

Polymer scaling laws of unfolded and intrinsically disordered proteins quantified with single-molecule spectroscopy

Hagen Hofmann^{a,1}, Andrea Soranno^a, Alessandro Borgia^a, Klaus Gast^b, Daniel Nettels^a, and Benjamin Schuler^{a,1}

^aBiochemisches Institut, Universität Zürich, Winterthurerstrasse 190, 8057 Zürich, Switzerland; and ^bPhysikalische Biochemie, Universität Potsdam, 14476 Potsdam, Germany

Edited by Ken A. Dill, Stony Brook University, Stony Brook, NY, and approved August 15, 2012 (received for review May 8, 2012)

The dimensions of unfolded and intrinsically disordered proteins are highly dependent on their amino acid composition and solution conditions, especially salt and denaturant concentration. However, the quantitative implications of this behavior have remained unclear, largely because the effective theta-state, the central reference point for the underlying polymer collapse transition, has eluded experimental determination. Here, we used single-molecule fluorescence spectroscopy and two-focus correlation spectroscopy to determine the theta points for six different proteins. While the scaling exponents of all proteins converge to 0.62 ± 0.03 at high denaturant concentrations, as expected for a polymer in good solvent, the scaling regime in water strongly depends on sequence composition. The resulting average scaling exponent of 0.46 ± 0.05 for the four foldable protein sequences in our study suggests that the aqueous cellular milieu is close to effective theta conditions for unfolded proteins. In contrast, two intrinsically disordered proteins do not reach the Θ -point under any of our solvent conditions, which may reflect the optimization of their expanded state for the interactions with cellular partners. Sequence analyses based on our results imply that foldable sequences with more compact unfolded states are a more recent result of protein evolution.

protein folding | single-molecule FRET | coil-globule transition | polymer theory

It has become increasingly clear that the structure and dynamics of unfolded proteins are essential for understanding protein folding (1–3) and the functional properties of intrinsically disordered proteins (IDPs) (4–6). Theoretical concepts from polymer physics (7–9) have frequently been used to describe the properties of unfolded polypeptide chains (4, 10, 11) with the goal to establish the link between protein folding and collapse (12–15). However, the methodology to test many of these concepts experimentally has only become available rather recently (2, 16, 17). A considerable body of experimental and theoretical work suggests that the dimensions of unfolded proteins depend on parameters such as amino acid composition (4), temperature (18), and solvent quality (3, 10, 15, 19). The continuous collapse of polymers has been treated exhaustively by a number of theories (20–24) based on general principles that relate the dimensions and the length of a chain to its free energy. However, a prerequisite for the quantitative application of these theories and their comparison to experimental results is that the dimensions of the Θ -state are known, which serves as an essential reference state. At the Θ -point*, chain–chain and chain–solvent interactions balance such that the polymer is at a critical point, at which the thermodynamic phase boundaries disappear. As a result, the polypeptide chain obeys the same length scaling as an ideal chain without excluded volume and intrachain interactions. However, the Θ -conditions for protein chains are unknown. Besides its importance for obtaining the correct thermodynamic parameters of the chain, such as excluded volume and interaction energies, the Θ -state for proteins has been suggested to be of special biological relevance since folding is predicted to occur most efficiently when the

Θ -point coincides with the transition midpoint for folding (9, 25, 26), while several previous results have been taken to suggest that unfolded proteins and folding intermediates are below the Θ -point under physiological conditions (27–30).

One way of obtaining this missing information is by means of scaling laws (20, 22) that relate the radius of gyration of the unfolded protein (R_G) to its length (N) via $R_G \propto N^\nu$. By determining the scaling exponent ν at different solvent conditions, the Θ -conditions are identified as the conditions for which $\nu = 1/2$. Here we used single-molecule Förster resonance energy transfer (smFRET) to systematically determine the dimensions of seventeen chain segments with different lengths in six different unfolded proteins at a wide range of denaturant concentrations, resulting in a large data set (Fig. 1A and *SI Appendix, Table S1*). To investigate the sequence dependence of the Θ -conditions, we chose four foldable proteins [cold shock protein, CspTm (3); cyclophilinA, hCyp (31); spectrin domains R15 and R17 (32)] and two more highly charged IDPs (prothymosin α , ProT α , and the N-terminal domain of HIV Integrase, IN) (4) (Fig. 1A and *SI Appendix, Table S1*). Estimates for the scaling exponent ν , the Θ -conditions, and the free energy of solvation could be obtained for all six proteins.

