

## Unfolded protein ensembles, folding trajectories, and refolding rate prediction

A. Das, B. K. Sin, A. R. Mohazab, and S. S. Plotkin

Citation: *J. Chem. Phys.* **139**, 121925 (2013); doi: 10.1063/1.4817215

View online: <http://dx.doi.org/10.1063/1.4817215>

View Table of Contents: <http://jcp.aip.org/resource/1/JCPSA6/v139/i12>

Published by the [AIP Publishing LLC](#).

---

### Additional information on *J. Chem. Phys.*

Journal Homepage: <http://jcp.aip.org/>

Journal Information: [http://jcp.aip.org/about/about\\_the\\_journal](http://jcp.aip.org/about/about_the_journal)

Top downloads: [http://jcp.aip.org/features/most\\_downloaded](http://jcp.aip.org/features/most_downloaded)

Information for Authors: <http://jcp.aip.org/authors>

## ADVERTISEMENT



Explore the **Most Cited**  
Collection in Applied Physics

**AIP**  
Publishing

# Unfolded protein ensembles, folding trajectories, and refolding rate prediction

A. Das, B. K. Sin, A. R. Mohazab,<sup>a)</sup> and S. S. Plotkin<sup>b)</sup>

*Department of Physics and Astronomy, University of British Columbia, 6224 Agricultural Road, Vancouver, British Columbia V6T 1Z1, Canada*

(Received 29 May 2013; accepted 17 July 2013; published online 23 August 2013)

Computer simulations can provide critical information on the unfolded ensemble of proteins under physiological conditions, by explicitly characterizing the geometrical properties of the diverse conformations that are sampled in the unfolded state. A general computational analysis across many proteins has not been implemented however. Here, we develop a method for generating a diverse conformational ensemble, to characterize properties of the unfolded states of intrinsically disordered or intrinsically folded proteins. The method allows unfolded proteins to retain disulfide bonds. We examined physical properties of the unfolded ensembles of several proteins, including chemical shifts, clustering properties, and scaling exponents for the radius of gyration with polymer length. A problem relating simulated and experimental residual dipolar couplings is discussed. We apply our generated ensembles to the problem of folding kinetics, by examining whether the ensembles of some proteins are closer geometrically to their folded structures than others. We find that for a randomly selected dataset of 15 non-homologous 2- and 3-state proteins, quantities such as the average root mean squared deviation between the folded structure and unfolded ensemble correlate with folding rates as strongly as absolute contact order. We introduce a new order parameter that measures the distance travelled per residue, which naturally partitions into a smooth “laminar” and subsequent “turbulent” part of the trajectory. This latter conceptually simple measure with no fitting parameters predicts folding rates in 0 M denaturant with remarkable accuracy ( $r = -0.95$ ,  $p = 1 \times 10^{-7}$ ). The high correlation between folding times and sterically modulated, reconfigurational motion supports the rapid collapse of proteins prior to the transition state as a generic feature in the folding of both two-state and multi-state proteins. This method for generating unfolded ensembles provides a powerful approach to address various questions in protein evolution, misfolding and aggregation, transient structures, and molten globule and disordered protein phases. © 2013 AIP Publishing LLC. [<http://dx.doi.org/10.1063/1.4817215>]

## I. INTRODUCTION

Structural biology has historically been grounded by several landmark studies wherein the atomic coordinates of several large molecules have been experimentally determined, giving insight into the mechanisms of their biological function. Principle relatively recent examples include the photosynthetic reaction center,<sup>1</sup> potassium channels<sup>2</sup> and aquaporins,<sup>3</sup> the ribosome,<sup>4–6</sup> the RNA polymerase II transcription complex,<sup>7</sup> and G protein-coupled receptors.<sup>8,9</sup> Despite the triumphs of the structure-function paradigm, there has been emergent evidence of the biological importance of intrinsically disordered proteins<sup>10–15</sup> for which atomic coordinates significantly fluctuate so that a three-dimensional structure is poorly defined, at least in the absence of binding partners<sup>16</sup> or osmotic stabilizing agents.<sup>17</sup>

While current nuclear magnetic resonance (NMR) measurements of chemical shifts, residual dipolar couplings (RDCs), and <sup>3</sup>J couplings can provide information on structural preferences, computer simulations can provide

critical information of disordered proteins through an explicit geometrical knowledge of the conformational ensemble. There have thus been recent efforts towards a computational characterization of both intrinsically disordered protein (IDP) ensembles<sup>18–23</sup> and chemically denatured ensembles of natively folded proteins.<sup>24,25</sup> One difficulty however is the time-scales necessary to sample a sufficiently large set of conformations to represent microscopic equilibrium properties of the unfolded ensemble.

In what follows, we first describe our method for generating a diverse, representative ensemble for the unfolded state of a protein. We apply this method here to several proteins, including 4 IDPs and 17  $\alpha$ ,  $\beta$ , or mixed natively folded proteins.

As an application of unfolded ensembles, we investigate transformations between unfolded and folded structures, for natively ordered proteins. We ask whether the distance covered during such transformations can predict folding kinetics, and we find several geometrical transformation measures that indeed correlate with folding kinetics for both 2-state and 3-state kinetic folders.

Numerous experiments have pointed to a kinetic molten globule – a semi-collapsed but hydrated state which may have significant secondary structure – as a generic feature of many

<sup>a)</sup>Present address: Recon Instruments #100, 1050 Homer Street, Vancouver, British Columbia V6B 3W9, Canada.

<sup>b)</sup>Electronic mail: [steve@phas.ubc.ca](mailto:steve@phas.ubc.ca)

unfolded proteins in the absence of denaturant (reviewed in Ref. 26). A collapsed ensemble prior to the transition state implies that reconfigurations in the presence of strong bonding interactions and significant steric constraints would be relevant to the folding barrier. Indeed such reconfiguration would be expected to be more substantial if more long range interactions were present, providing a potential explanation for the success of contact order in determining folding rates<sup>27</sup> that is somewhat distinct from explanations involving Flory entropy in loop closure. It is possible, however, that more accurately characterizing such motions could lead to even stronger correlates with folding barriers. We explore this notion and find such a characterization in the distance on average that the protein must travel to adopt the native structure. Taking a portion of the transformation distance where the trajectories become “turbulent” gives a remarkably strong correlation with folding rates in 0 M denaturant.

We organize this paper by first describing our method for generating a conformationally diverse unfolded ensemble,

which applies to both intrinsically folded and intrinsically unfolded proteins. An extension of the method to proteins containing disulfide bonds is described and applied. Properties of the unfolded ensemble, including chemical shifts, residual dipolar couplings, structural diversity, clustering properties, and scaling of the radius of gyration with polymer length are described. We then explore minimal transformations between the unfolded ensemble and the native structure for intrinsically foldable proteins, and investigate the correlation with folding rates of such minimal transformations.

## II. METHODS

In this section, we first describe our method for generating unfolded ensembles. We then describe our method for generating collections of optimized pathways from these unfolded configurations to the folded structure. The proteins analyzed in this paper are given in Table I. They consist of 8 2-state folders, 9 3-state folders, 6  $\alpha$ -helix proteins,

TABLE I. Proteins and their properties used in this study.<sup>a</sup>

PDB	kin <sup>b</sup>	SS <sup>c</sup>	$\log(k_f)$ <sup>d</sup>	$\log(k_u)$ <sup>d</sup>	$\log(k_{mp})$ <sup>e</sup>	N <sup>f</sup>	$\nu$ <sup>g</sup>	ACO	$\langle \text{RMSD} \rangle$ <sup>h</sup>	$\langle \text{TM} \rangle$	$\langle \text{GDT} \rangle$	$\langle \mathcal{D}_{NC} \rangle$	$\langle \mathcal{D} \rangle$	$\langle \mathcal{D}^{(lam)} \rangle$	$\langle \mathcal{D}^{(turb)} \rangle$
1L2Y	2	$\alpha$	12.5	11.5	13	20	0.48	3.7	6.38	0.173	0.462	...	...	...	...
1ENH	2	$\beta$	10.5	7.6	8.1	54	0.57	7.4	14.81	0.147	0.218	14.1	27.0	15.3	11.7
1SHG	2	$\alpha$	1.1	-4.8	-3.7	57	0.6	10.9	16.68	0.133	0.259	16.7	33.5	11	22.5
2CRO	3	$\beta$	3.7	-0.5	0.3	65	0.66	7.3	14.27	0.144	0.218	14.6	30.2	11.7	18.5
1CSP	2	$\beta$	6.5	2.3	2.7	67	0.57	11	18.01	0.129	0.160	18.3	31.5	16.8	14.7
1VII	2	$\alpha$	9.4	5.3	10.6	36	0.57	4	9.68	0.149	0.347	8.2	19.2	13.1	6.1
2PDD	2	$\alpha$	9.8	5.4	9.8	43	0.62	4.8	11.63	0.155	0.302	11	21.3	13.9	7.4
1BNI	3	$\alpha\beta$	2.6	-9.1	-4.3	108	0.57	12.3	21.5	0.127	0.132	22.8	39.8	18.4	21.4
1APS	2	$\alpha\beta$	-1.6	-9	-3.3	98	0.58	21.8	24.36	0.124	0.140	24.8	44.2	17.6	26.6
1A6N	3	$\alpha$	1.1	-3.8	-1.4	151	0.53	14	28.01	0.131	0.113	28.1	45.7	20.9	24.8
1CBI	3	$\beta$	-3.2	-9.8	-6.7	136	0.62	18.8	27.09	0.112	0.112	27.2	45.5	20.1	25.4
1TIT	3	$\beta$	3.6	-7.6	-6.9	89	0.49	15.8	20.39	0.126	0.135	20.4	36.5	19.6	17.0
1IMQ	2	$\alpha$	7.3	-1.9	-1.4	86	0.68	10.4	18.3	0.137	0.203	17	32.7	17.9	14.8
1PSF	3	$\beta$	3.2	...	...	69	0.61	11.7	18.43	0.131	0.254	17.8	...	...	...
2A5E	3	$\alpha\beta$	3.5	0.2	0.4	156	0.52	8.3	23.26	0.127	0.120	24.5	41.1	18.2	22.9
2RN2	3	$\alpha\beta$	0.1	-12	-4.6	155	0.54	19.3	28.69	0.129	0.106	31	43.0	18.5	24.4
1RA9	3	$\alpha\beta$	-2.5	-6.1	-5.2	159	0.52	22.3	24.73	0.140	0.110	27.9	47.8	19.1	28.7
IN <sup>i</sup>	...	...	...	...	...	60	0.57	2.7	12.7	0.143	0.233	...	...	...	...
1IYT	...	...	...	...	...	42	0.58	2.5	10.4	0.143	0.446	...	...	...	...
1XQ8	...	...	...	...	...	140	0.67	2.9	22.1	0.121	0.098	...	...	...	...
proT $\alpha$	...	...	...	...	...	129	0.63	2.5	21.9	0.129	0.130	...	...	...	...

<sup>a</sup>Proteins include Trp-cage miniprotein (1L2Y), Engrailed homeodomain (1ENH), src-homology 3 (SH3) domain (1SHG), phage 434 cro protein (2CRO), cold shock protein (1CSP), chicken villin headpiece (1VII), peripheral subunit-binding domain of dihydrolipoamide acetyltransferase (2PDD), barnase (1BNI), acylphosphatase (1APS), deoxy-myoglobin (1A6N), apo-cellular retinoic acid binding protein I (1CBI), titin, IG repeat 27 (1TIT), colicin E9 immunity protein IM9 (1IMQ), photosystem I protein (1PSF), tumor suppressor P16INK4A (2A5E), ribonuclease H (2RN2), dihydrofolate reductase (1RA9), N-terminal domain of HIV Integrase (IN; IUP structure generated from sequence<sup>28</sup>), amyloid beta-peptide (1-42) (1IYT), alpha-synuclein (1XQ8), Prothymosin alpha (proT $\alpha$ ; IUP structure generated from sequence<sup>28</sup>). 1L2Y and 1PSF were used only for unfolded ensemble generation; 1PSF does not have published unfolding or transition midpoint rates, and 1L2Y is temperature-denatured rather than chemically denatured.

<sup>b</sup>Indicates 2-state or 3-state kinetics.

<sup>c</sup>Secondary structure content.

<sup>d</sup>Natural logarithm of experimentally determined refolding and unfolding rates in 0 M denaturant.<sup>29</sup>

<sup>e</sup>Midpoint experimental relaxation rate.<sup>29</sup>

<sup>f</sup>Chain length.

<sup>g</sup>Scaling exponent of the radius of gyration with chain length  $N$ .

<sup>h</sup>All distance and alignment metrics are averaged over the equilibrium unfolded ensemble, as indicated by angle brackets. RMSD = root mean squared deviation in Å, TM = score from template modeling alignment, GDT = global distance test- total score,  $\mathcal{D}_{NC}$  = Distance between an unfolded conformation and the native, accounting for polymer non-crossing,<sup>30</sup>  $\mathcal{D}$  = Geometrical pathways (GP) generated distance,<sup>31</sup>  $\mathcal{D}^{(lam)}$  = Laminar component of the GP distance,  $\mathcal{D}^{(turb)}$  = Turbulent component of the GP distance. All distance metrics are in units of Å and are an average *per residue*.

<sup>i</sup>For IUP proteins, numbers for ACO indicate an average over all structures in the unfolded ensemble. Numbers for  $\langle \text{RMSD} \rangle$ ,  $\langle \text{TM} \rangle$ , and  $\langle \text{GDT} \rangle$  indicate an average over all pairs in the unfolded ensemble.

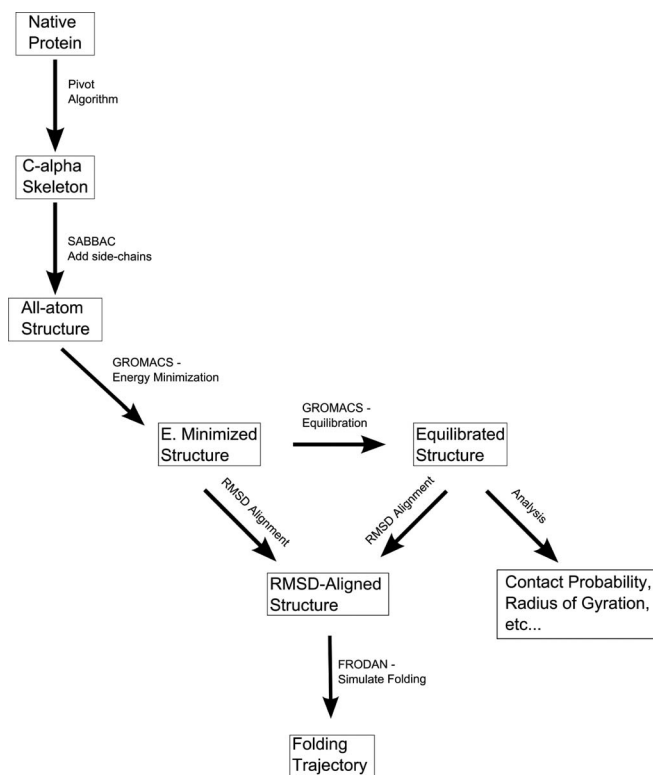


FIG. 1. Overview of algorithm.

6  $\beta$ -sheet proteins, 5 mixed  $\alpha/\beta$  proteins, and 4 intrinsically disordered proteins. This dataset is not large, but it is diverse, spanning a wide range of rates, size, and structural classes. Native structural homologs as defined through TM-score were not included in the dataset; otherwise, no additional pruning or selection of proteins was made.

### A. Generating diverse ensembles of unfolded configurations

We generate unfolded ensembles by employing the following method:

- Generate a diverse coarse-grained (CG) ensemble
- “Foliate” each structure by adding backbone/side-chain degrees of freedom
- Equilibrate each foliated structure for a short time.

The steps in the method are shown schematically in Figure 1. First, if the native crystal or NMR structure existed, it was used as a starting point. If a protein or peptide was intrinsically unfolded or if no pdb file was available, the proteins initial structure was generated *a priori* by submitting the sequence to the I-TASSER structure prediction server,<sup>28</sup> then minimizing and equilibrating the structure. In all cases, the initial structure was used solely as a starting point for the structure generation algorithm.

