaling analysis. In particular, we find that the neutral evolution scaling ( $\bar{\ell} \sim n_{\text{seq}} = k^L$ ,  $\ell_{\text{var}} \sim \bar{\ell}^2 \sim k^{2L}$ ) holds even when selection is present, provided that it is not too strong ( $Ns \leq 1$ , figure 1a,b,c). This means that the average number of substitutions to a global fitness maximum, even in the presence of weak selection, grows exponentially with  $L$ . Strong selection of course enables populations to find the global maximum much faster: the mean path length scales with the logarithm of sequence space size, and the distribution of path lengths is approximately Poisson rather than exponential. However, extremely strong selection ( $Ns \approx 100$ , figure 1e) is required for this more efficient behavior to take over. Selection of this magnitude may be produced by sudden environmental changes, as in our model of protein adaptation [18]. When selection is of more moderate strength ( $Ns \approx 10$ ), mean path length is not a simple function of sequence space size (figure 1d). We expect the more



**Figure 5.** Example landscapes of protein adaptation with selection for folding stability (cases 2 and 3). (a) Direct selection for folding only ( $f_{ub} = 1$ ,  $f_{uf} = 0$ ) where both the old and new targets have potentially strong binding ( $E_{b1}^{\min} = E_{b2}^{\min} = -12$  kcal/mol); (b) same selection as (a) but where the old target has weak binding ( $E_{b1}^{\min} = 0$ ) and the new target binds strongly ( $E_{b2}^{\min} = -12$  kcal/mol); (c) same selection as (a) but where the old target has strong binding ( $E_{b1}^{\min} = -12$  kcal/mol) and the new target binds weakly ( $E_{b2}^{\min} = 0$ ); and (d) direct selection for both binding and folding ( $f_{ub} = 0.9$ ,  $f_{uf} = 0$ ) with marginal stability and binding ( $E_f^{\text{ref}} = E_{b1}^{\min} = E_{b2}^{\min} = -10$  kcal/mol). All symbols are the same as in figure 2. The coarse-grained sequence parameters are  $L = 6$  and  $k = 5$ , with energies rescaled according to (15) with  $L_{\text{phys}} = 12$ .

complex relation in this case to depend on the specific details of the landscape and evolutionary dynamics.

These insights are valuable for general random walks on complex landscapes, e.g., for spin models where  $L$  is the number of spins and  $k$  is the number of individual spin states. The scaling properties of first-passage paths have been well-studied for random walks in the absence of an energy or fitness landscape [29, 43]. However, the effects of a landscape on scaling are less well known. Although the substitution dynamics of (2) considered here are different from the typical dynamics used in spin models and other random walks (e.g., Monte Carlo) [21], we expect our qualitative findings to remain valid. Thus we expect the pure random walk scaling ( $T = \infty$ ) to hold even for temperatures down to the size of the largest energy differences on the landscape. There is a non-trivial crossover regime at temperatures around the size of these landscape features, and at small  $T$  the  $T = 0$  scaling takes over. Investigating the nature of this crossover in both



**Figure 6.** Scaling of path properties for protein adaptation with selection for folding (cases 2 and 3). Panels are the same as figure 4 but with numerical data calculated using energy and fitness parameters matching examples in figure 5: red circles are for case 2 (direct selection for folding only) proteins with strong binding to both old and new targets (figure 5a); blue squares are for case 2 proteins with weak binding to the old target but strong binding to the new one (figure 5b); green triangles are for case 2 proteins with strong binding to the old target but weak binding to the new one (figure 5c); and purple crosses are for case 3 (direct selection for both folding and binding) proteins (figure 5d).

evolutionary and physical models is an important topic for future work.

## Acknowledgments

AVM was supported by an Alfred P. Sloan Research Fellowship.