Results

To probe the dimensions of the unfolded states of the six proteins, we attached AlexaFluor 488 as a donor and AlexaFluor 594 as an acceptor chromophore at different positions within the polypeptide chains (*SI Appendix, Table S1*). The labeled proteins were investigated with confocal smFRET while freely diffusing in solution. In the resulting transfer efficiency histograms for each protein and variant, up to three peaks are observed: The peak at very high transfer efficiency (E) results from folded molecules, and the peak at $E \approx 0$ results from molecules lacking an active acceptor dye (Fig. 1B and *SI Appendix, Figs. S1–S3*). We focus exclusively on the peak at intermediate transfer efficiencies, which results from unfolded molecules (Fig. 1B). The use of smFRET allows us to discriminate this population of unfolded molecules from folded molecules even in the virtual absence of denaturant (*SI Appendix, Figs. S1–S3*). With increasing concentration of the denaturant GdmCl, the transfer efficiency distributions of the unfolded subpopulations of all variants show a pronounced shift to lower E values, corresponding to an expansion of the polypeptide

Author contributions: H.H. and B.S. designed research; H.H. and K.G. performed research; H.H., A.S., A.B., K.G., and D.N. contributed new reagents/analytic tools; H.H., A.S., K.G., and D.N. analyzed data; and H.H. and B.S. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

*The critical point for heteropolymers is an effective Θ -point (24), but for convenience, we will use the term Θ -point also for heteropolymers.

¹To whom correspondence may be addressed. E-mail: schuler@bioc.uzh.ch or h.hofmann@bioc.uzh.ch.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1207719109/-DCSupplemental.

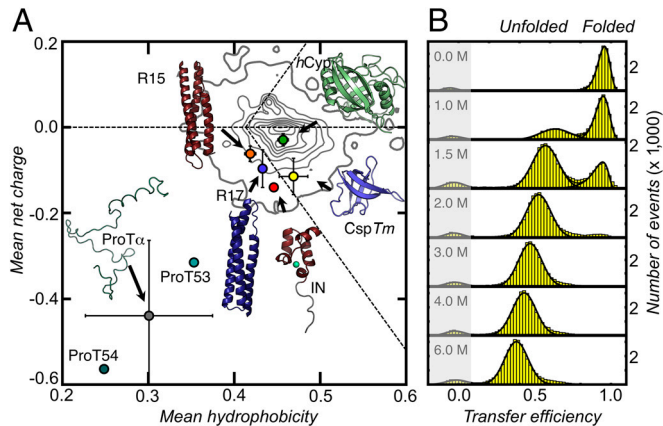


Fig. 1. Structures and amino acid compositions of the proteins used in this study (A) and single-molecule FRET efficiency histograms for CspTm (Csp66, *SI Appendix, Table S1*) at different concentrations of GdmCl (B). (A) Mean net charge, including the charges of the attached fluorophores, versus mean hydrophobicity per residue for hCyp, CspTm, R15, R17, IN, and ProTα (variants ProT53 and ProT54, *SI Appendix*) (circles). Error bars are standard deviations of mean net charge and mean hydrophobicity of the different variants of each protein. The density plot represents the distribution of 10,905 monomeric proteins with a sequence similarity $\leq 30\%$ taken from the Protein Data Bank. The horizontal dashed line indicates a mean net charge of zero. Diagonal dashed lines indicate the separation line between intrinsically disordered and folded proteins suggested by Uversky et al. (48).

chains (Fig. 1B and *SI Appendix, Figs. S1–S3*), as observed previously for a broad range of proteins and peptides (3, 10, 15, 19, 33).