#### 1. Pivot and crankshaft moves

We coarse grain the initial structure by retaining only the  $C_\alpha$  coordinates. This CG structure is then altered by employing a generalization of the pivot algorithm,<sup>19,30,32–34</sup> an effi-

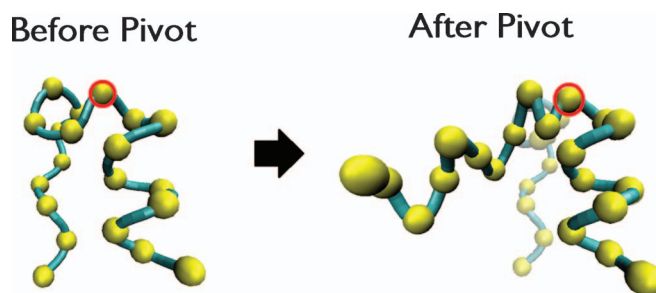


FIG. 2. Illustration of an example pivot move for PDB 1L2Y.

cient algorithm for generating self-avoiding random walk ensembles. We have previously implemented this generalization to generate CG  $C_\alpha$ -model unfolded ensembles.<sup>30</sup> For the CG structure, a pivot move selects a particular bond angle and its corresponding dihedral angle at random, and then resamples them from a native-centric Boltzmann distribution<sup>30</sup> (see Figure 2). The corresponding phenomenological energy function contains approximate angle and dihedral stiffness parameters. If after a pivot move the chain sterically interferes with itself, the move is discarded. For a chain of length  $N$ , pivot moves are repeatedly attempted until  $\mathcal{O}(N)$  successful pivot moves are implemented, such that on average, one pivot move per residue is obtained.

We also considered generated ensembles for proteins with disulfide bonds present. In this case, the pivot algorithm cannot be directly implemented because the disulfide constraint correlates the position of two remote parts of the chain. In this case, we implemented an alternative move for residues bounded by those participating in the disulfide bond. We implemented this procedure for human superoxide dismutase (SOD1), a 153 aa protein that contains a disulfide bond between C57 and C146.

We start by picking a residue at random. If, for a disulfide-bonded protein such as SOD1 (1HL5) (Fig. 3), the residue is before 57 or after 146, a pivot move is implemented for the part of the chain N-terminal or C-terminal to the selected residue. This preserves the coordinates of the disulfide loop. If, however, the selected residue is between 57 and 146, it lies inside the disulfide loop; in this case a non-local “crankshaft” move is implemented between the selected residue and another randomly chosen residue inside the

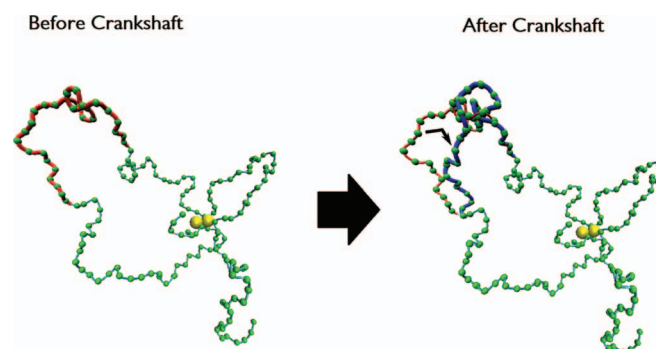


FIG. 3. Crankshaft move for SOD1 (PDB 1HL5), a protein with a long-range disulfide bond between C57 and C146. A minimized, non-equilibrated configuration is shown.



disulfide loop. This “large step” move, which is similar in spirit to the “small-step” biased gaussian steps in torsional space developed by Favrin *et al.*,<sup>35</sup> randomizes the dihedral angles between the two residues, subject to the constraint that the end point  $C_\alpha$  coordinates are fixed. An example crankshaft move is shown in Figure 3. Fixed end-point moves have been studied in detail previously by several authors.<sup>35–39</sup> One advantage of the method we employ here is that the moves result in large global changes of an arbitrarily large collection of contiguous residues. A disadvantage, however, is that the algorithm in its current form is not fast. The trade-off between speed and global reconfigurational change was acceptable for our methodology.

Consider two randomly selected  $C_\alpha$  atoms  $i, j$  with  $j > i$ . The deviation in end point position  $\delta \mathbf{r}_j$  may be Taylor expanded as a function of the  $n = j - i - 1$  dihedral angles in the coarse-grained structure between  $i$  and  $j$  whose rotations would alter the position  $\mathbf{r}_j$

$$\delta \mathbf{r}_j = \sum_{k=1}^n \frac{\partial \mathbf{r}_j}{\partial \phi_k} \delta \phi_k, \quad (1)$$

so that the square of the end point position is given by

$$\delta \mathbf{r}_j^2 = \sum_{k,\ell=1}^n \delta \phi_k G_{k\ell} \delta \phi_\ell, \quad (2)$$

where the  $n \times n$  non-negative, symmetric matrix  $\mathbf{G}$  has elements

$$G_{k\ell} = \frac{\partial \mathbf{r}_j}{\partial \phi_k} \frac{\partial \mathbf{r}_j}{\partial \phi_\ell}. \quad (3)$$

The 3 constraints corresponding to  $\delta \mathbf{r}_j = 0$  means that 4 or more dihedrals are required to have a mode corresponding to an eigenvector of  $\mathbf{G}$  with zero eigenvalue, i.e.,  $j - i = 5$  or greater. However, for  $j - i \geq 5$ , there will be combinations of dihedrals that leave  $\mathbf{r}_j$  unchanged. We randomly select one of these eigenvectors and implement a small rotation  $\delta \phi$  along it. The end point position generally moves slightly because Eq. (2) is only zero for infinitesimal displacements. We repeatedly implement multiple rotation vectors  $\delta \phi$  that leave the end point position  $\mathbf{r}_j$  nearly unchanged. Occasionally, the deviation in the end point position becomes appreciable (fractions of an Å). We then correct the end point deviation  $\delta \mathbf{r}_j$  by inverting Eq. (1) to find the vector  $\delta \phi$ , with components

$$\delta \phi_k = - \left( \frac{\partial \mathbf{r}_j}{\partial \phi_k} \right)^{-1} \cdot \delta \mathbf{r}_j. \quad (4)$$

This gives a set of dihedral rotations that rotates the end point back to its original position. In total, we implement rotations roughly 20 times per bond angle, i.e.,  $20(j - i)$  times, correcting as needed. The net result is to change the coordinates in a crankshaft fashion, as shown in Figure 3. This constitutes one update of the configuration. As mentioned above,  $\sim N$  updates are taken before sampling a new distinct configuration for foliation by adding the remaining side-chain and backbone atoms. We have illustrated non-local crankshaft moves here for the case of a disulfide bonded protein. However, for the remainder of the analysis we consider only non-disulfide-bonded proteins. In principle, one could implement a combination of pivot moves and crankshaft moves for all proteins. Because our unfolded ensembles were not dense, it was more efficient to implement pivot moves; for dense systems however such as polyglutamine repeats for example, it is an interesting future topic to explore combined move sets.

## 2. Foliation, minimization, and equilibration

Once a distinct CG conformation is obtained by pivot/crankshaft moves, side-chains and backbone are added. We examined two methods of adding atoms, PULCHRA<sup>40</sup> and SABBAC,<sup>41</sup> which yielded equivalent results; PULCHRA is available as an executable program, while calls to SABBAC must be uploaded to a server – in an automated way for a large number of configurations.

After side-chain and backbone atoms are added, the protein is energy minimized in GROMACS using a steepest descent algorithm to eliminate steric clashes; Figure 4 shows a rendering of the process for Trp cage (PDB 1L2Y). Though the configuration itself does not change much (RMSD values are approximately 1.4 Å), the energy typically decreases by several orders of magnitude. The radius of gyration typically increases slightly during this process – zero temperature energy minimization tends to favor extended chains when starting from a random initial structure.

The SABBAC and PULCHRA algorithms occasionally place side chains in sterically clashing positions that are not ameliorated by minimization. Roughly 15% of the initial states are not viable for equilibration and are discarded.

For disulfide bonded proteins, the residues involved in the disulfide bond, along with the adjacent residues in the primary sequence (e.g., residues 56–58 and 145–147 for SOD1), are all held fixed during crankshaft and pivot moves. This allows the side-chains of the disulfide bond (between residues 57 and 146 for SOD1) to be reconstructed from the initial configuration. Essentially, no conformations are lost by this

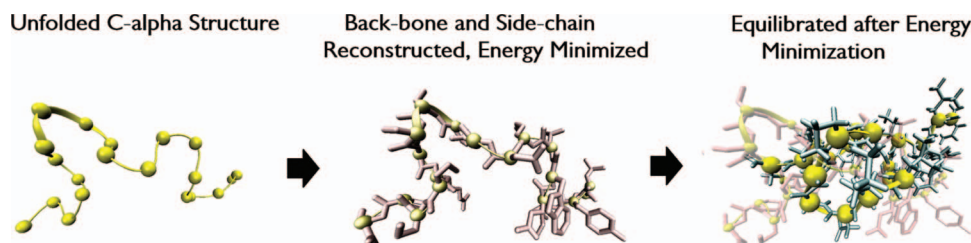


FIG. 4. Schematic figures indicating the processes of backbone and side-chain addition, energy minimization, and 1 ns thermal equilibration.

constraint after equilibration – the arbitrary origin of the coordinates can be thought of as centered around the disulfide bond. PULCHRA and SABBAC may also occasionally protonate histidines on a nitrogen that is likely unprotonated in the native structure. This may be corrected by applying a patch to properly protonate histidines in the unfolded conformational ensemble. We have generally not implemented such a patch, except for the ensemble of SOD1 conformations that are provided in the supplementary material.<sup>47</sup>

Viable all-atom configurations were then thermally equilibrated by molecular dynamics (MD) in explicit SPC water with CHARMM27 force field, using the GROMACS simulation package. All protein atoms are initially at least 20 Å from the faces of a cuboid aligned with the axes of the gyration tensor and having periodic boundary conditions. Simulation conditions were in the *NPT* ensemble at  $T = 300$  K using the modified Berendsen (V-rescale) weak coupling thermostat, and 1 bar using the Parrinello-Rahman barostat. The Particle Mesh Ewald method was applied for long-range electrostatics, and a 10 Å cut-off was used for non-bonded electrostatic and van der Waals (VDW) interactions. Covalent bonds lengths to hydrogen atoms were constrained using the LINCS algorithm. The integration time-step was 2 fs, and coordinates were saved every 100 ps. Each initial configuration is simulated generally for 1 ns. For  $\alpha$ -synuclein, proT $\alpha$ , and SOD1, configurations were simulated for 5 ns. A sample of 1000 equilibrated conformations from the unfolded ensemble of the disulfide-bonded protein superoxide dismutase (WT SOD1) is provided in the supplementary material.<sup>47</sup>

Several previous approaches have been used to generate unfolded ensembles. Zagrovic and Pande<sup>42</sup> have run thousands of independent MD trajectories in implicit GB/SA solvent for 3 small proteins: Villin, Trp zipper, and BBA5. Each of these simulations all started from the fully extended state ( $\phi = -135^\circ$ ,  $\psi = 135^\circ$ ), mandating individual simulations at least 10 ns long for convergence of quantities such as the radius of gyration. One advantage we found for generating random initial ensembles was that many global properties such as the radius of gyration often came to equilibrium on the 1 ns time scale. Pappu and colleagues<sup>19</sup> have run 50 iterations of 12 ns simulations for 90 randomly generated initial configurations of Gln<sub>5</sub> and Gln<sub>15</sub>, in explicit TIP4P solvent. Their procedure is similar to the one we employed here, except that there was no coarse-graining step, as we have employed. Other methods of generating unfolded ensembles have involved extracting conformational fragments from databases of high-resolution crystal structures. Sosnick and colleagues<sup>24</sup> find the conformation of an amino acid (accounting for nearest neighbor effects) using a statistical potential for regions outside of helices, sheets, and turns. If there is a steric clash, the conformation is nudged until the clash is resolved. This process certainly perturbs the ensemble from the equilibrium one,<sup>34,43,44</sup> and as well the resulting conformations are not subsequently equilibrated. Still, the model shows excellent agreement with experimental RDCs, for the experimental scenario of chemically denatured proteins in highly anisotropic media. Blackledge and colleagues<sup>18,25</sup> employ a similar approach of extracting conformational fragments from a database. Chains are grown by randomly selecting  $\phi/\psi$  an-

gles from a database of sequence-specific fragments outside of helices or sheets (turns are now included). If in the growth of the chain a steric clash ensues, those angles are rejected and another angle pair is selected until no steric clash is found. This process also perturbs the ensemble from the equilibrium one,<sup>34,43,44</sup> and as above the conformations are not equilibrated. Nevertheless, this model also shows excellent agreement with experimental RDCs in conditions of high denaturant and high anisotropy.

## B. Chemical shifts and residual dipolar couplings

Chemical shifts were obtained using the program CAMSHIFT.<sup>45</sup> Residual dipolar couplings were obtained using the program PALES.<sup>46</sup> Chemical shifts and RDCs for the proteins in this study are tabulated in the supplementary material.<sup>47</sup>

## C. Investigating the correlation between geometrical folding pathways and folding kinetics

We seek direct transformations between an unfolded configuration and the native configuration, to give a measure of average distance between the unfolded ensemble and the folded structure. We considered two generated ensembles: one energy minimized and equilibrated for 1 ns, and another only energy minimized. We quantified the unfolded-folded distance in several ways:

- RMSD to the native structure: calculated for  $C_\alpha$  atoms, and then averaged over the unfolded ensemble, which we denote by  $\langle RMSD \rangle$ .<sup>48</sup>
- Ensemble-averaged TM-score:  $\langle TM\text{-score} \rangle$
- Ensemble-averaged global distance test – total score:  $\langle GDT\text{-TS} \rangle$
- The distance accounting for polymer non-crossing constraints ( $\langle D_{NC} \rangle$ ), as calculated by the algorithm developed by Mohazab and Plotkin.<sup>30</sup>
- The distance corresponding to the RMSD-minimized trajectories generated by the Geometrical Pathways (GP) algorithm of Farrell *et al.*<sup>31</sup>
- A variation of the GP distance where only the smooth “laminar” part of the trajectories are taken, as described below.
- A variation of the GP distance where only the fluctuating “turbulent” part of the trajectories are taken, as described below.
- A variation of the GP distance, wherein trajectories are smoothed by applying a weighting function given by  $(\frac{1}{3}, \frac{2}{3}, 1, \frac{2}{3}, \frac{1}{3})$  to consecutive sets of 5 points along the trajectory. This procedure eliminates jagged edges along the trajectory.

### 1. Order parameters for unfolded structures

The ensemble-averaged order parameters described in this section are all given in Table I for the proteins in this study.

To examine RMSD, structures in the unfolded ensemble were RMSD-aligned to the native structure, and the average (residual) RMSD was calculated, using the software program VMD.<sup>49</sup>

TM-score was calculated with TM-ALIGN,<sup>50</sup> and is given by

$$\text{TM-score} = \max \frac{1}{L_{\text{targ}}} \sum_{i=1}^{L_a} \left[ 1 + \left( \frac{d_i}{d_o(L_{\text{targ}})} \right)^2 \right]^{-1},$$

where  $L_{\text{targ}}$  is the length of the target protein that another protein structure is aligned to,  $L_a$  is the number of template-aligned  $C_\alpha$  pairs,  $d_i$  is the distance between the  $i$ th pair of aligned pairs, and  $d_o(L_{\text{targ}}) = 1.24(L_{\text{targ}} - 15)^{1/3} - 1.8$  is a distance parameter that ensures that the average TM-score is not dependent on protein length. max indicates the maximization of this quantity by alignment.

GDT-TS<sup>51</sup> was calculated by file upload to the K<sub>0</sub>BaMIN web server,<sup>52</sup> and is calculated by

$$\text{GDT-TS} = \max \frac{1}{4N} (C_{d/4} + C_{d/2} + C_d + C_{2d}),$$

where  $N$  is the chain length,  $d = 4 \text{ \AA}$  is a distance threshold, and  $C_{d/4}$ , for example, is the number of residues superposed below a threshold of  $d/4$  after alignment.

Distance  $\mathcal{D}_{NC}$  accounting for polymer non-crossing was calculated using the method in Ref. 30, which is based on the calculation of minimal distance trajectories that we have developed previously.<sup>53–55</sup> This method calculates the approximate distance undertaken by all  $C_\alpha$  atoms in transforming between two structures, e.g., an unfolded structure and the native structure, while accounting for polymer non-crossing constraints. The method involves a depth first tree search algorithm to find the shortest distance trajectories between two conformations for a linear self-avoiding polymer. We apply this method here between the native structure and 200 coarse-grained unfolded structures. After coarse-graining (smoothing) conformations by sampling every other bead, each structure was transformed to the folded state by the algorithm discussed in Ref. 30, and the minimal distance cost was found.

