# Path-Based Approach to Random Walks on Networks Characterizes How Proteins Evolve New Functions

Michael Manhart[1] and Alexandre V. Morozov[1,2,*]

[1]*Department of Physics and Astronomy, Rutgers University, Piscataway, New Jersey 08854, USA*
[2]*BioMaPS Institute for Quantitative Biology, Rutgers University, Piscataway, New Jersey 08854, USA*

We develop a path-based approach to continuous-time random walks on networks with arbitrarily weighted edges. We describe an efficient numerical algorithm for calculating statistical properties of the stochastic path ensemble. After demonstrating our approach on two reaction rate problems, we present a biophysical model that describes how proteins evolve new functions while maintaining thermodynamic stability. We use our methodology to characterize dynamics of evolutionary adaptation, reproducing several key features observed in directed evolution experiments. We find that proteins generally fall into two qualitatively different regimes of adaptation depending on their binding and folding energetics.

Random walks on networks are ubiquitous across physics, chemistry, and biology, including molecular evolution [1–3], protein folding [4], chemical reactions [5], transport and search in complex media [6,7], stochastic phenotypes [8], and cell-type differentiation [9–11]. Each node on the network is assigned a value of the objective function (for example, energy or fitness) which defines the rates of jumping to the neighboring nodes. Statistical properties of random walks determine quantities of interest such as mean first-passage times and path length distributions. Characterizing the diversity of stochastic paths is a particularly central issue in evolutionary theory [1–3,12,13].

Analytical treatments of random walks on networks tend to be limited to simple models with equally weighted edges [6,7,14,15], while direct simulations can be computationally expensive, especially when rare events are considered. In reaction rate theory, ensembles of stochastic trajectories may be built by transition path sampling [16–19]; however, this method involves considerable computational costs in complex systems. Another alternative, transition path theory [4,20,21], requires a numerical solution of the backward equation. Neither approach directly addresses the diversity of stochastic paths.

Here we develop a systematic and numerically efficient path-based approach to stochastic processes. Our method is applicable to continuous-time random walks [22] on networks with arbitrary edge weights. The approach is well suited for obtaining statistics that describe the diversity of paths, such as the distribution of path lengths and path entropy. We use it to study adaptive dynamics of proteins evolving a new function while maintaining thermodynamic stability [12,23–26], a phenomenon of central interest in both natural and directed evolution (the latter aimed at engineering proteins with novel enzymatic activities [26,27]).

A continuous-time random walk (semi-Markov jump process) on the discrete state space $\mathcal{S}$ is defined by a set of jump probabilities $\langle \sigma' | \mathbf{Q} | \sigma \rangle$ for $\sigma \to \sigma'$ ($\sigma, \sigma' \in \mathcal{S}$), and probability distributions $\psi_\sigma(t)$ of waiting time $t$ in state $\sigma$ before making a jump [22]. We assume that $\psi_\sigma(t)$ has finite mean $w(\sigma)$ for all $\sigma \in \mathcal{S}$. Note that $\mathcal{S}$ equipped with the jump matrix $\mathbf{Q}$ defines a network with directed, weighted edges.

Define a path $\varphi$ as a sequence of states $\{\sigma_0, \sigma_1, \ldots, \sigma_\ell\}$. The time-independent probability of the system taking the path $\varphi$ is $\Pi[\varphi] = \pi(\sigma_0) \prod_{i=0}^{\ell-1} \langle \sigma_{i+1} | \mathbf{Q} | \sigma_i \rangle$, where $\pi(\sigma_0)$ is the initial state probability. Let $\Phi$ be an ensemble of first-passage paths from a set of initial states $\mathcal{S}_i$ to a set of final states $\mathcal{S}_f$. The partition function for this ensemble is $Z_\Phi = \sum_{\varphi \in \Phi} \Pi[\varphi]$, and the entropy is $S_\Phi = -Z_\Phi^{-1} \sum_{\varphi \in \Phi} \Pi[\varphi] \times \log(\Pi[\varphi]/Z_\Phi)$. Let $\mathcal{L}[\varphi]$ be the length (number of jumps) of path $\varphi$, and let $\mathcal{T}[\varphi] = \sum_{i=0}^{\ell-1} w(\sigma_i)$ be the average time of the path. We also define the average time the path spends in state $\sigma$, $\mathcal{T}_\sigma[\varphi] = \sum_{i=0}^{\ell-1} \delta_{\sigma,\sigma_i} w(\sigma_i)$ ($\delta$ is the Kronecker delta), and the indicator functional $\mathcal{I}_\sigma[\varphi]$, which equals 1 if $\varphi$ contains $\sigma$ and equals zero otherwise.

The average time of paths in the ensemble is then given by $\bar{\tau}_\Phi = \langle \mathcal{T} \rangle_\Phi = Z_\Phi^{-1} \sum_{\varphi \in \Phi} \mathcal{T}[\varphi] \Pi[\varphi]$. The average path length is $\bar{\ell}_\Phi = \langle \mathcal{L} \rangle_\Phi$, and the path length distribution is given by $\rho_\Phi(\ell) = Z_\Phi^{-1} \sum_{\varphi \in \Phi} \delta_{\ell, \mathcal{L}[\varphi]} \Pi[\varphi]$ [19]. Let $\ell_\Phi^{\text{sd}}$ be the standard deviation of path lengths, $\langle \mathcal{I}_\sigma \rangle_\Phi$ the spatial density of paths (the probability of paths in $\Phi$ that visit state $\sigma$), and $\langle \mathcal{T}_\sigma \rangle_\Phi / \bar{\tau}_\Phi$ the fraction of time spent in state $\sigma$. We can also construct multipoint correlation functions such as $\langle \mathcal{I}_{\sigma'} \mathcal{I}_\sigma \rangle_\Phi$.

Let $|\pi\rangle$ be the vector of initial state probabilities and $|\sigma\rangle$ be the vector with 1 at position $\sigma$ and zero otherwise. For each step $\ell$ and intermediate state $\sigma$, we can recursively calculate $P_\ell(\sigma) = \langle \sigma | \mathbf{Q}^\ell | \pi \rangle$, $T_\ell(\sigma)$, and $\Gamma_\ell(\sigma)$, which

are the total probability, average time, and entropy, respectively, of all paths that end at $\sigma$ in $\ell$ steps

$$P_\ell(\sigma') = \sum_{\text{nn } \sigma \text{ of } \sigma'} \langle \sigma' | \mathbf{Q} | \sigma \rangle P_{\ell-1}(\sigma),$$

$$T_\ell(\sigma') = \sum_{\text{nn } \sigma \text{ of } \sigma'} \langle \sigma' | \mathbf{Q} | \sigma \rangle [T_{\ell-1}(\sigma) + w(\sigma) P_{\ell-1}(\sigma)],$$