## Appendix A. Numerical algorithm for statistics of the path ensemble

We calculate statistical properties of the evolutionary paths using an exact algorithm based on transfer matrices [26, 27]. Let  $Q(\sigma'|\sigma)$  be the jump probability defined by a rate matrix as in (5). For each substitution  $\ell$  and intermediate genotype  $\sigma$ , we calculate  $P_\ell(\sigma)$ , the total probability of all paths that end at  $\sigma$  in  $\ell$  substitutions, as well as  $\Gamma_\ell(\sigma)$ , their total entropy. These quantities obey the following recursion relations:

$$P_\ell(\sigma') = \sum_{\sigma} Q(\sigma'|\sigma) P_{\ell-1}(\sigma),$$

$$\Gamma_\ell(\sigma') = \sum_{\sigma \sim \sigma'} Q(\sigma'|\sigma) [\Gamma_{\ell-1}(\sigma) - (\log Q(\sigma'|\sigma)) P_{\ell-1}(\sigma)], \quad (\text{A.1})$$

where  $P_0(\sigma) = 1$  if  $\sigma$  is the initial state and  $P_0(\sigma) = 0$  otherwise, and  $\Gamma_0(\sigma) = 0$  for all  $\sigma$ . The sums are over the  $L(k-1)$  nearest mutational neighbors  $\sigma$  of  $\sigma'$ , and the final states are treated as absorbing to ensure that only first-passage paths are counted. We use these transfer-matrix objects to calculate the path ensemble quantities described in the text:

$$\rho(\ell) = \sum_{\sigma \in \mathcal{S}_f} P_\ell(\sigma), \quad S_{\text{path}} = \sum_{\ell=1}^{\Lambda} \sum_{\sigma \in \mathcal{S}_{\text{final}}} \Gamma_\ell(\sigma), \quad (\text{A.2})$$

where  $\mathcal{S}_{\text{final}}$  is the set of final states. The sums are calculated up to a path length cutoff  $\Lambda$ , which we choose such that  $1 - \sum_{\ell=1}^{\Lambda} \rho(\ell) < 10^{-6}$ . The time complexity of the algorithm scales as  $\mathcal{O}(\gamma n \Lambda)$  [26], where  $\gamma$  is the average connectivity and  $n$  is the total size of the state space.

## Appendix B. Mean path length in the strong-selection limit

Since sites can be considered independent in the strong-selection limit, we need only calculate the mean path length for a single site with  $k$  possible alleles. A path begins at  $A_1$ , and initially all  $k$  alleles are of equal or higher fitness and are therefore accessible. The first substitution can go to any  $A_j \in \{A_2, \dots, A_k\}$  with equal probability  $(k-1)^{-1}$ , after which there are  $k-j+1$  remaining alleles. Thus the mean path length  $\bar{\ell}_k$  for  $k$  alleles must satisfy the recursion relation

$$\bar{\ell}_k = 1 + \frac{1}{k-1} \sum_{j=2}^k \bar{\ell}_{k-j+1}, \quad (\text{B.1})$$

where  $\bar{\ell}_1 = 0$ . This is satisfied by

$$\bar{\ell}_k = H_{k-1}, \quad (\text{B.2})$$

where  $H_n$  is the  $n$ th harmonic number defined by

$$H_n = \sum_{j=1}^n \frac{1}{j}. \quad (\text{B.3})$$

To prove this, we first note that

$$\begin{aligned} \sum_{j=1}^n H_n &= n + \frac{n-1}{2} + \frac{n-2}{3} + \dots + \frac{1}{n} \\ &= \sum_{j=1}^n \frac{n+1-j}{j} \\ &= (n+1)H_n - n \\ &= (n+1)H_{n+1} - (n+1), \end{aligned} \quad (\text{B.4})$$

where we have used the property  $H_{n+1} = H_n + (n+1)^{-1}$ . Now we substitute  $\bar{\ell}_j = H_{j-1}$  on the right side of (B.1) and invoke (B.4) to obtain

$$\begin{aligned} 1 + \frac{1}{k-1} \sum_{j=2}^k H_{k-j} \\ &= 1 + \frac{1}{k-1} \sum_{j=1}^{k-2} H_j \\ &= 1 + \frac{1}{k-1} ((k-1)H_{k-1} - (k-1)) \\ &= H_{k-1}. \end{aligned} \quad (\text{B.5})$$

This proves (B.2) is the solution to the recursion relation.