The geometrical pathways<sup>31</sup> distance  $\mathcal{D}$  was found using the program FRODAN, which calculates a stereochemically acceptable transformation between two all-atom structures, by following a steepest descent pathway that minimizes RMSD. From such transformations, we calculate the distance that  $C_\alpha$  atoms have moved. Sample folding trajectories of 5  $C_\alpha$  atoms have moved. The total arc length that all  $C_\alpha$  residues have travelled is accumulated, to obtain the total distance for one conformation pair. Similarly to TM-score and GDT-TS, this total is then divided by chain length, which yields the mean distance travelled per residue, for one conformation pair,

$$\mathcal{D} = \frac{1}{N} \sum_{i=1}^N \int_{r_i^{(\alpha)}}^{r_i^{(N)}} |d\mathbf{r}_i|. \quad (5)$$

Here, the sum is over the  $N$   $C_\alpha$  atoms, and the integral sums up arc-length increments from initial to final position for each  $C_\alpha$  atom (see Figure 5). For the proteins in Table I, the mean distance is about 36  $\text{\AA}$ .

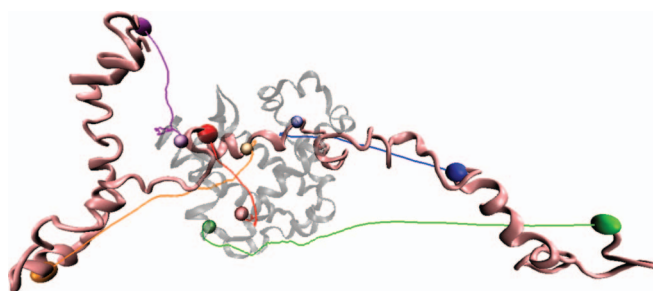


FIG. 5. Example optimal folding trajectories for 5  $C_\alpha$  atoms in apo-myoglobin (1A6N). Unfolded and folded structures are also shown.

Particle trajectories obtained from the GP method tended to be delineated by two parts, an early smooth “laminar” segment, and a late rugged “turbulent” segment (Figs. 6(a) and 6(b)). Figure 6(c) plots the distance travelled per step as a function of step index along a  $C_\alpha$  trajectory. The turbulent

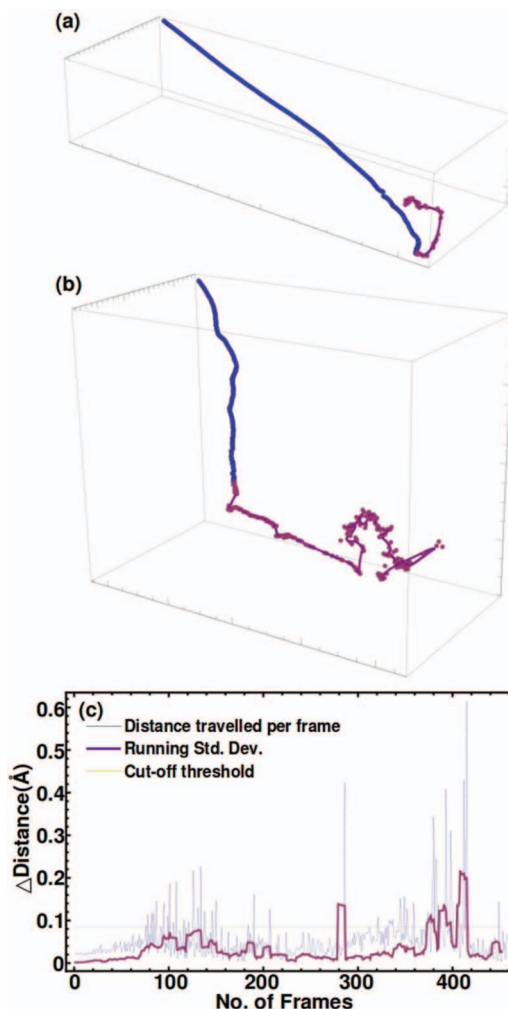


FIG. 6. Each  $C_\alpha$  trajectory is divided into a smooth “laminar” and rugged “turbulent” part. Panels (a) and (b) show sample trajectories for  $C_\alpha(4)$  and  $C_\alpha(75)$  of apo-myoglobin. Panel (a) is predominantly laminar – the corresponding distances are  $\mathcal{D}^{(\text{lam})} = 51 \text{ \AA}$ ,  $\mathcal{D}^{(\text{turb})} = 7 \text{ \AA}$ . Panel (b) is predominantly turbulent – the corresponding distances are  $\mathcal{D}^{(\text{lam})} = 4.6 \text{ \AA}$ ,  $\mathcal{D}^{(\text{turb})} = 12 \text{ \AA}$ . (c) Criterion for determining the transition from laminar to turbulent trajectories. When the root variance in the distance travelled per step jumps above a threshold given by 7 times the baseline value, the trajectory from then on is defined as turbulent.



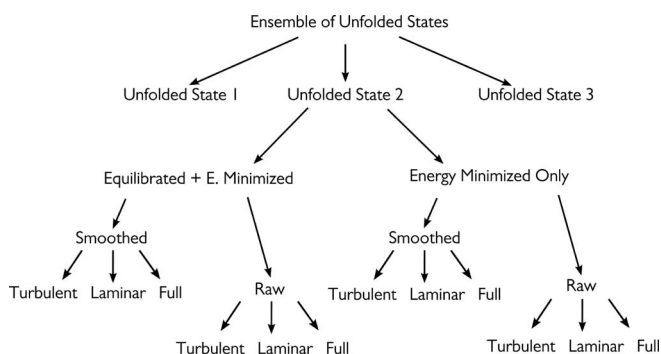


FIG. 7. Different ensembles considered in this study to compare with protein folding kinetics.

segment is characterized by heterogeneous jumps of variable distance. The root variance (standard deviation) in step length, averaged over 10 frames, is also plotted. We determine a baseline standard deviation by averaging the first 50 frames, then when the standard deviation exceeds a threshold which we set to 7 times the baseline, the trajectory is deemed turbulent from then on. We thus partition each trajectory up into laminar and turbulent parts, and define the corresponding accumulated distances from the unfolded structure ( $\alpha$ ) to the native structure ( $N$ )

$$\mathcal{D}^{(lam)} = \frac{1}{N} \sum_{i=1}^N \int_{r_i^{(\alpha)}}^{r_i^{(cut)}} |dr_i|,$$

$$\mathcal{D}^{(turb)} = \frac{1}{N} \sum_{i=1}^N \int_{r_i^{(cut)}}^{r_i^{(N)}} |dr_i|. \quad (6)$$

The transition to turbulence can also be signalled by an abrupt increase in the curvature of the trajectory as shown in Figure S1 of the supplementary material.<sup>47</sup>

Choosing either the minimized or equilibrated ensembles, the smoothed or raw trajectories, and the laminar, turbulent, or full trajectories, gives 12 different measures of distance, as depicted in Figure 7.

Both 2 and 3 state proteins were selected from proteins with known kinetics using the webserver KineticDB.<sup>29</sup> Proteins were ensured to be non-homologous by TM-score, as compared to a non-redundant (NR) protein database<sup>56</sup> (see Fig. 8).

### III. RESULTS AND DISCUSSION

#### A. Chemical shifts and residual dipolar couplings

Chemical shift values for  $C_\alpha$  atoms are obtained using the program CAMSHIFT.<sup>45</sup> These values are plotted for  $A\beta_{1-42}$  in Figure 9, which show good agreement with the experimental values of Hou *et al.*<sup>57</sup> Chemical shift values were obtained from a generated ensemble of 773 structures. These numbers agree with those from 1  $\mu$ s explicit water simulations.<sup>22</sup> That said, even better agreement with experimental values is obtained from chemical shifts compiled from databases of loop regions in protein structures (e.g., CAMCOIL,<sup>58</sup>  $r = 0.99$ ).

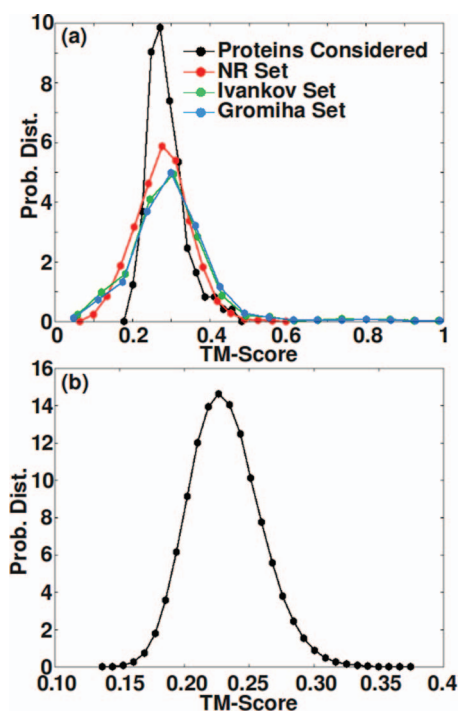


FIG. 8. (Panel (a)) TM-score distributions between *native* structures, showing homology of our dataset compared to a NR dataset,<sup>56</sup> and other datasets used for protein folding kinetics analysis.<sup>84,118</sup> One can see some homologous protein pairs in other datasets. (Panel (b)) TM-score distribution between 1299 unfolded states for  $\alpha$ -synuclein. Similar distributions are obtained for other proteins.

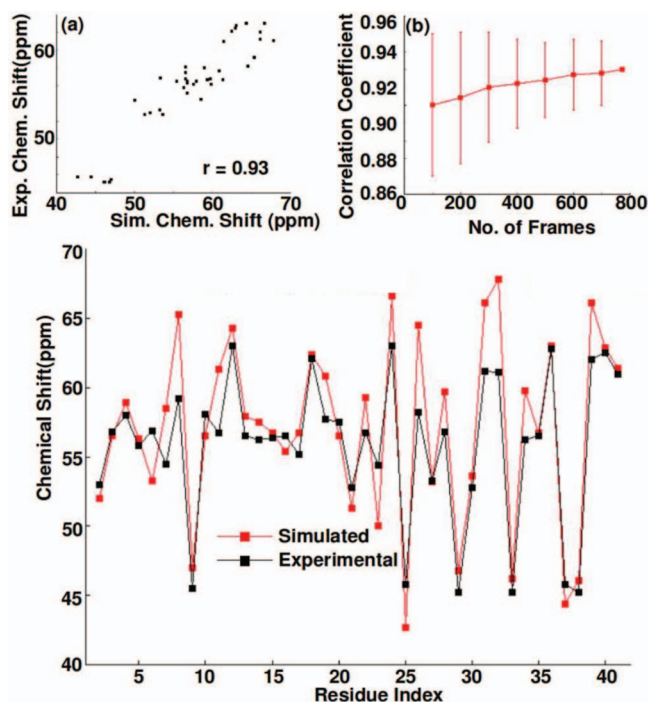


FIG. 9. Comparison between experimental and simulated  $^{13}C_\alpha$  chemical shift values, for  $A\beta_{1-42}$ . (Main panel) Black data points are experimental values from Ref. 57, red data points are those from the simulated ensemble of 773 conformations, using CAMSHIFT. (Inset (a)) Scatter plot of experimental vs simulated chemical shifts ( $r = 0.93$ ). (Panel (b)) Convergence study of the correlation coefficient between experimental and simulated data. Mean correlation coefficient is shown; vertical bars indicate the standard deviation of correlation coefficient values when random subsets with a given number of frames are taken from the total dataset.



RDCs measure the nuclear dipole coupling between two spin 1/2 nuclei such as the amide H and labelled  $^{15}\text{N}$ . This coupling between nuclei  $i$  and  $j$  depends on the angle the internuclear bond makes with the orientation of the whole molecule, and is given in the principle axis frame by

$$D_{ij}(\theta, \phi) = \frac{D_{max}}{2} \left[ S_{zz}(3 \cos^2 \theta - 1) + (S_{xx} - S_{yy}) \sin^2 \theta \cos 2\phi \right], \quad (7)$$

where  $S_{zz}$  is the axial component of the alignment tensor, and  $(S_{xx} - S_{yy})$  is the rhombic component of the alignment tensor, and  $|S_{zz}| > |S_{yy}| \geq |S_{xx}|$ . The angles  $\theta$  and  $\phi$  give the orientation of the vector in the principle basis, and  $D_{max}$  sets the scale of dipolar interactions and is given by  $D_{max} = \gamma_i \gamma_j \mu_o h / (8\pi^3 \bar{r}_{ij}^3)$ , where  $\bar{r}_{ij}$  is the effective internuclear distance accounting for libration of internuclear vector,<sup>59</sup>  $\gamma_i, \gamma_j$  are gyromagnetic ratios of nuclei  $i, j$ ,  $\mu_o$  is the magnetic permeability of vacuum, and  $h$  is Planck's constant.

RDCs were obtained from simulated ensembles using PaLes;<sup>46</sup> we have provided tables of amide NH RDC values for our generated ensembles in the supplementary material.<sup>47</sup> While simulated RDCs of native or near-native structures show good correlation with experimental RDCs, simulated RDCs of unfolded ensembles have not shown strong correlation with experimental RDCs in previous studies;<sup>60</sup> as well, we also did not see significant correlation with experimental values.

There are many potential reasons for this. Experimental RDC values change significantly depending on the degree of extension of the orienting liquid crystal (e.g., polyacrylamide), and also vary with denaturant concentration and pH. The chemically denatured ensemble may significantly differ from the unfolded ensemble. Unfolded ensembles may show partial native structure,<sup>61</sup> and at least partial collapse,<sup>62,63</sup> depending on net charge or hydrophobicity. Urea or GuHCl-denatured ensembles on the other hand tend to exhibit relatively simple self-avoiding walk behavior.<sup>64–66</sup> Finally, because steric obstacles represent absorbing boundaries for the probability distribution of a disordered polymer, the effect of steric hindrance itself, due to the aligning liquid crystals, significantly modifies the structure of the unfolded ensemble. Thus, the experimental conditions under which RDCs are obtained may result in very different ensembles than those in simulation, which model an isolated protein in the absence of denaturant at  $\text{pH} \sim 7$ .

In spite of this, several recent models of the unfolded state show remarkable agreement with experimental RDCs in denaturant and in stretched polyacrylamide.<sup>18,24,25</sup> These models reconstruct the unfolded ensemble from unstructured elements in the protein data bank, and evidently are very good models of unfolded proteins in aligning media. It is not clear on the other hand how such models would predict RDCs as external conditions were varied, e.g., as the polyacrylamide were relaxed or compressed, without some phenomenological adjustment. It may be best to think of the RDC values we obtain here as “unperturbed” RDCs arising solely due to cor-

relations between bond vectors and the inherent anisotropy of the polymer, rather than the induced anisotropy.

## B. Polymer scaling laws and persistence length

After equilibration, the radius of gyration  $R_G$  was obtained for all subsequences with length  $\leq N$ , for each protein in Table I. The slope of the log-log plot gives the exponent of the scaling law  $R_G = r_1 N^\nu$ , where in 3-dimensions  $\nu = 3/5$  for a self-avoiding random walk,  $1/2$  in the  $\Theta$ -state, and  $1/3$  for a compact globule state.<sup>67</sup> Figure 10(a) shows a plot of the radius of gyration vs timestep for several trajectories along with the average over trajectories, for the 129 aa intrinsically unfolded protein (IUP) proT $\alpha$ . Figure 10(b) shows the scaling law for the radius of gyration obtained by taking subsequences of the full length protein and averaging  $R_G$  for those lengths over the 5 ns equilibrium ensemble.