$$\Gamma_\ell(\sigma') = \sum_{\text{nn } \sigma \text{ of } \sigma'} \langle \sigma' | \mathbf{Q} | \sigma \rangle [\Gamma_{\ell-1}(\sigma)$$
$$- \log \langle \sigma' | \mathbf{Q} | \sigma \rangle P_{\ell-1}(\sigma)], \tag{1}$$

where $P_0(\sigma) = \pi(\sigma)$ and $T_0(\sigma) = \Gamma_0(\sigma) = 0$ for all $\sigma \in \mathcal{S}$, and the sums run over all nearest neighbors (nn) $\sigma$ of $\sigma'$ ($\sigma \in \mathcal{S}_f$ are treated as absorbing states). Therefore,

$$Z_\Phi = \sum_{\ell=1}^\infty \sum_{\sigma \in \mathcal{S}_f} P_\ell(\sigma), \qquad \rho_\Phi(\ell) = \frac{1}{Z_\Phi} \sum_{\sigma \in \mathcal{S}_f} P_\ell(\sigma),$$

$$\bar{\tau}_\Phi = \frac{1}{Z_\Phi} \sum_{\ell=1}^\infty \sum_{\sigma \in \mathcal{S}_f} T_\ell(\sigma), \qquad S_\Phi = \frac{1}{Z_\Phi} \sum_{\ell=1}^\infty \sum_{\sigma \in \mathcal{S}_f} \Gamma_\ell(\sigma). \tag{2}$$

Similarly, we can calculate state-dependent quantities such as $\langle \mathcal{I}_\sigma \rangle_\Phi$ and $\langle \mathcal{T}_\sigma \rangle_\Phi$. The recursion relations are

$$P_\ell(\sigma'; \sigma) = \begin{cases} \displaystyle\sum_{\text{nn } \sigma'' \text{ of } \sigma'} \langle \sigma' | \mathbf{Q} | \sigma'' \rangle P_{\ell-1}(\sigma''; \sigma), & \sigma' \neq \sigma, \\ P_\ell(\sigma), & \sigma' = \sigma, \end{cases}$$

$$T_\ell(\sigma'; \sigma) = \sum_{\text{nn } \sigma'' \text{ of } \sigma'} \langle \sigma' | \mathbf{Q} | \sigma'' \rangle [T_{\ell-1}(\sigma''; \sigma)$$
$$+ \delta_{\sigma,\sigma''} w(\sigma'') P_{\ell-1}(\sigma''; \sigma)], \tag{3}$$

with the initial conditions $P_0(\sigma'; \sigma) = T_0(\sigma'; \sigma) = 0$ for all $\sigma, \sigma' \in \mathcal{S}$, $\sigma \neq \sigma'$ [$P_0(\sigma; \sigma) = \pi(\sigma)$, $T_0(\sigma; \sigma) = 0$]. Then $\langle \mathcal{I}_\sigma \rangle_\Phi = Z_\Phi^{-1} \sum_{\ell=1}^\infty \sum_{\sigma' \in \mathcal{S}_f} P_\ell(\sigma'; \sigma)$ and $\langle \mathcal{T}_\sigma \rangle_\Phi = Z_\Phi^{-1} \sum_{\ell=1}^\infty \sum_{\sigma' \in \mathcal{S}_f} T_\ell(\sigma'; \sigma)$. Finally, we can calculate mean path divergence that characterizes the spatial diversity of the paths in $\Phi$ [13]

$$\mathcal{D}_\Phi = \sum_{\ell=1}^\infty \sum_{\sigma, \sigma' \in \mathcal{S}} d(\sigma, \sigma') P_\ell(\sigma) P_\ell(\sigma'), \tag{4}$$

where $d(\sigma, \sigma')$ is a distance metric on $\mathcal{S}$.

Our algorithm allows for very general definitions of the path ensemble $\Phi$ without having to explicitly enumerate paths. For instance, $\Phi$ can include paths that begin and end at arbitrary sets of states, or are disallowed from passing through arbitrary sets of intermediate states. The time complexity of our algorithm is $\mathcal{O}(\gamma N \Lambda)$ for $Z_\Phi$, $\rho_\Phi(\ell)$, $\bar{\tau}_\Phi$, $S_\Phi$, and $\mathcal{O}(\gamma N^2 \Lambda)$ for $\langle \mathcal{I}_\sigma \rangle_\Phi$, $\langle \mathcal{T}_\sigma \rangle_\Phi$, $\mathcal{D}_\Phi$, where $\gamma$ is the average number of nn, $N$ is the number of states visited by paths in $\Phi$, and $\Lambda \sim \bar{\ell}_\Phi$ is the cutoff path length. For simple random walks, $\bar{\ell}_\Phi \sim N^{d_w/d_f}$ for $d_w \gtrsim d_f$ and $\bar{\ell}_\Phi \sim N$ for $d_w < d_f$, where $d_w$ is the dimension of the walk and $d_f$ is the fractal dimension of the space [7,15].

Therefore, the algorithm scales as $\mathcal{O}(\gamma N^{1+d_w/d_f})$ for $d_w \gtrsim d_f$ and $\mathcal{O}(\gamma N^2)$ for $d_w < d_f$, automatically accounting for the sparseness of network connections.

We now illustrate our approach on two reaction rate problems [28]. First we consider a 2D double-well potential $V(x, y)$ on a square lattice $\mathcal{S}$ [See Supplemental Material [29], Fig. S1(a)]. The potential has two metastable states $A$ and $B$ with boundaries $\partial A$ and $\partial B$. The initial states on $\partial A$ and $\partial B$ are weighted by the equilibrium distribution $\pi(x, y) = e^{-\beta V(x,y)} / \sum_{(x,y) \in \mathcal{S}} e^{-\beta V(x,y)}$, where $\beta = 1/T$ is the inverse temperature. Let us define the ensemble of transition paths (TPs) between $A$ and $B$: these paths begin on either $\partial A$ or $\partial B$ and end on the opposite boundary without crossing any boundaries in between [17,18]. Similarly, we define the ensemble of return paths (RPs) which come back to the boundary on which they started.