## Appendix C. Distribution of path lengths in the strong-selection limit

Here we address the whole path length distribution  $\rho(\ell)$  for a single site in the strong-selection limit. With alleles ordered by fitness rank, a path of  $\ell$  substitutions is of the form  $A_1 \rightarrow A_{j_1} \rightarrow \dots \rightarrow A_{j_{\ell-1}} \rightarrow A_k$ , where  $1 < j_1 < \dots < j_{\ell-1} < k$ . Since all beneficial substitutions are equally likely in this limit, the jump probability out of allele  $A_j$  is  $(k-j)^{-1}$ . Therefore the probability of taking a path of length  $\ell$  is

$$\rho(\ell) = \frac{1}{k-1} \sum_{j_1=2}^{k-(\ell-1)} \frac{1}{k-j_1} \sum_{j_2=j_1+1}^{k-(\ell-2)} \frac{1}{k-j_2} \dots \sum_{j_{\ell-1}=j_{\ell-2}+1}^{k-1} \frac{1}{k-j_{\ell-1}}. \quad (\text{C.1})$$

The mean of this distribution is exactly  $\bar{\ell} = H_{k-1}$  as shown in Appendix B. Here we obtain an approximate form for the whole distribution. Define  $\epsilon = k^{-1}$  and  $x_i = j_i/k$ . For  $k \gg 1$  ( $\epsilon \ll 1$ ) we can take the continuum limit of the exact expression to obtain

$$\begin{aligned} \rho(\ell) &\approx \frac{1}{k-1} \int_{2\epsilon}^{1-(\ell-1)\epsilon} \frac{dx_1}{1-x_1} \int_{x_1+\epsilon}^{1-(\ell-2)\epsilon} \frac{dx_2}{1-x_2} \dots \\ &\quad \int_{x_{\ell-2}+\epsilon}^{1-\epsilon} \frac{dx_{\ell-1}}{1-x_{\ell-1}}. \end{aligned} \quad (\text{C.2})$$

By changing variables to  $y_i = x_i - (i+1)\epsilon$ , we rewrite this as

$$\begin{aligned} \rho(\ell) &\approx \frac{1}{k-1} \int_0^{1-(\ell+1)\epsilon} \frac{dy_1}{1-y_1-2\epsilon} \\ &\quad \int_{y_1}^{1-(\ell+1)\epsilon} \frac{dy_2}{1-y_2-3\epsilon} \dots \end{aligned}$$

$$\int_{y_{\ell-2}}^{1-(\ell+1)\epsilon} \frac{dy_{\ell-1}}{1-y_{\ell-1}-\ell\epsilon}. \quad (\text{C.3})$$

Each integral is dominated by its integrand's value near the upper limit. However, because the domain of integration requires ordering of the  $y_i$  variables ( $0 < y_1 < y_2 < \dots < y_{\ell-1} < 1-(\ell+1)\epsilon$ ), the integrand for  $y_{\ell-1}$  has the greatest support near its upper limit. Since the integrands are all similar near their lower limits, we thus approximate each integrand by the one for  $y_{\ell-1}$ :

$$\frac{1}{1-y_i-(i+1)\epsilon} \approx \frac{1}{1-y_i-\ell\epsilon}. \quad (\text{C.4})$$

This approximation allows us to use the identity

$$\int_a^b dx_1 f(x_1) \int_{x_1}^b dx_2 f(x_2) \cdots \int_{x_{n-1}}^b dx_n f(x_n) = \frac{1}{n!} \left( \int_a^b dx f(x) \right)^n. \quad (\text{C.5})$$