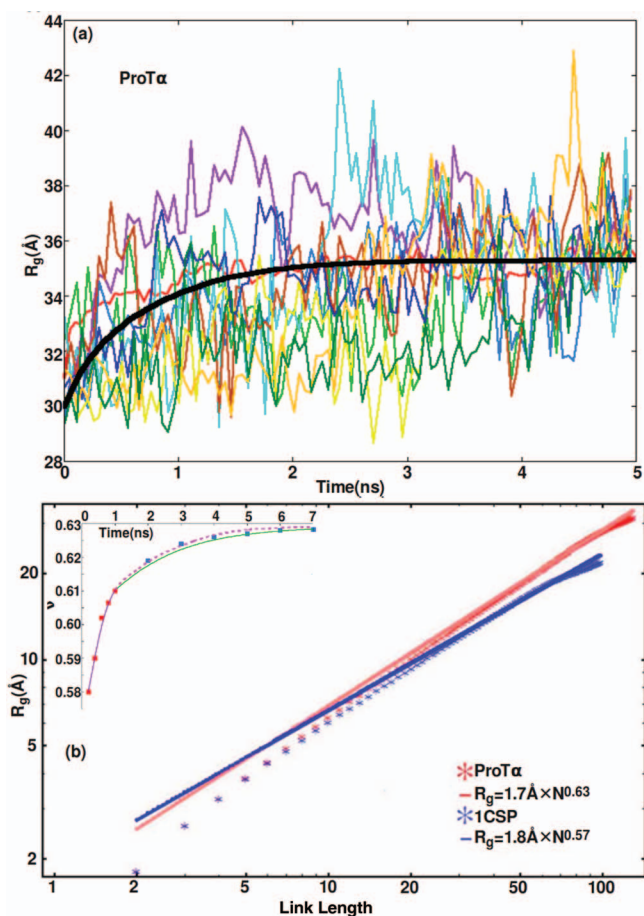


FIG. 10. (a) Radius of gyration vs. time (equilibration process), for proT $\alpha$ : a highly charged, intrinsically disordered protein. The relaxation time is about 0.8 ns, and the asymptotic value of the radius of gyration  $R_G$  is about 35.5 Å. (b) Scaling of the radius of gyration  $R_G$  with chain length, obtained by taking all subsections of a given length and finding the ensemble averaged radius of gyration. (Inset) Extrapolation procedure to find the asymptotic value of the scaling exponent  $\nu$ . The value of  $\nu$  is obtained for ensembles at a given equilibration time. This value converges exponentially to the  $t \rightarrow \infty$  value. Extrapolation from ensembles with  $t \leq 1$  ns gives an asymptotic value of 0.633, while extrapolation from ensembles with  $t \leq 5$  ns gives an asymptotic value of 0.631. A similar conclusion was obtained from extrapolation of the data for  $\alpha$ -syn. Thus, extrapolation of  $\nu$  from  $t \leq 1$  ns ensembles is likely to be sufficiently accurate in general.

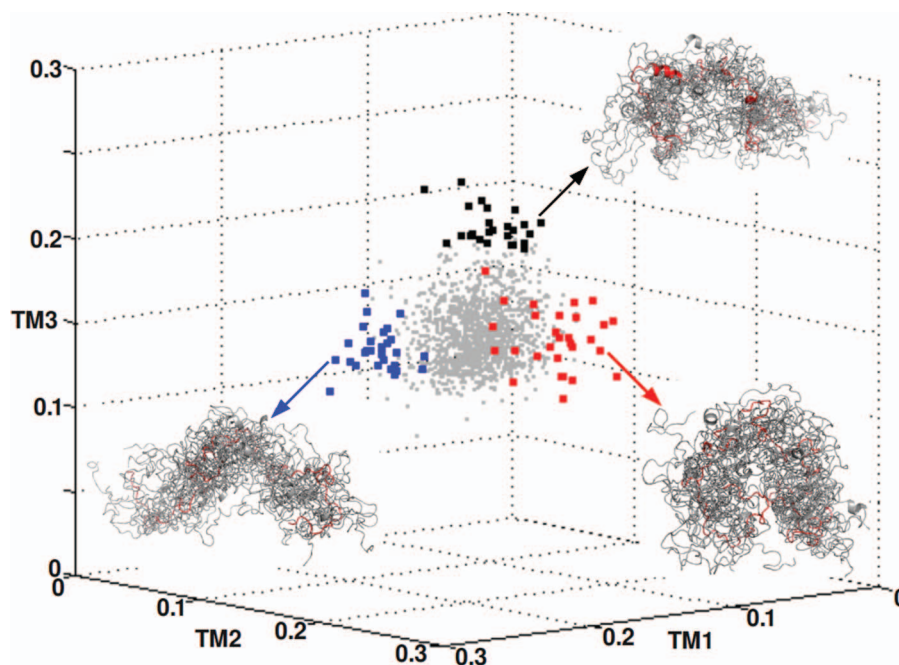


FIG. 11. Nearest neighbor clustering using TM-score of 1299 structures of  $\alpha$ -synuclein, projected onto the TM-scores to the centroid structures of the largest three clusters (blue, red, and black, respectively). Representative conformations in each cluster are shown. The lack of distinct clustering indicates diverse sampling of the unfolded ensemble.

We obtain the scaling exponent  $\nu$  for the ensembles at several times, shown in the inset of Figure 10(b). The exponent converges exponentially to an asymptotic value at  $t \rightarrow \infty$ . Extrapolating from times  $\leq 1$  ns is sufficient to obtain this asymptotic value to high accuracy. From this procedure, a scaling exponent that is slightly larger than a self-avoiding random walk (SAW) is obtained for proT $\alpha$ , likely because of the high charge density of this protein. We generally see scaling exponents for the unfolded ensemble with values between those in  $\Theta$ -solvents and those in good solvents (Table I). The IUPs in our study tended to have larger scaling exponents than foldable proteins (0.62 vs. 0.57 on average) – some of the IUPs such as proT $\alpha$  are highly charged; none fall into the class of collapsed globule IUPs.<sup>19,68</sup> Our observations on scaling exponents are generally consistent with recent experimental observations by Schuler and colleagues<sup>63</sup> from single-molecule spectroscopy experiments.

An estimate for the persistence length in the unfolded state may be obtained from the prefactor  $r_1$  as  $\ell_p = (2\nu + 1)(2\nu + 2)r_1^2/(2b)$ , where  $b = 3.8$  Å is the  $C_\alpha$ - $C_\alpha$  distance;<sup>69</sup> a SAW distribution has been assumed in this estimate. From this estimate, the persistence lengths we measured varied from 2 to 4 Å, averaging around 3 Å; these values were somewhat smaller than those obtained experimentally from force spectroscopy studies: typically around 3.5-4.0 Å (see, e.g., Refs. 70 and 71).

### C. Clustering analysis

We undertook a clustering analysis for  $\alpha$ -synuclein using the program `maxcluster` with nearest neighbor clustering.<sup>72</sup> The lack of strong clustering would be evidence of the success of our method to generate a diverse unfolded

ensemble. Figure 11 shows the distribution of 1299 configurations of  $\alpha$ -synuclein, projected along the TM-scores to the centroid structures of the three largest clusters. The elements of the clusters are indeed not well differentiated from the other structures, and the elements of the top three clusters are themselves fairly unrelated to the cluster centroid, with low TM-scores (cluster centroids have TM-scores of unity along their respective axis and are not shown – the centroids barely pull structures from the bulk ensemble).

We found that, for both the dominant clusters and the unfolded ensemble, the ends of the 140 aa protein tend to be closer on average than 140 aa stretches of other proteins that we had investigated (end to end distances  $r_{ee}(\alpha\text{syn}) \approx 89$  Å, whereas  $r_{ee}(1A6N_{140}) \approx 96$  Å), consistent with previous experimental NMR data that indicated aggregation-inhibiting, long-range tertiary interactions between the N- and C-termini.<sup>73</sup> It is a valid question as to whether longer simulation times would result in enhancement of long-range tertiary interactions. Our configurations for  $\alpha$ -syn were equilibrated for 5 ns; taking the equilibrated radius of gyration  $R_G$  of  $\alpha$ -syn for our ensemble and treating the polymer as a self-avoiding chain in a good solvent, the longest Rouse-Zimm-like relaxation times are<sup>74</sup>  $\tau_r \approx \frac{1}{3} \frac{\eta_s}{k_B T} R_G^3 \approx 10$  ns. This number is significantly longer than the relaxation time we observed for the radius of gyration for  $\alpha$ -syn:  $\approx 0.8$  ns, most likely because the slowest Rouse-Zimm modes have already come to equilibrium by construction of the unfolded ensemble. On the other hand, single molecule Förster resonance energy transfer (FRET) measurements in the denatured state of cold shock protein give longer global reconfiguration times,  $\approx 50$  ns.<sup>75</sup> As well, specific, tertiary contacts are observed to have significantly slower formation rates,<sup>76,77</sup> e.g., for the naphthalene-xanthone labelled 56 aa peptide poly(GS)<sub>28</sub>, the

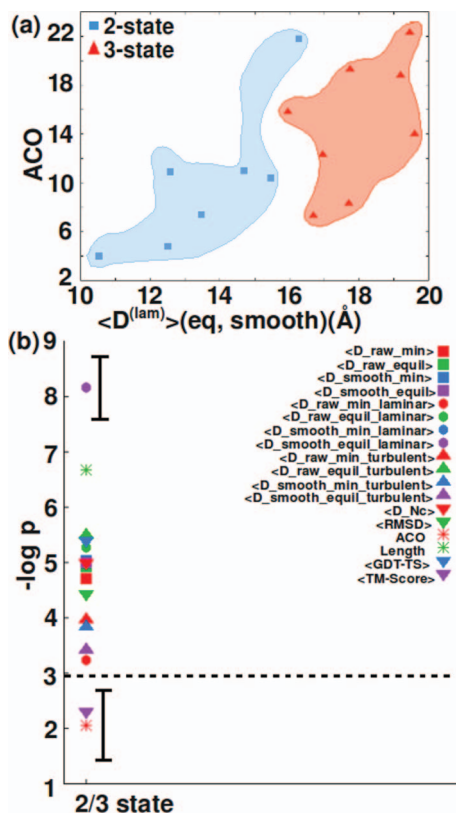


FIG. 12. (a) Scatter plot of the absolute contact order (ACO) and average laminar distance (equilibrium ensemble, with smoothed trajectories), for the 15 natively folded proteins in Table I. 2-state proteins (blue squares) and 3-state proteins (red triangles) are well-clustered by  $D^{(lam)}$ , but not by ACO, as can be seen by inspection, i.e., by projecting data onto each order parameter. Closed curves circumscribing each class of protein are a guide to the eye. (b) Statistical significance (p-values) that the various metrics for 2-state and 3-state folders arise from different distributions, as determined by t-test.<sup>30</sup>  $-\log(p)$  is plotted, so that a higher number indicates better ability to distinguish between the two classes. The dashed black horizontal line indicates a threshold of 5% for statistical significance. Only ACO and maxcluster-determined TM-score fail to distinguish 2-state from 3-state folders. Error bars for ACO and  $D^{(lam)}$  are obtained by removing 1 data point at random from the dataset, recomputing  $-\log(p)$ , and then calculating the standard deviation for the resulting collection of values. Notation used in this panel is further described in Figure 14.

time constant for the formation of specific end-to-end contacts is  $\approx 170$  ns.<sup>76</sup> Exploration of such contact dynamics in  $\alpha$ -syn and other proteins is an interesting topic of future research.

We had found previously that a coarse-grained measure of the mean distance  $\mathcal{D}$  in Eq. (5), as well as  $\langle RMSD \rangle$  and chain length, significantly discriminates 2- and 3-state folders.<sup>30</sup> Here, we investigated which order parameters in Sec. II C cluster 2-state folders sufficiently separate from 3-state folders, such that they may be discriminated from each other. Figure 12(a) shows a scatter plot of the clustering along ACO and  $D^{(lam)}$  (with smoothed trajectories and for the equilibrium ensemble). Figure 12(b) plots the negative logarithm of the statistical significance, based on a t-test,<sup>30</sup> that each order parameter distinguishes 2-state from 3-state proteins (see Table S3 for listed p-values). More significant distinguishers have larger values on this plot.

## D. Correlations between geometrical folding pathways and folding kinetics

Having generated statistically diverse, quasi-equilibrium, unfolded ensembles, we turned to the question as to whether transformations between such an ensemble and the native state could address folding kinetics. We investigated optimal folding trajectories both by our previous non-crossing method,<sup>30</sup> and by the GP method of Thorpe and colleagues;<sup>31</sup> these are described in Sec. II C. We focused primarily on the GP method because it applied to all-atom systems, and because the transformation could be visualized at all intermediate stages.

Information on the folding mechanism is gained from determining which quantity correlates with rate for a given structural or kinetic class of protein. For example, the fact that absolute contact order (ACO)<sup>27,78</sup> or extensions such as long range order<sup>79–82</sup> or total contact distance<sup>83</sup> correlate well with rate for 2-state proteins indicates a dominance of the process of loop closure, through the formation of native contacts, as the rate limiting step in folding. The fact that these quantities do not strongly correlate with rates for 3-state proteins, and that chain length does,<sup>84</sup> indicates that other, perhaps more subtle mechanisms embodying topological complexity may play a role in determining folding barrier heights.

Since the early studies of contact order, many subsequent studies have investigated correlates with protein folding kinetics across protein classes. Extensions of contact order have been developed from polymer theory,<sup>85</sup> which support early polymer physics models predicting rate determining barriers scaling as  $\sim N^{1/2}$ .<sup>86</sup> Mean field theory had predicted folding barriers increasingly linearly with protein chain length.<sup>87</sup> Rates were found to correlate with thermodynamic properties, including native stability<sup>88</sup> and heterogeneity of contact formation probability and  $\phi$ -values.<sup>89–91</sup> Combinations of contact order and length have been taken.<sup>92</sup> Rates were observed to correlate with contact clustering<sup>93</sup> and to anticorrelate with the number of tightly packed contacts defined through Delaunay edges.<sup>94</sup> Folding rates for 2-state proteins were observed to correlate with helix, turn, and hairpin secondary structure propensity,<sup>95</sup> and inversely with chain length, in a model with 4 weighting coefficients. Consistent with this, both 2-state and 3-state folding rates anticorrelate with residual length after helical segments are renormalized to 3 residues in length,<sup>96</sup> implying significant secondary structure formation in the transition state. Folding times for two-state proteins were also observed to correlate with chain length within structural classes, but not across them;<sup>80</sup> the difficulty for ACO as a predictor across structural classes is also born out in statistical physics models.<sup>97</sup> Other studies have found that rates of 3-state proteins significantly anticorrelated with  $\alpha$  and  $\beta$  secondary structure length, and that 2-state protein rates anticorrelate with  $\beta$ +loop secondary structure length.<sup>98</sup> Rates of both 2- and 3-state proteins can now be predicted to very high-accuracy from sequence alone,<sup>94,99–101</sup> but such approaches generally involve a large number of parameters – up to 20 for each amino acid type, or up to 49 amino acid properties – so that there is not much prospect of uncovering an underlying mechanism.



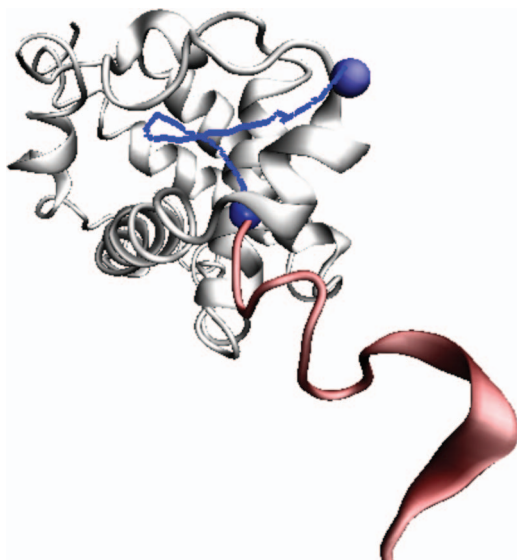


FIG. 13. Optimal folding trajectory of  $C_\alpha(50)$  in apo-myoglobin (1A6N). The trajectory is curved, due to steric constraints with the remainder of the protein.  $C_\alpha(50)$  is shown as blue spheres in the initial and final states. The region of protein N-terminal to  $C_\alpha(50)$  in the initial unfolded state is shown in red. This transforms to the short helix N-terminal to  $C_\alpha(50)$  in the final position.

Here, we asked the simple, physically based question as to whether the distance covered by optimal trajectories, generated as described in Sec. II C, would anticorrelate with folding rate across proteins (both 2- and multi-state). We also asked if simpler alignment measures such as RMSD, TM-SCORE, and GDT-TS would anticorrelate (for RMSD) or correlate (for TM-SCORE and GDT-TS) with rates, when averaged over the unfolded ensembles generated as described in Sec. II A.

In finding quantities related to distance  $\mathcal{D}$  travelled, we found that the GP trajectories often meandered significantly from straightline motion, indicating that even for minimal transformations, steric hindrance induces curvilinear motion for folding trajectories (Fig. 13). We record the laminar and turbulent components of the distance between all unfolded conformations and the native structure, as described in Sec. II C 1.