Many TP and RP statistics, such as the distribution of path lengths $\rho_{\text{TP}}(\ell)$, mean path divergence $\mathcal{D}_{\text{TP+RP}}$, the density of states $p(x, y | \text{TP}) = \langle \mathcal{T}_{(x,y)} \rangle_{\text{TP}} / \bar{\tau}_{\text{TP}}$, the total TP flux $\lambda$, and reaction rates of transitions between $A$ and $B$, can be calculated straightforwardly with our method. The density of states on transition paths $p(x, y | \text{TP})$ shows two symmetric channels by which most reactions between $A$ and $B$ occur [See Supplemental Material [29], Fig. S1(b)]. To determine the cutoff path length $\Lambda$, we recall that $\rho_\Phi(\ell) \sim e^{-\alpha \ell / \bar{\ell}_\Phi}$ for sufficiently large $\ell$, where $\alpha = \mathcal{O}(1)$ [See Supplemental Material [29], Fig. S2(a)] [15]. Other path statistics, such as the average time $\bar{\tau}_\Phi(\ell)$ of paths up to length $\ell$ [See Supplemental Material [29], Fig. S2(b)], also show exponential asymptotic behavior. Therefore, in practice, one need only consider paths with $\ell < \Lambda$ and infer the contributions of all longer paths from an exponential fit to the tail, taking advantage of the fact that information about longer paths is already contained in the properties of shorter paths.

In general, we expect paths to increase in length and diversity at higher temperatures. However, between $\beta = 5$ and $\beta = 1$ the paths become shorter and less diverse as $T$ increases [See Supplemental Material [29], Figs. S2(c), S2(d)]. This is a signature of entropic switching [30]: at a critical value of $\beta$, the two most energetically favored pathways that dominate the low-$T$ behavior become less favorable than the shorter path through the middle. Entropic switching is reflected in plots of the relative path divergence, $\bar{\tau}_{\text{TP}}$, $\bar{\ell}_{\text{TP}}$, and $S_{\text{TP}}$ [See Supplemental Material [29], Figs. S2(c), S2(d)], which readily generalize to arbitrary network spaces.

We can also calculate the continuous-space limit of the TP flux $\lambda$ and the reaction rates. We analytically continue $\lambda$ as a function of the lattice spacing $\Delta x$ [See Supplemental Material [29], Fig. S2(a), inset], yielding continuous-limit rates of $k_{A \to B} = k_{B \to A} \approx 1.3 \times 10^{-4}$. Therefore, one need only calculate $\lambda$ at a few finite lattice spacings in order to infer continuous-limit rates. Our approach can be straightforwardly extended to reactions on more complex

structures such as fractals, which serve as models of transport in disordered media [6] (See Supplemental Material [29], Fig. S3).

We now apply our methodology to study evolution of protein function; here the function is defined as binding a target molecule such as an enzymatic substrate or another protein. Let $E_f$ be the protein folding free energy (i.e., the free energy difference between its folded and unfolded conformations), and $E_b$ the free energy of binding relative to the chemical potential of the target. Then the protein has the probability of folding $1/(1 + e^{\beta E_f})$ and, independently, the probability of binding $1/(1 + e^{\beta E_b})$, where $\beta = 1.7$ (kcal/mol)$^{-1}$ is inverse room temperature. We assume that the protein contributes fitness 1 if it both folds and binds, and fitness $f_0 < 1$ otherwise [31]. Then fitness averaged over all proteins in an organism is given by

$$\mathcal{F}(E_f, E_b) = \frac{1 + f_0(e^{\beta E_f} + e^{\beta E_b} + e^{\beta(E_f + E_b)})}{(1 + e^{\beta E_f})(1 + e^{\beta E_b})}. \quad (5)$$

The folding and binding energies are functions of the amino acid sequence $\sigma$. We assume that the protein has a small number $L$ of "hot-spot" residues at the binding interface [32], and that each residue makes an additive contribution to the total energy [33]: $E_f(\sigma) = E_f^0 + \sum_{\mu=1}^L \epsilon_f(\mu, \sigma^\mu)$, $E_b(\sigma) = E_b^0 + \sum_{\mu=1}^L \epsilon_b(\mu, \sigma^\mu)$, where $E_f^0$, $E_b^0$ are overall offsets and $\epsilon_f(\mu, \sigma^\mu)$, $\epsilon_b(\mu, \sigma^\mu)$ are the folding and binding energy contributions of amino acid $\sigma^\mu$ at position $\mu$. The offset $E_f^0$ is a fixed contribution to folding energy from the rest of the protein, which we assume to be perfectly adapted; $\epsilon_f$'s are sampled from a Gaussian with 1.25 kcal/mol mean and 1.6 kcal/mol standard deviation [34]. Since binding hot spots typically have a minimum penalty of 1–3 kcal/mol for mutations away from the wild-type amino acid [35], we set $\epsilon_b(\mu, \sigma_{bb}^\mu) = 0$ for all $\mu$ [$\sigma_{bb}$ is the best-binding sequence: $E_b(\sigma_{bb}) = E_b^0$], and sample the rest of $\epsilon_b$'s from an exponential distribution defined in the range of $(1, \infty)$ kcal/mol, with 2 kcal/mol mean [36]. Here we consider $L = 5$ hot-spot residues and a reduced alphabet of $k = 8$ amino acids grouped by physicochemical properties, resulting in $8^5 = 32768$ unique sequences. The exact choice of these parameters has little effect on the overall qualitative features of the model.

Our fitness landscape is nonlinear and, thus, epistatic: the fitness effect of a given mutation depends on the entire background sequence [1–3]. However, the landscape is correlated [37] (as $k^L$ sequences are determined by $2Lk$ $\epsilon_f$ and $\epsilon_b$ parameters), and, thus, differs from completely random landscapes [10] in a manner consistent with experimental studies [3,13]. Our model naturally incorporates evolutionary tradeoffs between function and stability [25,26,38], even though binding and folding energetics are uncorrelated [24].

We sample one set of $\epsilon_f$'s and two sets of $\epsilon_b$'s, one for the old binding target and one for the new one. This procedure defines two fitness landscapes, $\mathcal{F}_1$ and $\mathcal{F}_2$, through Eq. (5) ($E_f^0$ and $E_b^0$ are fixed). Initially, each organism in the population has the sequence with maximum fitness under $\mathcal{F}_1$. The population then adapts to binding the new target on $\mathcal{F}_2$. Mutations occur at a rate $LNu \ll (\log N)^{-1}$, where $N$ is the effective population size and $u$ is the mutation rate per amino acid, making the population effectively monomorphic [39]. We assume the strong-selection limit: beneficial mutations are guaranteed to fix, while deleterious and neutral mutations are rapidly eliminated [40]. With Markovian waiting times, the jump probabilities are $\langle \sigma' | \mathbf{Q} | \sigma \rangle = 1/b(\sigma)$ if $\mathcal{F}(\sigma') > \mathcal{F}(\sigma)$ and zero otherwise, where $b(\sigma)$ is the number of beneficial mutations possible from $\sigma$. Note that in this limit our results are independent of $f_0$, and $Nu$ only affects the overall time scale. The path ensemble consists of all adaptive paths (APs) (first-passage paths to local maxima
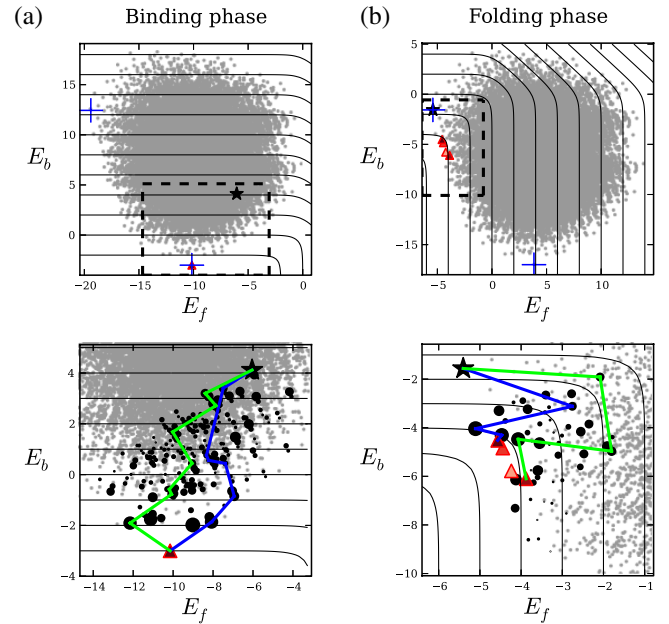