Therefore,

$$\begin{aligned} \rho(\ell) &\approx \frac{1}{k-1} \frac{1}{(\ell-1)!} \left( \int_0^{1-(\ell+1)\epsilon} \frac{dy}{1-y-\ell\epsilon} \right)^{\ell-1} \\ &= \frac{\log^{\ell-1}(k-\ell)}{(\ell-1)!(k-1)}. \end{aligned} \quad (\text{C.6})$$

In the limit of  $k \gg 1$  and  $\ell/k \ll 1$ ,

$$\begin{aligned} \rho(\ell) &\approx \frac{(\log k + \log(1-\ell/k))^{\ell-1}}{(\ell-1)!(k-1)} \\ &\approx \frac{(\log k)^{\ell-1}}{(\ell-1)!} e^{-\log k}. \end{aligned} \quad (\text{C.7})$$

Thus  $\rho(\ell)$  is approximately a Poisson distribution with mean and variance  $\log k$ . This is consistent with the exact solution since  $\bar{\ell} = H_{k-1} \approx \log k$  for large  $k$ .

#### Appendix D. Size and connectivity of sequence space in the strong-selection limit

Starting from the sequence with minimum fitness, all  $k^L$  sequences are accessible in the strong-selection limit. More generally, if the population begins at sequence  $\sigma$ , there are  $\prod_{j=1}^k (k-j+1)^{n_j(\sigma)}$  accessible sequences, including  $\sigma$  itself. If the initial sequence is chosen at random, then the average number of accessible sequences is

$$n_{\text{seq}} = \sum_{n_1, \dots, n_k} \frac{1}{k^L} \binom{L}{n_1, \dots, n_k} \prod_{j=1}^k (k-j+1)^{n_j(\sigma)}$$

$$\begin{aligned} &= \sum_{n_1, \dots, n_k} \binom{L}{n_1, \dots, n_k} \prod_{j=1}^k \left( 1 - \frac{(j-1)}{k} \right)^{n_j(\sigma)} \\ &= \left( \sum_{j=1}^k \left( 1 - \frac{(j-1)}{k} \right) \right)^L \\ &= \left( \frac{k+1}{2} \right)^L. \end{aligned} \quad (\text{D.1})$$

Similarly, each sequence  $\sigma$  has  $\sum_{j=1}^k (k-j)n_j(\sigma)$  possible beneficial mutations. Thus the connectivity averaged over all sequences is

$$\begin{aligned} \gamma &= \sum_{n_1, \dots, n_k} \binom{L}{n_1, \dots, n_k} \sum_{j=1}^k (k-j)n_j \\ &= \sum_{j=1}^k (k-j) \frac{L}{k} \\ &= \frac{1}{2} L(k-1). \end{aligned} \quad (\text{D.2})$$

We can also determine the average connectivity of the accessible sequences starting from a random initial sequence. We first consider a single site. The initial allele  $A_j$  is chosen with probability  $1/k$ , leaving  $k-j+1$  accessible alleles. Thus the average connectivity of this accessible space is

$$\begin{aligned} \gamma &= \sum_{j=1}^k \frac{1}{k} \sum_{i=j}^k \frac{1}{k-j+1} (k-i) \\ &= \frac{1}{4} (k-1). \end{aligned} \quad (\text{D.3})$$

Since multiple sites contribute additively to the connectivity, the total average connectivity of the accessible space is  $L(k-1)/4$ .