Figure 14 shows a matrix of the Pearson correlation coefficients (upper triangle) and statistical significance (lower triangle), between the order parameters investigated here, including the 12 variants of distance described in Figure 7,  $\langle \text{RMSD} \rangle$ ,  $\langle \text{TM-SCORE} \rangle$  obtained from the program `max-cluster`,  $\langle \text{GDT-TS} \rangle$ , ACO, and protein length. Also included in the table are the log experimental folding and unfolding rates  $k_f$  and  $k_u$ , along with the log of the midpoint relaxation rate  $k_{mp}$ . A corresponding table with Kendall values is given in the supplementary material.<sup>47</sup>

Smoothing the trajectories did not significantly change an order parameter's correlation with rates. On the other hand, in nearly all cases, equilibration *decreased* an order parameter's correlation with rates. We suspect that this may be an indication of either the fast-mixing experimental protocol often used to dilute an initially high-denaturant state, or of problems with the force fields for the SPC solvent model, in which case more

refined solvent models such as TIP4P may yield improved results. On the other hand, it is fortunate that the minimized, non-equilibrated ensemble performs so well – the computationally expensive procedure of equilibration may then not be necessary.

Alignment metrics to the native structure, such as the ensemble-averaged  $\langle \text{RMSD} \rangle$ ,  $\langle \text{TM-SCORE} \rangle$ , and  $\langle \text{GDT-TS} \rangle$ , all showed significant correlation with folding rates across 2- and 3-state proteins. It is perhaps surprising that a quantity as simple as RMSD has not been tested in folding rate prediction, but this may be because information on pairs of structures rather than a single (native) structure is needed to calculate it.

The mean distances obtained by our previous polymer noncrossing method for coarse-grained  $C_\alpha$ -model polymers<sup>30</sup> correlate very strongly with those obtained from the GP method ( $r = 0.99$ ,  $p = 6 \times 10^{-12}$ ). Moreover, these mean distances  $\langle \mathcal{D} \rangle$  correlate remarkably strongly with the mean RMSD of the unfolded ensemble  $\langle \text{RMSD} \rangle$  ( $r = 0.94$  on average), so that this latter quantity may be used as a crude proxy for either distance calculation. This is a fortunate result in the sense that the calculation of  $\langle \text{RMSD} \rangle$  is less computationally intensive than distance calculations by either method. It also means that the qualitative result is captured by simple polymer models. We will see however that significant quantitative effects are observed that depend on the all-atom steric volume of the protein, and its role in obstructing or guiding folding.

All variants of the distance travelled – minimized/equilibrated and smoothed/raw – showed significant correlation with folding rates. However, the laminar component of the distance travelled does not correlate strongly with rates, and is insignificant for rates in water. We have checked the convergence of the correlation between various distance metrics and folding rates; convergence is achieved fairly quickly, after about 100 unfolded configurations (Fig. S2 in the supplementary material<sup>47</sup>).

Figure 15 gives a synopsis of the correlation between various quantities and folding rates, in terms of minus log base 10 of their statistical significance (e.g., a significance of  $10^{-4}$  would give a value of 4 on the plot). The most striking feature is the degree of correlation shown by the turbulent component of the distance travelled with folding rates in water ( $r = -0.95$ ,  $p = 1 \times 10^{-7}$ ). A scatter plot of folding rate vs.  $\mathcal{D}^{(\text{turb})}$  is shown in Figure 16(a). The ensemble-averaged RMSD also correlates significantly with folding rate – comparable in general to ACO (Fig. 16(b)).

The turbulent motion involves nonlinear docking and registering motions between at least partially formed secondary structures. It appears as “late-stage” reconfiguration of structured elements, and may be thought of as measuring the difficulty in fitting secondary structured units together. The observation that these motions appear to govern the barrier that determines rates implies a transition state ensemble with significant native structure present.

Of the quantities we investigated, unfolding rates at 0 M denaturant anticorrelate strongest with ACO ( $r = -0.84$ ,  $p = 1 \times 10^{-4}$ ); the entropy of loop closure governs the unfolding barrier as it does the folding barrier, and implies a mechanism for kinetic stability of native structures:



Pearson Corr. Coeff. p-value	<D_raw_min>	<D_raw_equil>	<D_smooth_min>	<D_smooth_equil>	<D_raw_min_laminar>	<D_raw_equil_laminar>	<D_smooth_min_laminar>	<D_smooth_equil_laminar>	<D_raw_min_turbulent>	<D_raw_equil_turbulent>	<D_smooth_min_turbulent>	<D_smooth_equil_turbulent>	<RMSD_min>	<RMSD_equil>	ACO	Length	<GDT_TS>	<TM_score>	ln_kf	ln_ku	ln_kmp
<D_raw_min>	1.	0.99	1.	0.99	0.75	0.71	0.87	0.9	0.86	0.91	0.97	0.97	0.97	0.96	0.86	0.9	-0.92	-0.84	-0.9	-0.81	-0.84
<D_raw_equil>	11.	1.	0.99	1.	0.76	0.7	0.87	0.91	0.84	0.93	0.95	0.98	0.98	0.97	0.83	0.93	-0.93	-0.89	-0.86	-0.8	-0.83
<D_smooth_min>	17.	11.	1.	0.99	0.76	0.71	0.88	0.92	0.85	0.91	0.96	0.97	0.97	0.96	0.85	0.92	-0.92	-0.84	-0.89	-0.79	-0.82
<D_smooth_equil>	11.	17.	11.	1.	0.78	0.72	0.89	0.92	0.92	0.92	0.94	0.98	0.98	0.98	0.82	0.94	-0.93	-0.89	-0.85	-0.78	-0.8
<D_raw_min_laminar>	2.8	3.	3.	3.2	1.	0.85	0.95	0.76	0.51	0.55	0.57	0.74	0.77	0.81	0.65	0.78	-0.83	-0.74	-0.47	-0.54	-0.55
<D_raw_equil_laminar>	2.5	2.4	2.5	2.6	4.3	1.	0.78	0.56	0.54	0.4	0.6	0.75	0.71	0.77	0.62	0.74	-0.7	-0.7	-0.48	-0.49	-0.48
<D_smooth_min_laminar>	4.5	4.6	4.8	5.	7.	3.2	1.	0.91	0.69	0.73	0.72	0.82	0.89	0.9	0.74	0.87	-0.9	-0.78	-0.65	-0.65	-0.68
<D_smooth_equil_laminar>	5.4	5.7	5.9	6.	3.	1.5	5.6	1.	0.82	0.89	0.83	0.83	0.91	0.88	0.69	0.91	-0.91	-0.81	-0.76	-0.65	-0.72
<D_raw_min_turbulent>	7.5	6.5	7.1	6.	1.3	1.4	2.4	3.8	1.	0.93	0.99	0.92	0.91	0.87	0.82	0.82	-0.81	-0.75	-0.95	-0.8	-0.84
<D_raw_equil_turbulent>	5.7	6.5	5.7	6.1	1.5	0.85	2.7	5.1	6.4	1.	0.92	0.89	0.9	0.87	0.75	0.82	-0.84	-0.79	-0.86	-0.78	-0.82
<D_smooth_min_turbulent>	8.7	7.4	8.1	6.9	1.6	1.8	2.6	3.9	13.	6.	1.	0.94	0.92	0.9	0.83	0.85	-0.83	-0.79	-0.93	-0.79	-0.82
<D_smooth_equil_turbulent>	8.7	10.	8.2	9.8	2.8	2.9	3.8	3.8	5.9	4.9	6.9	1.	0.96	0.97	0.84	0.9	-0.89	-0.88	-0.84	-0.8	-0.8
<RMSD_min>	8.6	9.6	9.1	10.	3.1	2.5	5.	5.7	5.5	5.4	6.	7.8	1.	0.99	0.85	0.94	-0.92	-0.88	-0.86	-0.8	-0.79
<RMSD_equil>	7.7	9.1	7.9	9.9	3.7	3.1	5.3	4.8	4.6	4.5	5.2	9.	11.	1.	0.82	0.92	-0.92	-0.89	-0.81	-0.79	-0.77
ACO	4.5	3.8	4.3	3.8	2.1	1.8	2.8	2.4	3.7	2.9	3.9	4.	4.2	3.8	1.	0.67	-0.77	-0.6	-0.87	-0.84	-0.8
Length	5.4	6.2	5.8	6.6	3.2	2.8	4.6	5.5	3.7	3.7	4.2	5.2	6.6	6.	2.2	1.	-0.87	-0.86	-0.74	-0.64	-0.66
<GDT_TS>	5.8	6.2	5.9	6.4	3.9	2.4	5.4	5.5	3.6	4.1	4.	5.	6.	5.9	3.1	4.5	1.	0.84	0.71	0.7	0.77
<TM_score>	4.	5.	4.1	5.1	2.7	2.5	3.2	3.5	2.9	3.3	3.4	4.8	4.9	4.9	1.8	7.6	4.	1.	0.65	0.63	0.62
ln_kf	5.3	4.4	5.	4.1	1.1	1.1	2.	3.	7.	4.5	6.5	4.1	4.4	3.6	4.5	2.7	2.5	2.1	1.	0.86	0.85
ln_ku	3.5	3.5	3.4	3.3	1.4	1.2	2.	2.1	3.4	3.2	3.4	3.5	3.5	3.3	4.	2.	2.5	1.9	4.4	1.	0.92
ln_kmp	4.1	3.8	3.8	3.5	1.5	1.2	2.3	2.6	4.	3.7	3.7	3.4	3.3	3.1	3.5	2.1	3.1	1.9	4.2	6.	1.

FIG. 14. Correlation matrix for all geometrical parameters, as well as experimental folding rates. The upper triangular elements are Pearson correlation coefficients. The lower triangular elements are the corresponding statistical significance values, which are represented as  $-\log_{10}$  so that, e.g., 4.5 corresponds to  $p = 10^{-4.5} = 3.2 \times 10^{-5}$ . Red represents strong positive correlation; blue represents strong negative correlation. “\_raw” indicates numbers taken from the raw trajectory, while “\_smooth” indicates numbers taken from the smoothed trajectory. Trajectories are further divided into “\_laminar” and “\_turbulent” parts. Initial ensembles are either equilibrated “\_equil.” or pre-equilibration (energy minimized only or “\_min”). Other parameters shown include ACO, protein length, GDT-TS, TM-score, natural log of the folding and unfolding rates in 0 M denaturant, and natural log of relaxation rate at the transition midpoint.

long-range contacts promote more cooperative unfolding barriers. The distance-based metrics are still significant: the correlation coefficient of  $\mathcal{D}$  with folding rate is  $-0.81$ .

The dominance of turbulent distance as a rate-predictor goes away at the transition midpoint (Fig. 15(c)). The total distance travelled along with the turbulent distance are still the strongest predictors however. Laminar components of the motion, while themselves weak predictors, become more significant in these conditions than they are in water.

This effect is consistent with expansion of the unfolded state as solvent conditions are varied by adding denaturant.<sup>75</sup> In the absence of denaturant, the unfolded state may generally be a significantly collapsed molten-globule like state.<sup>26,102</sup> Early kinetic folding intermediates of apoMb are nearly as collapsed as the native structure.<sup>103</sup> A pre-collapsed, compact unfolded state, driven primarily by hydrophobic interactions, is seen as well in two-state folding reactions for SH3,<sup>104</sup> thermostable variants of cold shock protein,<sup>105</sup> destabilizing

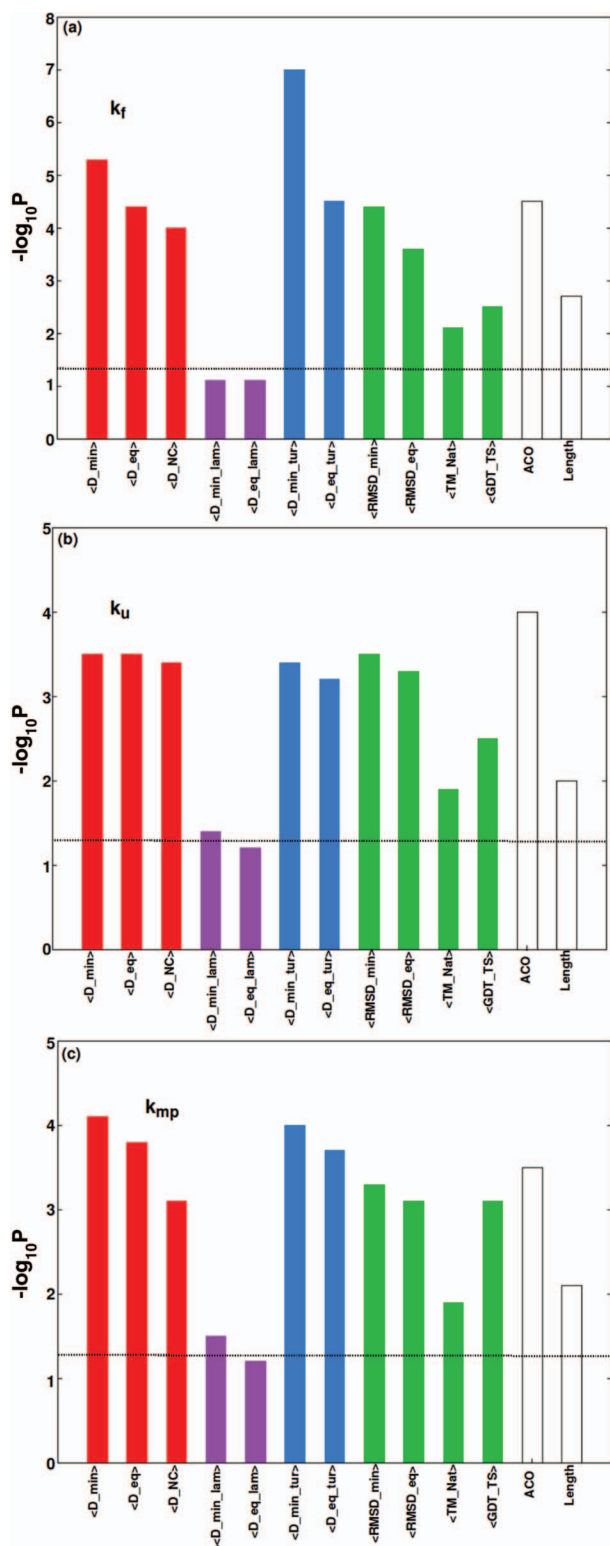


FIG. 15. (Panel (a)) Correlation of various distance metrics with experimental refolding rate in water, for the dataset of proteins listed in Table I. Raw (rather than smoothed) data are taken here. Minus the log base 10 of the statistical significance is plotted, and the horizontal dashed line gives the threshold of statistical significance ( $p = 0.05$ ). The best predictor of folding rates in water, the turbulent distance, has a significance of  $10^{-7}$ . Each integer below this value in the plot corresponds to a decrease in significance by an order of magnitude. (Panel (b)) Same as panel (a) but for experimental unfolding rate in water. Here, ACO emerges as the strongest correlator of unfolding rate. (Panel (c)) Same as panel (a) but for relaxation rate at the transition midpoint. Here, several variants of the distance travelled correlate best with relaxation rate, e.g., both  $D$  and  $D^{(\text{turb})}$  have a correlation coefficient  $r = -0.84$ .