FIG. 1 (color online). Two phases of adaptive protein evolution. (a) Binding phase, with $E_f^0 = -17$ kcal/mol and $E_b^0 = -3$ kcal/mol. (b) Folding phase, with $E_f^0 = -3$ kcal/mol and $E_b^0 = -17$ kcal/mol. Note that $\epsilon_f$'s and the two sets of $\epsilon_b$'s (for $\mathcal{F}_1$ and $\mathcal{F}_2$) are the same in (a) and (b). Top panels show the global distribution of all $8^5$ sequences in energy space according to $\mathcal{F}_2$. Blue crosses indicate the best-folding and best-binding sequences, red triangles correspond to local fitness maxima on $\mathcal{F}_2$ (shaded according to their commitment probabilities), and black stars indicate the initial state for adaptation (global maximum on $\mathcal{F}_1$). Black lines are contours of constant fitness on $\mathcal{F}_2$. In the bottom panels, only the region of energy space accessible to APs (outlined by dashed lines in the top panels) is shown. Representative APs are traced in blue and green; black circles indicate intermediate states along APs, sized according to the AP density $\langle \mathcal{I}_\sigma \rangle_{AP}$; small gray circles are states inaccessible to APs.
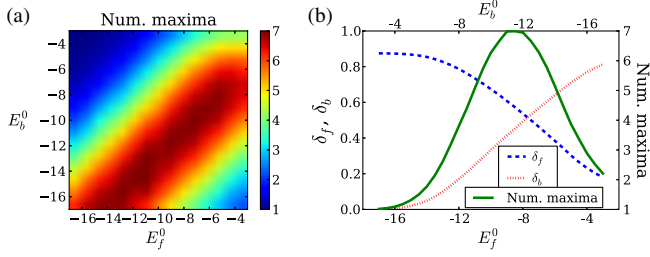
FIG. 2 (color online). (a) Average number of local fitness maxima as a function of the energy offsets $E_f^0$ and $E_b^0$. (b) Average number of local fitness maxima (solid green line), average Hamming distance $\delta_f$ between the maxima and the best-folding sequence (dashed blue line), and average Hamming distance $\delta_b$ between the maxima and the best-binding sequence (dotted red line) for the parameter subspace $E_f^0 + E_b^0 = -20$ kcal/mol. Note that the distance between two random sequences is $1 - 1/k = 0.875$, where $k = 8$ is the alphabet size. All data points are averages over $5 \times 10^3$ realizations; realizations with no adaptation are excluded.

on $\mathcal{F}_2$). Figure 1 shows two examples of $\mathcal{F}_2$ with representative APs.

Since our fitness landscapes [Eq. (5)] are randomly generated, we focus on their generic properties averaged over many realizations of $\epsilon_f$ and $\epsilon_b$ (Figs. 2 and 3). Scans of the $E_f^0 - E_b^0$ parameter space reveal the existence of two qualitatively different phases of adaptation. One, which we call the binding phase, is observed when $E_f^0$ is low and $E_b^0$ is high [see Fig. 1(a) for an example]. In this case, the mean number of local fitness maxima is very low (Fig. 2) and $\delta_f$, the average Hamming distance between these maxima and the best-folding sequence (with the lowest $E_f$), is large [Fig. 2(b)]. In contrast, $\delta_b$, the average Hamming distance to the best-binding sequence, is close to zero. Thus, in this phase the need to bind dominates adaptation.

In the opposite limit (high $E_f^0$ and low $E_b^0$) [see Fig. 1(b) for an example], the folding phase is observed in which the mean number of local maxima is also low (Fig. 2), but these maxima are much closer to the best-folding rather than the best-binding sequence [Fig. 2(b)]. Here, the need to preserve protein stability dominates adaptive dynamics. In the crossover regime between these two phases, there are many local maxima and, therefore, the most epistasis. Epistasis is also reflected in the fact that the fraction of local maxima accessible from the initial state and the probability that the global maximum has the largest commitment probability (i.e., probability that a given maximum is reached from the initial state) are lower, while the fraction of sequence space accessible to APs is higher in this regime compared to the binding and folding phases [Figs. 3(a) and 3(b)]. In the crossover regime, the evolutionary tradeoff between binding and folding alone can result in proteins with marginal folding stability, in contrast with previous hypotheses that explain marginal stability
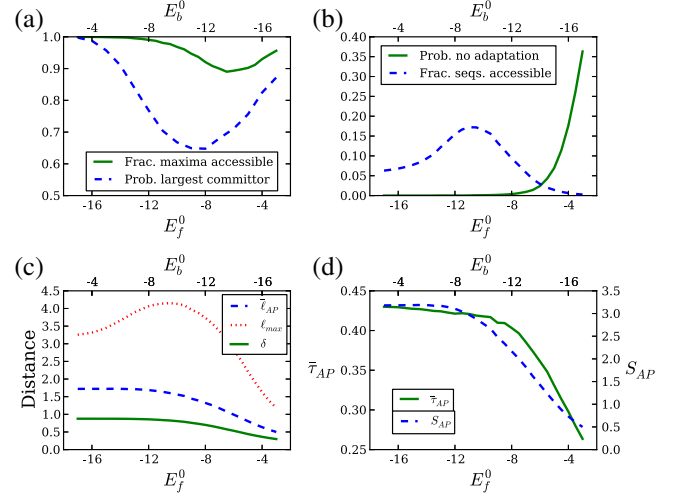


FIG. 3 (color online). (a) Fraction of local fitness maxima accessible from the initial state (solid green line), and probability that the global maximum has the largest commitment probability (committor) among all local maxima (dashed blue line). (b) Probability that the initial sequence starts at a local maximum resulting in no adaptation (solid green line), and fraction of sequence space accessible to APs (dashed blue line). (c) Mean path length $\bar{\ell}_{AP}$ (dashed blue line), maximum possible path length $\ell_{max}$ (dotted red line), and the average net distance $\delta$ between the initial state and final states (solid green line). On average, proteins undergo twice as many substitutions as the net distance $\delta$, and the maximum number of substitutions is three times larger than $\delta$. (d) Mean adaptation time $\bar{\tau}_{AP}$ [in units of $(Nu)^{-1}$] (solid green line), and entropy $S_{AP}$ (dashed blue line). All quantities in (c) and (d) are per-residue. The probability of no adaptation in (b) is an average over $2 \times 10^4$ landscape realizations; all other data points are averages over $5 \times 10^3$ realizations, and realizations with no adaptation are excluded.

with mutational entropy [23] or a fitness function that disfavors hyperstable proteins [41].