#### References

- [1] J. Maynard Smith. Natural selection and the concept of a protein space. *Nature*, 225:563–564, 1970.
- [2] S. Kauffman and S. Levin. Toward a general theory of adaptive walks on rugged landscapes. *J Theor Biol*, 128:11–45, 1987.
- [3] H. Flyvbjerg and B. Lautrup. Evolution in a rugged fitness landscape. *Phys Rev A*, 46:6714–6723, 1992.
- [4] D. R. Rokyta, C. J. Beisel, and P. Joyce. Properties of adaptive walks on uncorrelated landscapes under strong selection and weak mutation. *J Theor Biol*, 243:114–120, 2006.
- [5] J. Franke, A. Klözer, J. A. G. M. de Visser, and J. Krug. Evolutionary accessibility of mutational pathways. *PLoS Comput Biol*, 7:e1002134, 2011.
- [6] J. F. C. Kingman. A simple model for the balance between selection and mutation. *J Appl Probab*, 15:1–12, 1978.
- [7] T. Aita, H. Uchiyama, T. Inaoka, M. Nakajima, T. Kokubo, and Y. Husimi. Analysis of a local fitness landscape with a model of the rough Mt. Fuji-type landscape:

- application to prolyl endopeptidase and thermolysin. *Biopolymers*, 54:64–79, 2000.
- [8] J. Neidhart, I. G. Szendro, and J. Krug. Adaptation in tunably rugged fitness landscapes: The rough mount fuji model. *Genetics*, 2014. genetics.114.167668.
- [9] D. M. Weinreich, N. F. Delaney, M. A. DePristo, and D. L. Hartl. Darwinian evolution can follow only very few mutational paths to fitter proteins. *Science*, 312:111–114, 2006.
- [10] I. G. Szendro, M. F. Schenk, J. Franke, J. Krug, and J. A. de Visser. Quantitative analyses of empirical fitness landscapes. *J Stat Mech*, page P01005, 2013.
- [11] I. G. Szendro, J. Franke, J. A. G. M. de Visser, and J. Krug. Predictability of evolution depends nonmonotonically on population size. *Proc Natl Acad Sci USA*, 110:571–576, 2013.
- [12] S. J. Gould. *Wonderful Life: The Burgess Shale and the Nature of History*. W. W. Norton and Company, New York, USA, 1990.
- [13] A. E. Lobkovsky and E. V. Koonin. Replaying the tape of life: quantification of the predictability of evolution. *Front Gene*, 3:246, 2012.
- [14] J. D. Bloom, S. T. Labthavikul, C. R. Otey, and F. H. Arnold. Protein stability promotes evolvability. *Proc Natl Acad Sci USA*, 103:5869–5874, 2006.
- [15] R. A. Neher and B. I. Shraiman. Competition between recombination and epistasis can cause a transition from allele to genotype selection. *Proc Natl Acad Sci USA*, 106:6866–6871, 2009.
- [16] A. E. Lobkovsky, Y. I. Wolf, and E. V. Koonin. Predictability of evolutionary trajectories in fitness landscapes. *PLoS Comput Biol*, 7:e1002302, 2011.
- [17] M. Heo, S. Maslov, and E. I. Shakhnovich. Topology of protein interaction network shapes protein abundances and strengths of their function and nonspecific interactions. *Proc Natl Acad Sci USA*, 108:4258–4263, 2011.
- [18] M. Manhart and A. V. Morozov. Protein folding and binding can emerge as evolutionary spandrels through structural coupling. 2014. arXiv:1408.3786.
- [19] L. D. Bogarad and M. W. Deem. A hierarchical approach to protein molecular evolution. *Proc Natl Acad Sci USA*, 96:2591–2595, 1999.
- [20] S. A. Kauffman and E. D. Weinberger. The NK model of rugged fitness landscapes and its application to maturation of the immune response. *J Theor Biol*, 141:211–245, 1989.
- [21] J. M. Yeomans. *Statistical Mechanics of Phase Transitions*. Oxford University Press, Oxford, 1992.
- [22] N. Champagnat. A microscopic interpretation for adaptive dynamics trait substitution sequence models. *Stoch Proc Appl*, 116:1127–1160, 2006.