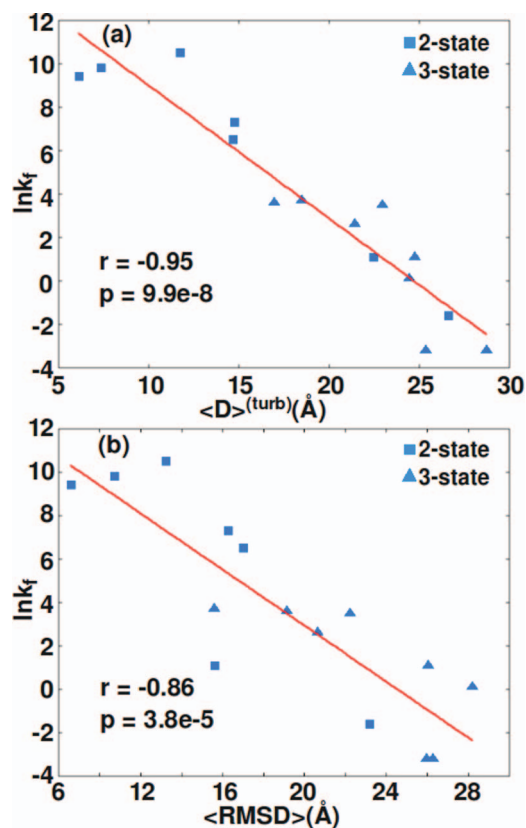


FIG. 16. (Panel (a)) Scatter plot of experimental folding rate at 0 M denaturant with the unfolded ensemble-averaged turbulent distance travelled during folding, corresponding to late-stage protein reconfiguration of structured elements. (Panel (b)) Scatter plot of the folding rate at 0 M denaturant with the ensemble-averaged RMSD between unfolded structures and the native. For both plots, the pre-equilibrated, energy-minimized, ensemble is taken, and raw rather than smoothed data are taken. Data for 2-state proteins are shown as squares, data for 3-state proteins are shown as triangles.

mutants of WT NTL9,<sup>106</sup> and Trp-cage miniprotein.<sup>107</sup> Fast initial collapse on the time-scale of tens of nanoseconds is seen directly in FRET measurements of BBL<sup>108</sup> and simulation of small proteins such as Villin.<sup>42,109</sup> Collapse emerges naturally in statistical field theories of heteropolymer collapse<sup>86,110–113</sup> as well as coarse grained computer simulations of folding.<sup>114–116</sup> Further along in the folding process, the transition state of CI2 has been interpreted as a globally collapsed, condensed nucleus with significant native structure, only slightly expanded, and lacking specific native packing interactions.<sup>117</sup> These experimental observations support the changing nature of the unfolded state with denaturant, which is reflected in which component of the minimal distance serves as the best predictor of folding rates.

#### IV. CONCLUSIONS

Here, we have connected the problem of generating unfolded ensembles with refolding kinetics, by applying transformations between unfolded conformations and the native structure. We developed a method to generate a diverse, quasi-equilibrium unfolded ensemble by employing coarse-grained sampling, foliation of the coarse-grained structure with side chains and backbone, and short equilibration of each

configuration. Ensembles for proteins with disulfide bonds can be generated as well, by employing non-local crankshaft-like moves.

Chemical shifts showed general agreement with experimental values, while residual dipolar couplings did not correlate with experimental values. We proposed some possible reasons for this discrepancy, including the fact that steric liquid crystal media that would induce rotational anisotropy in experiments would themselves modify the distribution of disordered conformations – a phenomenon that is not a factor for folded proteins.

Distance metrics as applied here between unfolded conformations and the native structure can also be applied to the ensemble of unfolded conformations, to obtain a general measure of the connectivity of the unfolded state. Distance metrics correlated strongly with common metrics of structural similarity, e.g., RMSD. The average RMSD between the folded structure and the unfolded ensemble correlated as strongly with folding rate as absolute contact order did.

The turbulent distance characterizes motion towards the folded structure that involves steric avoidance, jostling, rearrangement, and ultimately docking of highly structured units. The amount of motion this involved correlated most significantly with folding time ( $r = 0.95$ ,  $p = 1 \times 10^{-7}$ ). The dominance of “turbulent” minimal trajectories in predicting folding rates in water, but not at the transition midpoint, is a manifestation of a largely collapsed unfolded state in water or similar conditions, and the importance of reconfigurational motions from this state in finding transition states conducive to rapid progress toward the native structure.

## ACKNOWLEDGMENTS

S.S.P acknowledges funding support from PrioNet Canada, NSERC, and computational support from the West-Grid high-performance computing consortium.

- <sup>1</sup>J. Deisenhofer, O. Epp, K. Miki, R. Huber, and H. Michel, *Nature (London)* **318**, 618 (1985).
- <sup>2</sup>D. A. Doyle, J. M. Cabral, R. A. Pfuetzner, A. Kuo, J. M. Gulbis, S. L. Cohen, B. T. Chait, and R. MacKinnon, *Science* **280**, 69 (1998).
- <sup>3</sup>P. Agre, L. S. King, M. Yasui, W. B. Guggino, O. P. Ottersen, Y. Fujiyoshi, A. Engel, and S. Nielsen, *J. Physiol.* **542**, 3 (2002).
- <sup>4</sup>N. Ban, P. Nissen, J. Hansen, P. B. Moore, and T. A. Steitz, *Science* **289**, 905 (2000).
- <sup>5</sup>A. P. Carter, W. M. Clemons, D. E. Brodersen, R. J. Morgan-Warren, B. T. Wimberly, and V. Ramakrishnan, *Nature (London)* **407**, 340 (2000).
- <sup>6</sup>F. Schluenzen, A. Tocilj, R. Zarivach, J. Harms, M. Gluehmann, D. Janell, A. Bashan, H. Bartels, I. Agmon, F. Franceschi, and A. Yonath, *Cell* **102**, 615 (2000).
- <sup>7</sup>A. L. Gnatt, P. Cramer, J. Fu, D. A. Bushnell, and R. D. Kornberg, *Science* **292**, 1876 (2001).
- <sup>8</sup>J. A. Pitcher, N. J. Freedman, and R. J. Lefkowitz, *Annu. Rev. Biochem.* **67**, 653 (1998).
- <sup>9</sup>V. Cherezov, D. M. Rosenbaum, M. A. Hanson, S. G. Rasmussen, F. S. Thian, T. S. Kobilka, H.-J. Choi, P. Kuhn, W. I. Weis, B. K. Kobilka, and R. C. Stevens, *Science* **318**, 1258 (2007).
- <sup>10</sup>P. E. Wright and H. J. Dyson, *J. Mol. Biol.* **293**, 321 (1999).
- <sup>11</sup>V. N. Uversky, J. R. Gillespie, and A. L. Fink, *Proteins: Struct., Funct., Bioinf.* **41**, 415 (2000).
- <sup>12</sup>A. K. Dunker, J. D. Lawson, C. J. Brown, R. M. Williams, P. Romero, J. S. Oh, C. J. Oldfield, A. M. Campen, C. M. Ratliff, K. W. Hipps, J. Ausio, M. Nissen, R. Reeves, C. Kang, C. R. Kissinger, R. W. Bailey, M. D. Griswold, W. Chiu, E. C. Garner, and Z. Obradovic, *J. Mol. Graphics Modell.* **19**, 26 (2001).
- <sup>13</sup>P. Tompa, *Trends Biochem. Sci.* **27**, 527 (2002).
- <sup>14</sup>H. J. Dyson and P. E. Wright, *Nat. Rev. Mol. Cell Biol.* **6**, 197 (2005).
- <sup>15</sup>A. K. Dunker, C. Oldfield, J. Meng, P. Romero, J. Yang, J. Chen, V. Vacic, Z. Obradovic, and V. Uversky, *BMC Genomics* **9**, S1 (2008).
- <sup>16</sup>P. E. Wright and H. J. Dyson, *Curr. Opin. Struct. Biol.* **19**, 31 (2009).
- <sup>17</sup>I. Baskakov and D. W. Bolen, *J. Biol. Chem.* **273**, 4831 (1998).
- <sup>18</sup>P. Bernadó, C. W. Bertocini, C. Griesinger, M. Zweckstetter, and M. Blackledge, *J. Am. Chem. Soc.* **127**, 17968 (2005).
- <sup>19</sup>X. Wang, A. Vitalis, M. A. Wyczalkowski, and R. V. Pappu, *Proteins: Struct., Funct., Bioinf.* **63**, 297 (2006).
- <sup>20</sup>N. G. Sgourakis, Y. Yan, S. A. McCallum, C. Wang, and A. E. Garcia, *J. Mol. Biol.* **368**, 1448 (2007).
- <sup>21</sup>A. Vitalis, X. Wang, and R. V. Pappu, *Biophys. J.* **93**, 1923 (2007).
- <sup>22</sup>O. O. Olubiyi and B. Strodel, *J. Phys. Chem. B* **116**, 3280 (2012).
- <sup>23</sup>S. E. Jónsson, S. Mohanty, and A. Irbäck, *Proteins: Struct., Funct., Bioinf.* **80**, 2169 (2012).
- <sup>24</sup>A. K. Jha, A. Colubri, K. F. Freed, and T. R. Sosnick, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 13099 (2005).
- <sup>25</sup>P. Bernadó, L. Blanchard, P. Timmins, D. Marion, R. W. H. Ruigrok, and M. Blackledge, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 17002 (2005).
- <sup>26</sup>J. B. Udgaonkar, *Arch. Biochem. Biophys.* **531**, 24 (2013).
- <sup>27</sup>K. W. Plaxco, K. T. Simons, and D. Baker, *J. Mol. Biol.* **277**, 985 (1998).
- <sup>28</sup>Y. Zhang, *BMC Bioinf.* **9**, 40 (2008).
- <sup>29</sup>N. S. Bogatyreva, A. A. Osypov, and D. N. Ivankov, *Nucleic Acids Res.* **37**, D342 (2009).
- <sup>30</sup>A. R. Mohazab and S. S. Plotkin, *PLoS ONE* **8**, e53642 (2013).
- <sup>31</sup>D. W. Farrell, K. Speranskiy, and M. F. Thorpe, *Proteins: Struct., Funct., Bioinf.* **78**, 2908 (2010).
- <sup>32</sup>M. Lal, *Mol. Phys.* **17**, 57 (1969).
- <sup>33</sup>N. Madras and A. D. Sokal, *J. Stat. Phys.* **50**, 109 (1988).
- <sup>34</sup>S. Hadizadeh, A. Linhananta, and S. S. Plotkin, *Macromolecules* **44**, 6182 (2011).
- <sup>35</sup>G. Favrin, A. Irbäck, and F. Sjunnesson, *J. Chem. Phys.* **114**, 8154 (2001).
- <sup>36</sup>N. Gö and H. A. Scheraga, *Macromolecules* **3**, 178 (1970).
- <sup>37</sup>L. R. Dodd, T. D. Boone, and D. N. Theodorou, *Mol. Phys.* **78**, 961 (1993).
- <sup>38</sup>D. Hoffmann and E.-W. Knapp, *Eur. Biophys. J.* **24**, 387 (1996), see <http://dx.doi.org/10.1007/BF00576711>.
- <sup>39</sup>M. R. Betancourt, *J. Chem. Phys.* **123**, 174905 (2005), see <http://link.aip.org/link/?JCP/123/174905/1>.
- <sup>40</sup>P. Rotkiewicz and J. Skolnick, *J. Comput. Chem.* **29**, 1460 (2008).
- <sup>41</sup>J. Maupetit, R. Gautier, and P. Tufféry, *Nucleic Acids Res.* **34**, W147 (2006).
- <sup>42</sup>B. Zagrovic, C. D. Snow, S. Khaliq, M. R. Shirts, and V. S. Pande, *J. Mol. Biol.* **323**, 153 (2002).
- <sup>43</sup>M. N. Rosenbluth and A. W. Rosenbluth, *J. Chem. Phys.* **23**, 356 (1955).
- <sup>44</sup>P. Grassberger, *Phys. Rev. E* **56**, 3682 (1997).
- <sup>45</sup>K. J. Kohlhoff, P. Robustelli, A. Cavalli, X. Salvatella, and M. Vendruscolo, *J. Am. Chem. Soc.* **131**, 13894 (2009).
- <sup>46</sup>M. Zweckstetter, *Nat. Protoc.* **3**, 679 (2008).
- <sup>47</sup>See supplementary material at <http://dx.doi.org/10.1063/1.4817215> for supporting figures, tables, and description.
- <sup>48</sup>RMSD values for backbone and all heavy atoms were also examined and gave similar results.
- <sup>49</sup>W. Humphrey, A. Dalke, and K. Schulten, *J. Mol. Graphics* **14**, 33 (1996).
- <sup>50</sup>Y. Zhang and J. Skolnick, *Nucleic Acids Res.* **33**, 2302 (2005).
- <sup>51</sup>A. Zemla, *Nucleic Acids Res.* **31**, 3370 (2003).
- <sup>52</sup>J. P. Rodrigues, M. Levitt, and G. Chopra, *Nucleic Acids Res.* **40**, W323 (2012).
- <sup>53</sup>S. S. Plotkin, *Proc. Natl. Acad. Sci. U.S.A.* **104**, 14899 (2007).
- <sup>54</sup>A. R. Mohazab and S. S. Plotkin, *J. Phys. Condens. Matter* **20**, 244133 (2008).
- <sup>55</sup>A. R. Mohazab and S. S. Plotkin, *Biophys. J.* **95**, 5496 (2008).
- <sup>56</sup>B. Thiruv, G. Quon, S. Saldanha, and B. Steipe, *BMC Struct. Biol.* **5**, 12 (2005).
- <sup>57</sup>L. Hou, H. Shao, Y. Zhang, H. Li, N. K. Menon, E. B. Neuhaus, J. M. Brewer, I.-J. L. Byeon, D. G. Ray, M. P. Vitek, T. Iwashita, R. A. Makula, A. B. Przybyla, and M. G. Zagorski, *J. Am. Chem. Soc.* **126**, 1992 (2004).
- <sup>58</sup>A. De Simone, A. Cavalli, S.-T. D. Hsu, W. Vranken, and M. Vendruscolo, *J. Am. Chem. Soc.* **131**, 16332 (2009).
- <sup>59</sup>G. Lipari and A. Szabo, *J. Am. Chem. Soc.* **104**, 4546 (1982).
- <sup>60</sup>N. G. Sgourakis, M. Merced-Serrano, C. Boutsidis, P. Drineas, Z. Du, C. Wang, and A. E. Garcia, *J. Mol. Biol.* **405**, 570 (2011).