Chain dimensions from FRET efficiencies. Quantitative information about the dimensions of the unfolded proteins can be obtained from the average values $\langle E \rangle$ of their transfer efficiency peaks. We used the coil-to-globule transition theory of Sanchez (21) to extract the chain dimensions from $\langle E \rangle$. The advantage of this theory is its ability to describe the dimensions of a chain under all solvent conditions by explicitly taking into account effects such as excluded volume, intrachain interactions, and multibody interactions (10, 11, 21). The theory provides an expression for the probability density function of the radius of gyration r_G in the form of a Boltzmann-weighted Flory–Fisk distribution (11, 34):

$$P(r_G, \varepsilon, R_{G\Theta}) = Z^{-1} r_G^6 \exp \left[-\frac{7r_G^2}{2R_{G\Theta}^2} + nq(\phi, \varepsilon) \right] \quad [1]$$

$$\text{with } q = \frac{1}{2} \varepsilon \phi - \frac{1 - \phi}{\phi} \ln(1 - \phi)$$

Here, $R_{G\Theta} \equiv \langle r_G^2 \rangle_{\Theta}^{1/2}$ is the root mean squared radius of gyration of the Θ -state; ε is the mean interaction energy between amino acids; ϕ is the volume fraction of the chain; n is the number of amino acids in the chain segment probed by FRET; Z is a normalization factor; and q is the excess free energy per monomer with respect to the ideal chain (11). An expression similar to Eq. 1 was also obtained in heteropolymer theories (12, 13), showing that Eq. 1 is not specific for homopolymers (*SI Appendix*). Note, however, that none of these descriptions take into account effects from sequence complexities; e.g., the patterning of residues.

In order to relate the distribution $P(r_G, \varepsilon, R_{G\Theta})$ to a segment end-to-end distance distribution $P(r, \varepsilon, R_{G\Theta})$, which is required to describe the transfer efficiencies of the polypeptide chains, we used the conditional probability density function $P(r|r_G)$ suggested by Ziv and Haran (11) (*SI Appendix, Eq. S1*). The observed mean transfer efficiency $\langle E \rangle$ is related to Eq. 1 by

$$\begin{aligned} \langle E \rangle &= \int_0^L E(r) P(r, \varepsilon, R_{G\Theta}) dr \\ &= \int_0^L E(r) \int_{R_C}^{L/2} P(r|r_G) P(r_G, \varepsilon, R_{G\Theta}) dr_G dr \\ &\text{with } E(r) = \frac{R_0^6}{R_0^6 + r^6}, \end{aligned} \quad [2]$$

where R_0 is the Förster radius (5.4 nm in our case) and L is the contour length of the protein segment probed. Importantly, the root mean squared radius of gyration of the chain segment, $R_G \equiv \langle r_G^2 \rangle^{1/2}$, is largely independent of the specific value of $R_{G\Theta}$ (*SI Appendix, Fig. S8*), which allows us to determine R_G for every protein segment from its mean transfer efficiency, $\langle E \rangle$. We then use the scaling of R_G with the number of peptide bonds in the unfolded protein segments, $R_G \propto N^\nu$, to determine $R_{G\Theta}$ from the conditions at which $\nu = 1/2$. With the correct value of $R_{G\Theta}$, we then determine ε exactly. $P(r|r_G)$ (*SI Appendix, Eq. S1*) assumes unfolded proteins to be spherical in shape, which is an approximation (35–37), but we investigated the accuracy of Eq. 2 by simulation and found the error in R_G to be $\leq 6\%$ (*SI Appendix, Fig. S5*).

The radius of gyration of polymers scales with the number of bonds (N) according to the power-law relation $R_G = \rho_0 N^\nu$. The specific value of ν depends on the dimensions of the chain, with a value of 3/5 for the expanded coil state (22), 1/2 for the Θ -state, and 1/3 for the most compact globule state (21, 35). In contrast, the value of the prefactor ρ_0 depends on the details of the monomer and the bond geometry. For a self-avoiding chain with scaling exponent ν , R_G is given by (38)

$$R_G = \rho_0 N^\nu = \sqrt{\frac{2l_p^* b}{(2\nu + 1)(2\nu + 2)}} N^\nu \quad [3]$$