- [23] M. Kimura. *The Neutral Theory of Molecular Evolution*. Cambridge University Press, Cambridge, UK, 1983.
- [24] M. Kimura. On the probability of fixation of mutant genes in a population. *Genetics*, 47:713–719, 1962.
- [25] J. H. Gillespie. Molecular evolution over the mutational landscape. *Evolution*, 38:1116–1129, 1984.
- [26] M. Manhart and A. V. Morozov. Path-based approach to random walks on networks characterizes how proteins evolve new functions. *Phys Rev Lett*, 111:088102, 2013.
- [27] M. Manhart and A. V. Morozov. Statistical physics of evolutionary trajectories on fitness landscapes. In R. Metzler, G. Oshanin, and S. Redner, editors, *First-Passage Phenomena and Their Applications*. World Scientific, Singapore, 2014.
- [28] E. M. Boltt and D. ben-Avraham. What is special about diffusion on scale-free nets? *New J Phys*, 7:26–47, 2005.
- [29] S. Condamin, O. Bénichou, V. Tejedor, R. Voituriez, and J. Klafter. First-passage times in complex scale-invariant media. *Nature*, 450:77–80, 2007.
- [30] A. Haldane, M. Manhart, and A. V. Morozov. Biophysical fitness landscapes for transcription factor binding sites. *PLoS Comput Biol*, 10:e1003683, 2014.
- [31] A. W. Serohijos and E. I. Shakhnovich. Merging molecular mechanism and evolution: theory and computation at the interface of biophysics and evolutionary population genetics. *Curr Opin Struct Biol*, 26:84–91, 2014.
- [32] T. E. Creighton. *Proteins: Structures and Molecular Properties*. W.H. Freeman and Company, New York, 1992.
- [33] T. Clackson and J. A. Wells. A hot spot of binding energy in a hormone-receptor interface. *Science*, 267:383–386, 1995.
- [34] I. S. Moreira, P. A. Fernandes, and M. J. Ramos. Hot spots — a review of the protein-protein interface determinant amino-acid residues. *Proteins*, 68:803–812, 2007.
- [35] J. A. Wells. Additivity of mutational effects in proteins. *Biochemistry*, 29:8509–8517, 1990.
- [36] N. Tokuriki, F. Stricher, J. Schymkowitz, L. Serrano, and D. S. Tawfik. The stability effects of protein mutations appear to be universally distributed. *J Mol Biol*, 369:1318–1332, 2007.
- [37] K. S. Thorn and A. A. Bogan. ASEdb: a database of alanine mutations and their effects on the free energy of binding in protein interactions. *Bioinformatics*, 17:284–285, 2001.
- [38] T. Aita and Y. Husimi. Fitness spectrum among random mutants on Mt. Fuji-type fitness landscape. *J Theor Biol*, 182:469–485, 1996.
- [39] T. Aita and Y. Husimi. Adaptive walks by the fittest among finite random mutants on a Mt. Fuji-type fitness landscape. *J Theor Biol*, 193:383–405, 1998.
- [40] M. Bucciantini, E. Giannoni, F. Chiti, F. Baroni, L. Formigli, J. Zurdo, N. Taddei, G. Ramponi, C. M. Dobson, and M. Stefani. Inherent toxicity of aggregates implies a common mechanism for protein misfolding diseases. *Nature*, 416:507–511, 2002.
- [41] D. A. Drummond and C. O. Wilke. Mistranslation-induced protein misfolding as a dominant constraint on coding-sequence evolution. *Cell*, 134:341–352, 2008.
- [42] K. A. Geiler-Samerotte, M. F. Dion, B. A. Budnik, S. M. Wang, D. L. Hartl, and D. A. Drummond. Misfolded proteins impose a dosage-dependent fitness cost and trigger a cytosolic unfolded protein response in yeast. *Proc Natl Acad Sci USA*, 108:680–685, 2011.
- [43] S. Redner. *A Guide to First-Passage Processes*. Cambridge University Press, Cambridge, 2001.