On average, paths in the binding phase are longer than those in the folding phase, and adaptation takes more time [Figs. 3(c) and 3(d)]. Paths in the binding phase have higher entropy, indicating that adaptation involves a more diverse set of pathways rather than a few dominant ones. In the folding phase, APs tend to be short since the initial sequence is often either close to, or already at a local maximum [Figs. 3(b) and 3(c)]. A similar situation is observed in directed evolution experiments where the initial sequence already has some affinity for the new ligand but cannot increase it any further [12,42]. In such cases, the sequence must first be mutated away from the local maximum. Furthermore, in the folding phase, folding energy tends to increase at the beginning of paths and decrease toward the end, as a consequence of the distribution of sequences in energy space relative to the fitness contours [Fig. 1(b)]. This is consistent with experiments in which folding stability is sacrificed first and recovered later en route to the new function [26].

Our model can be extended to account for binding-mediated stability, in which binding stabilizes an otherwise disordered protein [43]. We can also incorporate chaperone-assisted folding [44] by modifying $E_f^0$ or the distribution of $\epsilon_f$'s. Furthermore, we can include "folding hot spots" away from the binding interface to see if they acquire stabilizing mutations as a buffer against destabilizing but function-improving mutations at the interface [25,26]. The role of neutral and weakly selected mutations can be studied as well by using substitution rates from more complex population genetics models [39,45], although we expect nonadaptive substitutions to play little role on short time scales. We look forward to studying these extensions in future work.

*morozov@physics.rutgers.edu

[1] D. M. Weinreich, N. F. Delaney, M. A. DePristo, and D. L. Hartl, Science **312**, 111 (2006).

[2] F. J. Poelwijk, D. J. Kiviet, D. M. Weinreich, and S. J. Tans, Nature (London) **445**, 383 (2007).

[3] M. Carneiro and D. L. Hartl, Proc. Natl. Acad. Sci. U.S.A. **107**, 1747 (2010).

[4] F. Noé, C. Schütte, E. Vanden-Eijnden, L. Reich, and T. R. Weikl, Proc. Natl. Acad. Sci. U.S.A. **106**, 19 011 (2009).

[5] P. G. Bolhuis, D. Chandler, C. Dellago, and P. L. Geissler, Annu. Rev. Phys. Chem. **53**, 291 (2002).

[6] D. ben-Avraham and S. Havlin, *Diffusion and Reactions in Fractals and Disordered Systems* (Cambridge University Press, Cambridge, England, 2000).

[7] S. Condamin, O. Bénichou, V. Tejedor, R. Voituriez, and J. Klafter, Nature (London) **450**, 77 (2007).

[8] D. M. Roma, R. A. O'Flanagan, A. E. Ruckenstein, A. M. Sengupta, and R. Mukhopadhyay, Phys. Rev. E **71**, 011902 (2005).

[9] C. H. Waddington, *The Strategy of the Genes. A Discussion of Some Aspects of Theoretical Biology* (Allen and Unwin, London, 1957).

[10] S. Kauffman, *The Origins of Order: Self-Organization and Selection in Evolution* (Oxford University Press, New York, 1993).

[11] T. Enver, M. Pera, C. Peterson, and P. W. Andrews, Cell Stem Cell **4**, 387 (2009).

[12] J. T. Bridgham, E. A. Ortlund, and J. W. Thornton, Nature (London) **461**, 515 (2009).

[13] A. E. Lobkovsky, Y. I. Wolf, and E. V. Koonin, PLoS Comput. Biol. **7**, e1002302 (2011).

[14] J. D. Noh and H. Rieger, Phys. Rev. Lett. **92**, 118701 (2004).

[15] E. M. Bollt and D. ben-Avraham, New J. Phys. **7**, 26 (2005).

[16] C. Dellago, P. G. Bolhuis, F. S. Csajka, and D. Chandler, J. Chem. Phys. **108**, 1964 (1998).

[17] C. Dellago, P. G. Bolhuis, and P. L. Geissler, Adv. Chem. Phys. **123**, 1 (2003).

[18] G. Hummer, J. Chem. Phys. **120**, 516 (2004).

[19] B. Harland and S. X. Sun, J. Chem. Phys. **127**, 104103 (2007).

[20] W. E and E. Vanden-Eijnden, J. Stat. Phys. **123**, 503 (2006).

[21] P. Metzner, C. Schütte, and E. Vanden-Eijnden, Multiscale Model. Simul. **7**, 1192 (2009).

[22] G. H. Weiss, *Aspects and Applications of the Random Walk* (North Holland, Amsterdam, 1994).

[23] K. B. Zeldovich, P. Chen, and E. I. Shakhnovich, Proc. Natl. Acad. Sci. U.S.A. **104**, 16 152 (2007).

[24] N. Tokuriki, F. Stricher, L. Serrano, and D. S. Tawfik, PLoS Comput. Biol. **4**, e1000002 (2008).

[25] N. Tokuriki and D. S. Tawfik, Curr. Opin. Struct. Biol. **19**, 596 (2009).

[26] J. D. Bloom and F. H. Arnold, Proc. Natl. Acad. Sci. U.S.A. **106**, 9995 (2009).

[27] T. A. Whitehead, A. Chevalier, Y. Song, C. Dreyfus, S. J. Fleishman, C. De Mattos, C. A. Myers, H. Kamisetty, P. Blair, I. A. Wilson, and D. Baker, Nat. Biotechnol. **30**, 543 (2012).

[28] P. Hänggi, P. Talkner, and M. Borkovec, Rev. Mod. Phys. **62**, 251 (1990).

[29] See Supplemental Material at http://link.aps.org/supplemental/10.1103/PhysRevLett.111.088102 for details of reaction rate results.

[30] P. Metzner, C. Schütte, and E. Vanden-Eijnden, J. Chem. Phys. **125**, 084110 (2006).