- <sup>61</sup>O. Tcherkasskaya and V. N. Uversky, *Proteins: Struct., Funct., Bioinf.* **44**, 244 (2001).
- <sup>62</sup>S. Doniach, *Chem. Rev.* **101**, 1763 (2001).
- <sup>63</sup>H. Hofmann, A. Soranno, A. Borgia, K. Gast, D. Nettels, and B. Schuler, *Proc. Natl. Acad. Sci. U.S.A.* **109**, 16155 (2012).
- <sup>64</sup>C. Tanford, K. Kawahara, and S. Lapanje, *J. Biol. Chem.* **241**, 1921 (1966).
- <sup>65</sup>D. K. Wilkins, S. B. Grimshaw, V. Receveur, C. M. Dobson, J. A. Jones, and L. J. Smith, *Biochemistry* **38**, 16424 (1999).
- <sup>66</sup>J. E. Kohn, I. S. Millett, J. Jacob, B. Zagrovic, T. M. Dillon, N. Cingel, R. S. Dothager, S. Seifert, P. Thiagarajan, T. R. Sosnick, M. Z. Hasan, V. S. Pande, I. Ruczinski, S. Doniach, and K. W. Plaxco, *Proc. Natl. Acad. Sci. U.S.A.* **101**, 12491 (2004).
- <sup>67</sup>P. J. Flory, *Principles of Polymer Chemistry* (Cornell University Press, Ithaca, 1953).
- <sup>68</sup>V. N. Uversky, *Protein Sci.* **11**, 739 (2002).
- <sup>69</sup>B. Hammouda, *Adv. Polymer Sci.* **106**, 87 (1993).
- <sup>70</sup>R. B. Best, B. Li, A. Steward, V. Daggett, and J. Clarke, *Biophys. J.* **81**, 2344 (2001).
- <sup>71</sup>D. Brockwell, E. Paci, R. Zinober, G. Beddard, P. Olmsted, D. Smith, R. Perham, and S. Radford, *Nat. Struct. Biol.* **10**, 731 (2003).
- <sup>72</sup>D. Shortle, K. T. Simons, and D. Baker, *Proc. Natl. Acad. Sci. U.S.A.* **95**, 11158 (1998); see also <http://www.sbg.bio.ic.ac.uk/maxcluster>.
- <sup>73</sup>C. W. Bertoncini, Y.-S. Jung, C. O. Fernandez, W. Hoyer, C. Griesinger, T. M. Jovin, and M. Zweckstetter, *Proc. Natl. Acad. Sci. U.S.A.* **102**, 1430 (2005).
- <sup>74</sup>M. Doi and S. F. Edwards, *The Theory of Polymer Dynamics* (Oxford University Press, Oxford, 1986), Chap. 4.
- <sup>75</sup>D. Nettels, I. V. Gopich, A. Hoffmann, and B. Schuler, *Proc. Natl. Acad. Sci. U.S.A.* **104**, 2655 (2007).
- <sup>76</sup>F. Krieger, B. Fierz, O. Bieri, M. Drewello, and T. Kiefhaber, *J. Mol. Biol.* **332**, 265 (2003).
- <sup>77</sup>J. Kubelka, J. Hofrichter, and W. A. Eaton, *Curr. Opin. Struct. Biol.* **14**, 76 (2004).
- <sup>78</sup>K. W. Plaxco, K. T. Simons, I. Ruczinski, and D. Baker, *Biochemistry* **39**, 11177 (2000).
- <sup>79</sup>M. M. Gromiha and S. Selvaraj, *J. Mol. Biol.* **310**, 27 (2001).
- <sup>80</sup>A. Y. Istomin, D. J. Jacobs, and D. R. Livesay, *Protein Sci.* **16**, 2564 (2007).
- <sup>81</sup>B. Harihar and S. Selvaraj, *Biopolymers* **91**, 928 (2009).
- <sup>82</sup>T. Zou and S. B. Ozkan, *Phys. Biol.* **8**, 066011 (2011).
- <sup>83</sup>H. Zhou and Y. Zhou, *Biophys. J.* **82**, 458 (2002).
- <sup>84</sup>D. N. Ivankov, S. O. Garbuzynskiy, E. Alm, K. W. Plaxco, D. Baker, and A. V. Finkelstein, *Protein Sci.* **12**, 2057 (2003).
- <sup>85</sup>M. Rustad and K. Ghosh, *J. Chem. Phys.* **137**, 205104 (2012).
- <sup>86</sup>D. Thirumalai, *J. Phys. I* **5**, 1457 (1995).
- <sup>87</sup>S. S. Plotkin, J. Wang, and P. G. Wolynes, *J. Chem. Phys.* **106**, 2932 (1997).
- <sup>88</sup>A. R. Dinner and M. Karplus, *Nat. Struct. Biol.* **8**, 21 (2001).
- <sup>89</sup>S. S. Plotkin and J. N. Onuchic, *J. Chem. Phys.* **116**, 5263 (2002).
- <sup>90</sup>B. Oztop, M. R. Ejtehadi, and S. S. Plotkin, *Phys. Rev. Lett.* **93**, 208105 (2004).
- <sup>91</sup>U. Bastolla, P. Bruscolini, and J. L. Velasco, *Proteins: Struct., Funct., Bioinf.* **80**, 2287 (2012).
- <sup>92</sup>N. Koga and S. Takada, *J. Mol. Biol.* **313**, 171 (2001).
- <sup>93</sup>C. Micheletti, *Proteins: Struct., Funct., Bioinf.* **51**, 74 (2003).
- <sup>94</sup>Z. Ouyang and J. Liang, *Protein Sci.* **17**, 1256 (2008).
- <sup>95</sup>H. Gong, D. G. Isom, R. Srinivasan, and G. D. Rose, *J. Mol. Biol.* **327**, 1149 (2003).
- <sup>96</sup>D. Ivankov and A. Finkelstein, *Proc. Natl. Acad. Sci. U.S.A.* **101**, 8942 (2004).
- <sup>97</sup>P. Bruscolini, A. Pelizzola, and M. Zamparo, *Phys. Rev. Lett.* **99**, 038103 (2007).
- <sup>98</sup>J.-T. Huang, J.-P. Cheng, and H. Chen, *Proteins: Struct., Funct., Bioinf.* **67**, 12 (2007).
- <sup>99</sup>M. M. Gromiha, *J. Chem. Inf. Model.* **45**, 494 (2005).
- <sup>100</sup>H.-B. Shen, J.-N. Song, and K.-C. Chou, *J. Biomed. Sci. Eng.* **2**, 136 (2009).
- <sup>101</sup>L. Chang, J. Wang, and W. Wang, *Phys. Rev. E* **82**, 051930 (2010).
- <sup>102</sup>A. D. Miranker and C. M. Dobson, *Curr. Opin. Struct. Biol.* **6**, 31 (1996).
- <sup>103</sup>D. Eliezer, P. A. Jennings, P. E. Wright, S. Doniach, K. O. Hodgson, and H. Tsuruta, *Science* **270**, 487 (1995).
- <sup>104</sup>A. Dasgupta and J. B. Udgaonkar, *J. Mol. Biol.* **403**, 430 (2010).
- <sup>105</sup>C. Magg and F. X. Schmid, *J. Mol. Biol.* **335**, 1309 (2004).
- <sup>106</sup>B. Anil, Y. Li, J.-H. Cho, and D. P. Raleigh, *Biochemistry* **45**, 10110 (2006).
- <sup>107</sup>K. H. Mok, L. T. Kuhn, M. Goez, I. J. Day, J. C. Lin, N. H. Andersen, and P. Hore, *Nature (London)* **447**, 106 (2007).
- <sup>108</sup>M. Sadqi, L. J. Lapidus, and V. Muñoz, *Proc. Natl. Acad. Sci. U.S.A.* **100**, 12117 (2003).
- <sup>109</sup>Y. Duan and P. A. Kollman, *Science* **282**, 740 (1998).
- <sup>110</sup>J. D. Bryngelson and P. G. Wolynes, *Biopolymers* **30**, 177 (1990).
- <sup>111</sup>S. S. Plotkin and J. N. Onuchic, *Q. Rev. Biophys.* **35**, 111 (2002).
- <sup>112</sup>S. S. Plotkin and J. N. Onuchic, *Q. Rev. Biophys.* **35**, 205 (2002).
- <sup>113</sup>G. Ziv, D. Thirumalai, and G. Haran, *Phys. Chem. Chem. Phys.* **11**, 83 (2009).
- <sup>114</sup>N. D. Succi and J. N. Onuchic, *J. Chem. Phys.* **103**, 4732 (1995).
- <sup>115</sup>A. M. Gutin, V. I. Abkevich, and E. I. Shakhnovich, *Biochemistry* **34**, 3066 (1995).
- <sup>116</sup>Z. Guo and D. Thirumalai, *Biopolymers* **36**, 83 (1995).
- <sup>117</sup>A. R. Fersht, *Curr. Opin. Struct. Biol.* **5**, 79 (1995).
- <sup>118</sup>M. Gromiha, A. Thangakani, and S. Selvaraj, *Nucleic Acids Res.* **34**, W70 (2006).



# Unfolded protein ensembles, folding trajectories, and refolding rate prediction: Supplementary Material

A. Das, B. K. Sin, A. R. Mohazab, and S. S. Plotkin

*Department of Physics and Astronomy, University of British Columbia, 6224 Agricultural Road, Vancouver, B.C. V6T 1Z1, Canada*

steve@phas.ubc.ca

(Dated: 15 July 2013)

## I. UNFOLDED ENSEMBLE FOR SUPEROXIDE DISMUTASE

A sample of 1000 5ns-equilibrated conformations from the unfolded ensemble of the disulfide-bonded protein superoxide dismutase (WT SOD1) is provided in two formats:

- A concatenated pdb file `sod1_all.pdb` (169 MB)
- An xtc trajectory file `sod1_all.xtc` (10.8 MB) along with a single pdb file `0001.pdb` (170K).

To view these files in VMD, issue the following commands from the terminal:

```
vmd -pdb sod1_all.pdb
```

or:

```
vmd -pdb 0001.pdb -xtc sod1_all.xtc
```

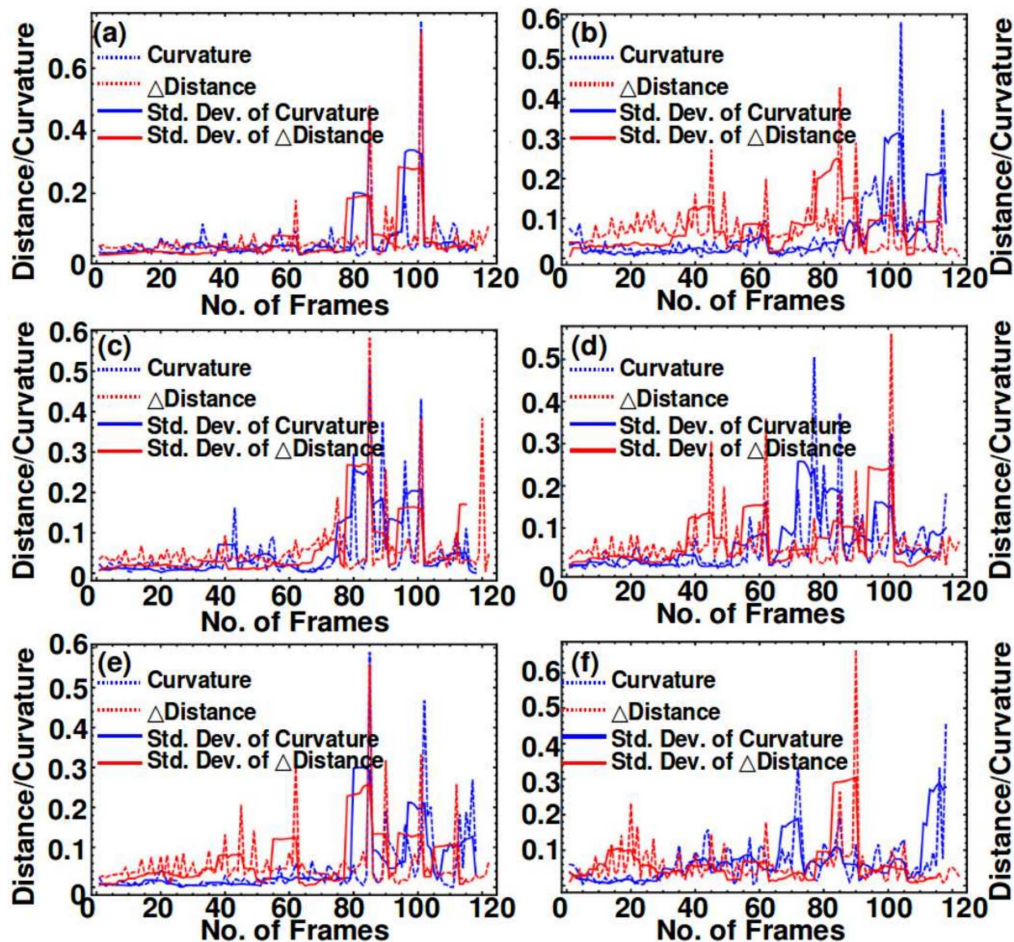


FIG. S1. **Transition from laminar to turbulent trajectories.** Each panel shows properties of the folding trajectory for selected  $C_\alpha$  atoms. (Red dashed) Distance  $\Delta\mathcal{D}$  traversed per frame, (Blue dashed) Curvature of the trajectory, (Red solid) running standard deviation of  $\Delta\mathcal{D}$  over 10 frames, (Blue solid) running standard deviation of the curvature over 10 frames. Panels (a), (c), and (e) show good correlation between the two methods, while panels (b), (d), and (f) show some discrepancies between the two methods.

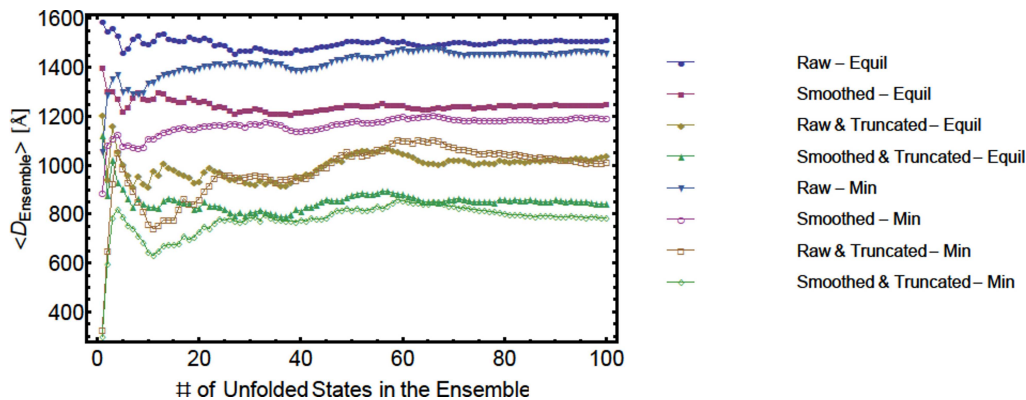


FIG. S2. **Convergence of mean distance travelled with ensemble size.** Convergence of various metrics for the total distance, as number of configurations in the unfolded ensemble is increased, for engrailed homeodomain 1ENH. Total distances converge after approximately 100 configurations.

Kendall Corr. Coeff. p-value	<D_raw_min>	<D_raw_equil>	<D_smooth_min>	<D_smooth_equil>	<D_raw_min_laminar>	<D_raw_equil_laminar>	<D_smooth_min_laminar>	<D_smooth_equil_laminar>	<D_raw_min_turbulent>	<D_raw_equil_turbulent>	<D_smooth_min_turbulent>	<D_smooth_equil_turbulent>	<RMSD_min>	<RMSD_equil>	ACO	Length	<GDT_TS>	<TM_score>	ln_kf	ln_ku	ln_kmp
<D_raw_min>	1.	0.89	1.	0.9	0.6	0.79	0.71	0.77	0.85	0.71	0.87	0.81	0.87	0.83	0.73	0.79	-0.75	-0.73	-0.74	-0.56	-0.59
<D_raw_equil>	5.4	1.	0.89	0.98	0.6	0.75	0.71	0.81	0.77	0.79	0.79	0.85	0.9	0.9	0.62	0.79	-0.78	-0.77	-0.66	-0.6	-0.56
<D_smooth_min>	6.7	5.4	1.	0.9	0.6	0.79	0.71	0.77	0.85	0.71	0.87	0.81	0.87	0.83	0.73	0.79	-0.75	-0.73	-0.74	-0.56	-0.59
<D_smooth_equil>	5.6	6.5	5.6	1.	0.62	0.77	0.73	0.79	0.75	0.77	0.77	0.87	0.92	0.92	0.64	0.81	-0.8	-0.79	-0.64	-0.58	-0.54
<D_raw_min_laminar>	2.7	2.7	2.7	2.9	1.	0.77	0.89	0.68	0.45	0.47	0.47	0.56	0.58	0.66	0.52	0.62	-0.67	-0.56	-0.39	-0.39	-0.54
<D_raw_equil_laminar>	4.4	4.	4.4	4.2	4.2	1.	0.81	0.79	0.64	0.54	0.66	0.68	0.77	0.77	0.6	0.81	-0.82	-0.75	-0.53	-0.43	-0.5
<D_smooth_min_laminar>	3.7	3.7	3.7	3.9	5.4	4.6	1.	0.79	0.56	0.58	0.58	0.68	0.7	0.77	0.64	0.73	-0.78	-0.68	-0.51	-0.47	-0.57
<D_smooth_equil_laminar>	4.2	4.6	4.2	4.4	3.4	4.4	4.4	1.	0.73	0.71	0.71	0.66	0.75	0.75	0.54	0.79	-0.78	-0.73	-0.59	-0.45	-0.48
<D_raw_min_turbulent>	5.	4.2	5.	4.	1.7	3.	2.5	3.9	1.	0.75	0.94	0.73	0.71	0.68	0.7	0.64	-0.63	-0.58	-0.82	-0.52	-0.56
<D_raw_equil_turbulent>	3.7	4.4	3.7	4.2	1.8	2.3	2.6	3.7	4.	1.	0.73	0.87	0.73	0.77	0.56	0.66	-0.65	-0.64	-0.63	-0.62	-0.54
<D_smooth_min_turbulent>	5.2	4.4	5.2	4.2	1.8	3.2	2.6	3.7	6.	3.9	1.	0.75	0.77	0.7	0.75	0.66	-0.69	-0.6	-0.84	-0.54	-0.57
<D_smooth_equil_turbulent>	4.6	5.	4.6	5.2	2.5	3.4	3.4	3.2	3.9	5.2	4.	1.	0.83	0.9	0.7	0.71	-0.73	-0.7	-0.64	-0.64	-0.56
<RMSD_min>	5.2	5.6	5.2	5.8	2.6	4.2	3.5	4.	3.7	3.9	4.2	4.8	1.	0.92	0.71	0.85	-0.88	-0.83	-0.68	-0.58	-0.54
<RMSD_equil>	4.8	5.6	4.8	5.8	3.2	4.2	4.2	4.	3.4	4.2	3.5	5.6	5.8	1.	0.68	0.81	-0.82	-0.79	-0.63	-0.58	-0.54
ACO	3.9	2.9	3.9	3.	2.2	2.7	3.	2.3	3.5	2.5	4.	3.5	3.7	3.4	1.	0.6	-0.65	-0.54	-0.7	-0.6	-0.67
Length	4.4	4.4	4.4	4.6	2.9	4.6	3.9	4.4	3.	3.2	3.2	3.7	5.	4.6	2.7	1.	-0.84	-0.94	-0.57	-0.5	-0.5
<GDT_TS>	4.	4.3	4.	4.5	3.3	4.7	4.3	4.3	3.	3.1	3.4	3.8	5.3	4.7	3.1	4.9	1.	0.82	0.6	0.54	0.53
<TM_score>	3.9	4.2	3.9	4.4	2.5	4.	3.4	3.9	2.6	3.	2.7	3.5	4.8	4.4	2.3	6.	4.7	1.	0.55	0.52	0.44
ln_kf	3.9	3.2	3.9	3.1	1.4	2.2	2.1	2.6	4.6	2.9	4.8	3.1	3.4	2.9	3.5	2.5	2.7	2.3	1.	0.68	0.6
ln_ku	2.5	2.7	2.5	2.6	1.4	1.6	1.8	1.7	2.2	2.9	2.3	3.	2.6	2.6	2.7	2.1	2.3	2.2	3.4	1.	0.77
ln_kmp	2.7	2.4	2.7	2.3	2.3	2.	2.5	1.9	2.4	2.3	2.5	2.4	2.3	2.3	3.3	2.	2.2	1.6	2.7	4.1	1.