[31] S. Mayer, S. Rüdiger, H. C. Ang, A. C. Joerger, and A. R. Fersht, J. Mol. Biol. **372**, 268 (2007).

[32] T. Clackson and J. A. Wells, Science **267**, 383 (1995).

[33] L. Serrano, A. G. Day, and A. R. Fersht, J. Mol. Biol. **233**, 305 (1993).

[34] N. Tokuriki, F. Stricher, J. Schymkowitz, L. Serrano, and D. S. Tawfik, J. Mol. Biol. **369**, 1318 (2007).

[35] A. A. Bogan and K. S. Thorn, J. Mol. Biol. **280**, 1 (1998).

[36] K. S. Thorn and A. A. Bogan, Bioinformatics **17**, 284 (2001).

[37] L. D. Bogarad and M. W. Deem, Proc. Natl. Acad. Sci. U.S.A. **96**, 2591 (1999).

[38] X. Wang, G. Minasov, and B. K. Shoichet, J. Mol. Biol. **320**, 85 (2002).

[39] M. Manhart, A. Haldane, and A. V. Morozov, Theor. Popul. Biol. **82**, 66 (2012).

[40] J. H. Gillespie, Evolution (Lawrence, Kans.) **38**, 1116 (1984).

[41] M. A. DePristo, D. M. Weinreich, and D. L. Hartl, Nat. Rev. Genet. **6**, 678 (2005).

[42] S. Bershtein, K. Goldin, and D. S. Tawfik, J. Mol. Biol. **379**, 1029 (2008).

[43] C. J. Brown, A. K. Johnson, A. K. Dunker, and G. W. Daughdrill, Curr. Opin. Struct. Biol. **21**, 441 (2011).

[44] S. L. Rutherford, Nat. Rev. Genet. **4**, 263 (2003).

[45] J. F. Crow and M. Kimura, *An Introduction to Population Genetics Theory* (Harper and Row, New York, 1970).

# Supplemental Material:
# Path-Based Approach to Random Walks on Networks Characterizes How Proteins Evolve New Functions

Michael Manhart[1] and Alexandre V. Morozov[1,2]

[1] *Department of Physics and Astronomy, Rutgers University, Piscataway, New Jersey 08854, USA*

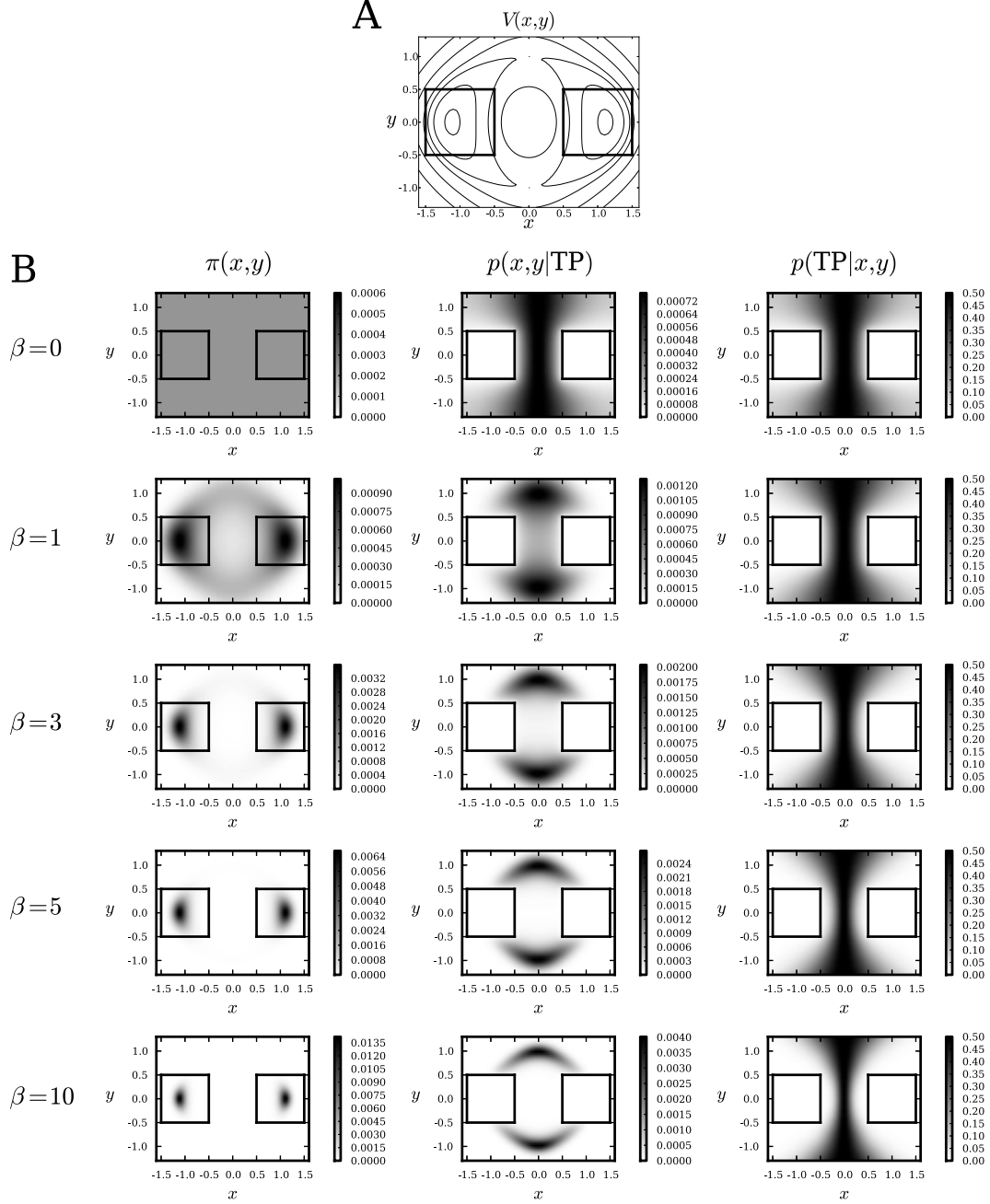[2] *BioMaPS Institute for Quantitative Biology, Rutgers University, Piscataway, New Jersey 08854, USA*