FIG. S3. Correlation matrix for all geometrical parameters as well as folding rates. The upper triangular elements are Kendall correlation coefficients. The lower triangular elements are the corresponding statistical significance values, which are represented as  $-\log_{10}$  so that e.g. 4.5 corresponds to  $p = 10^{-4.5} = 3.2e-5$ . Red represents strong positive correlation; blue represents strong negative correlation. “raw” indicates numbers taken from the raw trajectory, while “smooth” indicates numbers taken from the smoothed trajectory. Trajectories are further divided into “laminar” and “turbulent” parts. Initial ensembles are either equilibrated “\_equil”, or pre-equilibration (energy minimized only or “\_min”). Other parameters shown include ACO, protein length, GDT-TS, TM-score, natural log of the folding and unfolding rates in 0M denaturant, and natural log of relaxation rate at the transition midpoint.

TABLE S1. Residual dipolar couplings for simulated unfolded ensembles

Residue # / Protein	ILZY	IUVI	IUYT	2PDD	IENH	ISHG	H-integrate	2CRO	1CSP	1PWF	ILMQ	ITTT	IAPS	IUNI	ProT	ICBI	IYQX	IAGN	35R2D	2A5E	1RA0
1																					
2	-1.20097	0.0605753	-0.094398				-1.437	-0.549046	-1.028568	1.21902	-1.879	0.524305		-4.25113	0.451055	-0.539663	-0.831028	2.59249	4.2314	7.85124	
3	0.38086	-1.14529	2.08989				1.134	-2.20707	1.07725	-0.556135	1.49642	3.91555	1.87431		2.81987	0.43064	1.80827	-4.88336		-4.60669	
4	1.4542	1.36295	-0.362355				-1.327	-0.211473	2.1211	-2.03885	-0.262484	2.76337	0.548189	-1.59272	1.00896	1.32216	1.43814	0.344984	0.222699	5.64418	-1.56911
5	1.65776	1.21078	1.52435	0.962537			1.207	-0.707974	0.957837	-1.7421	0.50487	-1.43547	0.87287	-0.81279	0.900887	-0.851428	-0.027378	-2.73822	-3.27278	1.51858	
6	4.57027	0.787021	-2.65346	-1.38748			1.984	-2.45333	-1.38082	4.83925	3.04432	-3.52421	3.85805	-0.92213	-0.496224	-2.76671	0.416556	1.95685	5.99053	1.17812	0.810242
7	2.30181	0.244609	-2.41515	2.18495	1.1237		0.915	0.0201088	-0.0094972	-0.122168	2.33804	2.3004	0.425677	-1.06234	-2.36406	-0.339959	-5.03671	1.57337	1.04303	-1.25106	
8	1.68903	0.756388	-0.275751	0.160691	1.76023		1.022	-0.52145	0.13327	-1.06649	0.26737	0.399551	-2.89353	1.58533	0.74053	0.524459	-2.8303	-3.75343	-3.887	-1.5787	
9	1.05568	0.560431	1.06729	0.250395	-1.78592	-3.87	0.0979159	-0.693773	-1.09392	-0.597314	-1.31998	-4.6491	3.4869	0.315381	-1.64945	0.310447	1.99627	-4.05729	6.09122	0.26396	
10	2.83272	0.985069	-1.40685	0.482511	2.07912	0.654	1.48821	0.62837	2.07295	-2.2316	0.400664	-1.51489	-0.78828	-2.73552	0.67533	-0.939435	0.0236107	0.999013	2.64339	-0.237741	
11	1.27626	0.226871	-0.488455	-0.493947	3.29125	0.177	0.937655	0.222886	0.68828	0.502791	-0.046752	-1.1434	-1.32657	2.26396	4.10214	0.342204	1.19834	-5.19705		4.30341	
12	2.82674	0.865412	1.39393	0.0177406	0.757117	-0.054	-1.73676	2.49655	1.37329	1.13832	-3.02228	1.85041	-1.68975	1.47545	-2.22790	0.534982	-5.71094	-3.63672	-2.99888	-1.91581	0.15357
13	4.63063	1.45933	1.8002	-2.08479	2.17388	0.389	0.817959	2.47528	2.23115	0.833239	0.0671421	-2.27786	-1.28421	4.00106	2.08453	-4.87993	0.84539	-8.81844	-0.537307	2.7709	
14	3.11858	1.69752	0.088692	-0.215748	3.11674	-0.742	-0.503529	2.77608	1.54966	-8.80354	-2.78362	1.04291	-0.104966	1.83578	-1.7706	0.962196	-6.65785	-1.82984	-6.63725	0.604975	
15	0.45573	0.79666	0.728662	2.3332	0.0692001	0.525	0.18941	0.074307	0.092728	1.85011	-1.68975	1.47545	-2.22790	0.534982	-5.71094	0.540494	-3.63672	-2.99888	-1.91581	0.15357	
16	3.58277	0.653864	1.76859	-0.162881	0.19981	0.060	-0.961419	1.45863	0.482047	-1.13329	2.30351	-1.16927	0.621825	1.73969	0.625576	1.05823	-0.754129	-3.55883	-2.61825	0.071763	
17	-	-0.204336	2.02308	0.213467	1.26376	0.385	1.6594	1.54568	-0.828776	3.02684	-3.4253	-5.1509	-2.56145	-2.5699	-0.186022	0.6462	-1.61196	1.38589	2.88854	1.24854	
18	0.74609	2.50734	0.373533	-1.26872	0.193	0.263489	1.54677	-1.07296	0.483591	-2.31093	1.58753	-0.59986	-2.00012	-2.59089	-1.04745	-2.29565	-1.04745	-2.29565	-4.21472	7.15182	
19	-	1.8997	1.45496	-1.76141	-0.548988	0.1	-0.4453082	-1.41394	-1.07056	-0.42187	-0.449866	0.285918	-2.97771	-1.18538	-0.13477	-1.18538	-0.13477	-1.18538	-0.500185	1.12895	
20	1.08074	0.263612	1.5738	-0.336626	0.126	1.146	-0.17794	-0.336626	3.61254	3.16646	-0.446654	0.864373	-8.5883	0.166646	-0.864373	-0.864373	-0.864373	-0.864373	-1.48469	-0.046612	
21	-	0.756773	0.143864	0.712803	-0.24159	1.225	3.14383	-0.29204	-2.17781	1.26035	3.22408	0.6816	1.13996	0.104798	-4.52411	1.33923	-4.4918	-0.796724	2.60002	0.75792	
22	-	0.32331	-0.435331	-1.59693	0.141536	0.61	-0.21749	-0.123488	-1.02904	0.57019	-0.32059	1.65984	-0.29015	0.4422	-0.467028	1.05249	1.89974	-1.89974	-2.26039	1.48082	
23	-	-0.944872	1.4295	-4.87786	-1.88826	0.934	-0.489017	2.74614	-1.2455	2.4165	0.0438072	2.24584	-1.78075	-0.497773	1.42112	-0.930477	1.42112	-0.930477	5.22734	3.70781	7.1247
24	-	0.018028	1.30643	0.281959	-2.12731	0.24	2.5734	2.07358	-0.399564	-0.0450783	-2.43738	1.47042	-0.298663	5.74571	-1.10195	-2.99081	-1.69806	3.3204	3.21	4.74807	
25	-	1.1906	2.57571	0.601429	1.42815	0.239	1.31716	1.17011	1.54719	0.584166	-1.52446	1.32597	0.315284	-0.259749	1.315284	-0.259749	1.315284	-0.259749	1.68	1.56236	
26	-	1.3134	-0.43921	0.043138	0.552822	0.476	-0.246647	2.97663	0.812759	0.629649	-0.36023	3.08057	-2.02422	3.22104	-0.932446	3.22104	-0.932446	3.22104	2.98922	6.29544	
27	-	0.952429	0.089692	1.41538	1.40301	0.459	-1.89882	2.66639	-0.218245	1.77479	-0.345191	-3.26205	1.50854	0.157703	-1.63017	-1.84997	3.06934	3.29991	1.00568	10.124	
28	-	0.3221	1.33563	0.0448883	-0.190229	0.354	0.31225	1.6615	0.66012	0.62218	1.26447	0.26746	-1.01312	-1.53895	0.26746	-1.01312	-1.53895	0.26746	2.419	0.04742	
29	-	0.0043216	3.4444	0.0873672	-2.04725	1.871	1.42733	-1.08026	2.42502	4.95101	0.736707	0.789465	-2.27682	-0.985964	-0.258894	0.618966	1.06859	-0.544225	-1.19016	2.36422	
30	-	-0.439239	3.1855	1.39021	1.25692	0.914	2.13359	2.23385	2.4472	2.50622	2.30806	3.13687	0.736478	-1.40581	-1.71975	0.40239	4.02466	-5.64276	-6.44768	-4.46908	
31	-	1.1492	1.52996	-0.174858	-3.463	0.186	0.753765	0.813651	4.63261	-1.45209	2.37027	4.4374	2.87001	1.41303	-3.85991	-0.606531	2.11534	-0.26287	-3.27289	2.03661	
32	-	0.9256	1.631	1.39472	-0.051971	0.697	0.62854	0.889623	0.253335	-1.33446	0.253335	7.28476	0.26449	-0.52876	0.26449	-0.52876	0.26449	-0.52876	0.26449	0.26449	
33	-	1.08833	0.589274	-0.463991	0.807099	-0.7	-0.890134	-0.463991	1.21941	-1.040817	1.95514	3.92273	1.69925	1.30155	0.265073	-3.72994	3.16325	-3.72994	3.16325	-0.70634	
34	-	1.05304	1.082641	0.647987	0.278458	0.167	1.66048	0.898082	3.36702	-2.73994	0.892425	1.90294	2.51012	-1.91232	-4.06068	0.302759	-6.63134	-1.42384	-0.967988	4.05797	
35	-	0.32967	1.91919	-0.174375	0.543708	1.797	5.29	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	0.0638099	
36	-	0.601581	-0.58417	0.545351	0.446876	0.5	-0.370419	-3.90537	2.06882	-1.64414	3.52551	-1.52937	2.08076	2.24152	1.16586	-1.80257	-0.307414	5.39485	6.31608	8.31608	
37	-	1.02857	1.34491	0.826614	2.65287	0.807	-0.493437	1.86489	3.09452	-3.49384	2.11697	-3.73577	6.28908	1.21567	3.18988	0.937943	-0.20674	2.97427	-2.97427	-6.19108	
38	-	1.81498	1.84689	2.6741	1.51649	1.181	0.67838	3.68278	0.680882	-7.22329	1.7392	2.3307	3.94962	1.80189	4.88271	1.80189	4.88271	1.80189	4.88271	4.88271	
39	-	0.821209	0.720452	-0.793705	-1.36112	0.602	-1.24296	-0.138507	2.3445	-0.294876	1.22354	0.407381	-1.73088	0.451055	-0.75627	-1.73088	0.451055	-0.75627	-1.73088	4.88271	
40	-	2.42732	1.51246	-0.698889	0.865549	0.561	2.57706	-2.36943	0.327144	2.95361	-2.92009	-2.73379	-1.40805	-0.0417923	0.0899798	1.40805	1.23057	-3.94396	-1.38113		
41	-	1.14918	0.331248	-0.646699	0.744237	0.213	0.597907	0.753691	3.73207	1.55613	0.854433	0.854433	-2.31994	3.78441	-0.485561	0.542907	1.52342	-1.74112	1.95756	3.14852	
42	-	4.24258		0.85857		0.030	2.21422	-1.50812	3.85731	1.74574	3.61847	-1.63986	-3.25236	2.99082	1.42657	-1.58645	0.886637	3.245	2.34816	0.388271	
43	-	1.07446		0.439529	1.63807	0.215	2.18284	-1.09705	1.7476	1.98357	-0.239993	-0.207922	-0.554893	-0.93631	1.78845	-0.927733	-0.10489	0.341779	-3.88849	10.5277	
44	-	-0.814577		-0.018685	1.5284	0.029	-0.129339	1.36887	0.686952	2.46483	2.58156	2.86322	0.00916	0.00916	0.00916	0.00916	0.00916	0.00916	0.00916	0.00916	
45	-	1.2838		-2.56636	-0.979336	1.089	-0.801379	-0.447792	0.460191	2.51136	-0.573601	1.63511	3.785	0.7978	-0.95759	3.91488	2.83472	-1.48122	-2.71291	-2.71291	
46	-	1.18082		-2.33549	-1.196017	0.784	1.27552	-0.132149	1.10811	0.77912	-0.720912	0.727386	0.676961	-0.677192	1.08546	-1.24438	1.01798	2.88829	1.28325	1.92382	
47	-	1.45109		1.29518	-3.92828	1.36	2.66901	-0.645948	1.51713	0.574139	-0.624922	1.73588	1.73588	1.73588	1.73588	1.73588	1.73588	1.73588	1.73588	1.73588	
48	-	1.55301		-0.201167	1.27708	0.642	1.58204	-2.69677	2.7911	0.519141	-0.447081	-2.19228	2.250715	2.67211	1.29333	0.93784	9.58696	-1.17326	-3.96994	-0.97703	
49	-	3.24052		-0.619423	-0.076228	1.286	-2.02964	-0.40191	-0.711653	2.4972	0.265622	-0.855438	2.59718	-2.39817	2.11211	1.61208	0.27884	6.39661	-7.57491	1.5461	
50	-</																				





TABLE S3. **Discriminating 2-state from 3-state folders**

Parameter	p value <sup>a</sup>
$\langle D_{\text{raw\_min}} \rangle$	9.0e-3
$\langle D_{\text{raw\_equil}} \rangle$	7.2e-3
$\langle D_{\text{smooth\_min}} \rangle$	6.5e-3
$\langle D_{\text{smooth\_equil}} \rangle$	6.8e-3
$\langle D_{\text{raw\_min\_laminar}} \rangle$	4.0e-2
$\langle D_{\text{raw\_equil\_laminar}} \rangle$	5.2e-3
$\langle D_{\text{smooth\_min\_laminar}} \rangle$	4.5e-3
$\langle D_{\text{smooth\_equil\_laminar}} \rangle$	3.0e-4
$\langle D_{\text{raw\_min\_turbulent}} \rangle$	1.9e-2
$\langle D_{\text{raw\_equil\_turbulent}} \rangle$	4.1e-3
$\langle D_{\text{smooth\_min\_turbulent}} \rangle$	2.2e-2
$\langle D_{\text{smooth\_equil\_turbulent}} \rangle$	3.2e-2
$\langle D_{\text{NC}} \rangle$	6.7e-3
$\langle \text{RMSD} \rangle$	1.2e-2
ACO	1.3e-1
Length	1.3e-3
$\langle \text{GDT-TS} \rangle$	4.5e-3
$\langle \text{TM-Score} \rangle$	1.0e-1

<sup>a</sup> Statistical significance (p-values) for each metric in this table is determined by a t-test with null hypothesis that values of a given metric for the 2-state and 3-state folders come from independent random samples with normal distributions having equal means and equal but unknown variances, against the alternative that the means are not equal. The values in this table are plotted in figure 12 of the main text.