Figure S1: (A) Double-well potential $V(x,y) = \frac{1}{6}(4(1-x^2-y^2)^2+2(x^2-2)^2+((x+y)^2-1)^2+((x-y)^2-1)^2-2)$ defined on a 2D square lattice $\mathcal{S} = [-1.6, 1.6] \times [-1.3, 1.3]$ with spacing $\Delta x = 0.05$. The two metastable states (outlined) are $A = [-1.5, -0.5] \times [-0.5, 0.5]$ and $B = [0.5, 1.5] \times [-0.5, 0.5]$. Jumps between nearest neighbors (nn) have Monte Carlo rates $\langle x', y'|\mathbf{W}|x, y\rangle = (\Delta x)^{-2} \min[1, e^{-\beta(V(x',y')-V(x,y))}]$. Note that time scales are rescaled by $(\Delta x)^2$ so that Brownian dynamics with a fixed diffusion constant is recovered in the $\Delta x \to 0$ limit. Mean waiting times are given by $w(x,y) = (\sum_{\text{nn} (x',y') \text{ of } (x,y)} \langle x', y'|\mathbf{W}|x, y\rangle)^{-1}$, and jump probabilities are $\langle x', y'|\mathbf{Q}|x, y\rangle = w(x,y)\langle x', y'|\mathbf{W}|x, y\rangle$. By definition, the first step of all TP and RP paths is from $\partial A$ or $\partial B$ to a point outside of $A$ and $B$, and the waiting time on $\partial A$ or $\partial B$ is not counted. (B) For the 2D double-well potential, the equilibrium distribution of states $\pi(x,y)$, density of states on TPs $p(x,y|\text{TP}) = \langle \mathcal{T}_{(x,y)}\rangle_{\text{TP}}/\bar{\tau}_{\text{TP}}$, and TP densities (given that the system is at $(x,y)$, the probability it is on a TP: $p(\text{TP}|x,y) = \mathcal{Z}_{\text{TP}}\langle \mathcal{I}_{(x,y)}\rangle_{\text{TP}}/(\mathcal{Z}_{\text{TP}}\langle \mathcal{I}_{(x,y)}\rangle_{\text{TP}} + \mathcal{Z}_{\text{RP}}\langle \mathcal{I}_{(x,y)}\rangle_{\text{RP}}))$ at different values of inverse temperature $\beta$.
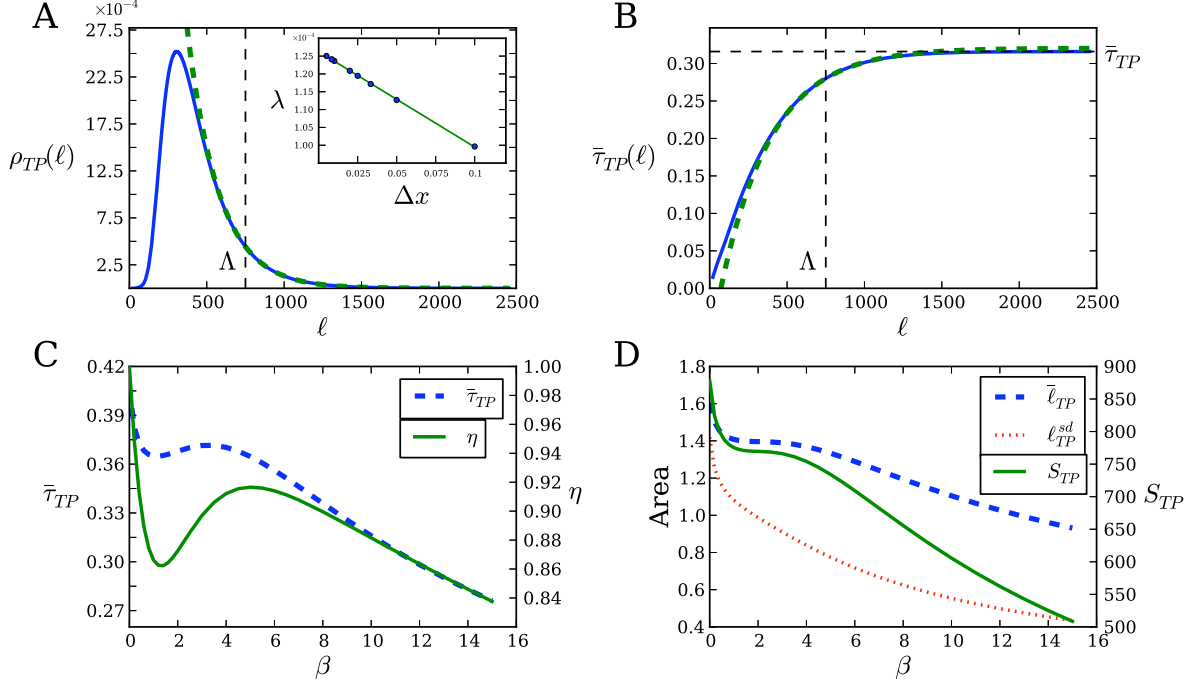
2

Figure S2: (A) For the ensemble of TPs in the 2D double-well potential, path length distribution $\rho_{\text{TP}}(\ell)$ (solid blue line) versus path length $\ell$, and exponential fit in the interval $[\Lambda - 50, \Lambda]$ (dashed green line) are shown. Inset: total TP flux $\lambda$ as a function of lattice spacing $\Delta x$. We approximate the TP flux as the probability of being on a TP divided by the average time of a TP (G. Hummer, J. Chem. Phys. **120**, 516 (2004)): $\lambda \approx p(\text{TP})/\bar{\tau}_{\text{TP}} = (1 - \pi_A - \pi_B)\mathcal{Z}_{\text{TP}}/(\mathcal{Z}_{\text{TP}}\bar{\tau}_{\text{TP}} + \mathcal{Z}_{\text{RP}}\bar{\tau}_{\text{RP}})$, where $\mathcal{Z}_{\text{TP}}, \bar{\tau}_{\text{TP}}$ and $\mathcal{Z}_{\text{RP}}, \bar{\tau}_{\text{RP}}$ are partition functions and average times for transition and return paths, and $\pi_A$ and $\pi_B$ are equilibrium probabilities of $A$ and $B$. The reaction rates are given by $k_{A \to B} = \lambda/(2\pi_A)$ and $k_{B \to A} = \lambda/(2\pi_B)$. The analytical continuation of $\lambda$ is given by $\lambda(\Delta x) = \lambda_0 + \lambda_1 \Delta x + \mathcal{O}(\Delta x^2)$, where $\lambda_0$ is the continuous-limit flux and $\Delta x$ should be smaller then the smallest length scale of the potential. Indeed, $\lambda(\Delta x)$ is linear as seen in the inset, yielding $\lambda_0$ and thus continuous-limit rates $k_{A \to B} = k_{B \to A} \approx 1.3 \times 10^{-4}$. (B) Mean time $\bar{\tau}_{\text{TP}}(\ell)$ of paths up to length $\ell$ (solid blue line). In the $\ell \to \infty$ limit, $\bar{\tau}_{\text{TP}}(\ell)$ converges to the total mean time $\bar{\tau}_{\text{TP}}$. Similar to $\rho_{\text{TP}}(\ell)$, $\bar{\tau}_{\text{TP}}(\ell)$ acquires a universal exponential form for sufficiently large $\ell$: $(\bar{\tau}_{\text{TP}} - \bar{\tau}_{\text{TP}}(\ell)) \sim e^{-a\ell/\bar{\ell}_{\text{TP}}}$, so that a fit in the range $\ell \in [\Lambda - 50, \Lambda]$ (dashed green line) closely matches the full calculation for $\ell > \Lambda$. In (A) and (B), $\Lambda = 750$ is used as the effective length cutoff, and inverse temperature is $\beta = 10$. (C) Relative mean path divergence $\eta = (\mathcal{D}_{\text{TP+RP}}(\beta)/\mathcal{D}_{\text{TP+RP}}(\beta = 0))^{1/2}$ (solid green line) and average time of TPs $\bar{\tau}_{\text{TP}}$ (dashed blue line) versus $\beta$. The divergence $\eta$ is calculated using Eq. (4) with $d(x, y; x', y') = (x - x')^2 + (y - y')^2$. (D) Mean length $\bar{\ell}_{\text{TP}}$ (dashed blue line), standard deviation $\ell_{\text{TP}}^{\text{sd}}$ (dotted red line), and entropy $S_{\text{TP}}$ of TPs (solid green line) versus $\beta$. Note that $\bar{\ell}_{\text{TP}}$ and $\ell_{\text{TP}}^{\text{sd}}$ have units of area since they are rescaled by $(\Delta x)^2$. In (A)–(D), the lattice spacing is $\Delta x = 0.05$ (except for the inset of (A)).
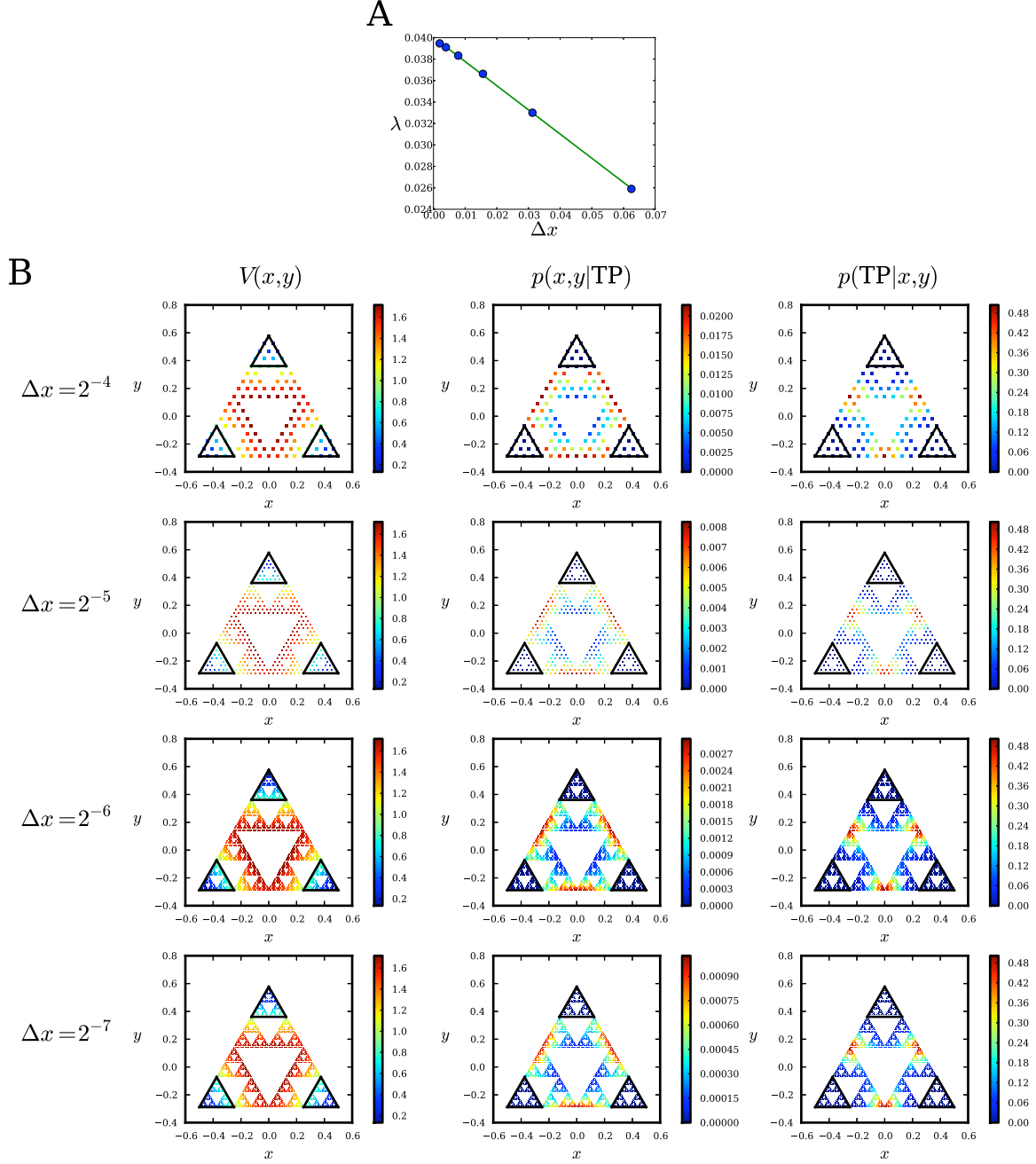
3

Figure S3: Sierpinski triangle embedded in a triple-well potential $V(x,y) = 10\sum_{i=1}^{3}((x - x_i)^2 + (y - y_i)^2)e^{-5(x-x_i)^2-5(y-y_i)^2}$, where $(x_1, y_1) = (0, 1/\sqrt{3})$, $(x_2, y_2) = (1/2, -1/(2\sqrt{3}))$, and $(x_3, y_3) = (-1/2, -1/(2\sqrt{3}))$. There are three metastable states outlined in black, one at each corner of the triangle. Monte Carlo jump rates are as in Fig. S1 but rescaled by $(\Delta x)^{d_w}$, where $\Delta x = 2^{-n}$ ($n$ is the fractal order) and $d_w = \log 5/\log 2$ is the dimension of a random walk on the Sierpinski triangle (D. ben-Avraham and S. Havlin, *Diffusion and Reactions in Fractals and Disordered Systems* (Cambridge University Press, Cambridge, England, 2000)). (A) Transition path flux $\lambda$ as a function of lattice spacing $\Delta x$. As with the double-well potential, analytic continuation of $\lambda(\Delta x)$ allows us to infer the reaction rate $k \approx \lambda/2 \approx 2.0\times10^{-2}$ between any pair of metastable states in an infinite-order fractal using a few finite-order realizations. (B) The potential $V(x, y)$, the density of states on TPs $p(x, y|\text{TP})$, and TP densities $p(\text{TP}|x, y)$ for $\beta = 6$.

4