(The derivation for a special case can also be found in ref. 34). Here, $b = 0.38$ nm (39) is the distance between two C_α -atoms, and l_p^* is the persistence length (*SI Appendix*). Values for ρ_0 from experiments (0.19 ± 0.03 nm and 0.2 ± 0.1 nm) (40, 41) and simulations (0.22 ± 0.02 nm, 0.24 nm, 0.198 ± 0.037 nm, and 0.199 nm) (42–45) obtained under good solvent conditions ($\nu = 3/5$) yield $l_p^* = 0.40 \pm 0.07$ nm, in agreement with persistence lengths from force spectroscopy experiments (39). Since the range of segment lengths accessible with smFRET is not broad enough to determine ρ_0 independently, we fixed l_p^* (but not ρ_0) to this value of 0.40 nm. For comparison, a free fit of the length scaling of R_G for 10,905 folded proteins selected from the Protein Data Bank results in $\nu = 0.34$ and a persistence length of $l_p^* = 0.53$ nm (Fig. 2) (35), but even using this value for our analysis as an upper bound does not change our conclusions (*SI Appendix*).

Identifying the Θ conditions from FRET and two-focus FCS. Previous measurements of the scaling exponent ν for unfolded proteins at high concentrations of denaturant resulted in values between 0.50 and 0.67 (40, 41, 46, 47). In the most extensive study, R_G for 28 proteins was determined by SAXS in the presence of high concentrations of GdmCl or urea (40). From this data set, $\nu = 0.598 \pm 0.028$ was obtained, indistinguishable from the theoretical prediction of 3/5 for an excluded volume chain (22), which indicates that unfolded proteins are in the coil-state and in good solvent at high concentrations of denaturant (Fig. 2). Under comparable solvent conditions (6 M GdmCl), we found the R_G values from smFRET to be in remarkable agreement with $R_G = 0.2 \text{ nm } N^{3/5}$, the scaling law obtained with SAXS (40) (Fig. 2). The scaling exponents we obtained at 6 M GdmCl range from 0.59 for hCyp to 0.63 for the hydrophilic IDP integrase. The high ν -value of prothymosin α ($\nu = 0.67$), a highly negatively

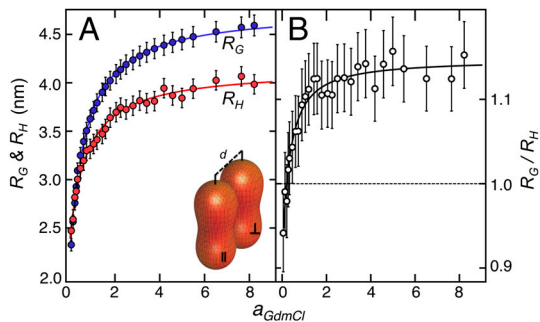


Fig. 4. Comparison between the radii of gyration and the hydrodynamic radii for *hCyp* as a function of GdmCl activity. (A) Radius of gyration, R_G , (blue circles) for Cyp163 (*SI Appendix, Table S1*) rescaled to the full length sequence ($N_{\text{bonds}} = 166 + 9$) according to the scaling laws shown in Fig. 2, and hydrodynamic radius (R_H) determined from 2fFCS (red circles) for the donor-labeled variant CypV2C as a function of the denaturant activity, a_{GdmCl} . Error bars for R_G were estimated from the change in I_p^* by $\pm 10\%$. Error bars for R_H represent the standard deviation of ± 0.1 nm estimated from the calibration of the instrument (*SI Appendix*). Solid lines are fits according to $y = y(0) + \gamma a_{\text{GdmCl}} / (K + a_{\text{GdmCl}})$, where y is R_G or R_H , respectively. *Inset*: Arrangement of the foci with parallel and vertical polarization in the 2f-FCS setup (51). (B) R_G/R_H as a function of the GdmCl activity. Error bars result from the error propagation of the uncertainties shown in A. The solid line is the ratio of the fits shown in A.

also by the scaling exponent of $\nu = 0.45 \pm 0.03$. These results support our estimates for the scaling exponents of unfolded *hCyp* from smFRET (Fig. 3A).

Interaction energies and the Tanford transfer model. The determination of the scaling exponents (Fig. 3A) now allows us to compute the absolute values of the intrachain interaction energies ϵ for the six unfolded proteins from the measured transfer efficiencies using Eq. 2. The radius of gyration of the Θ -state, which we found to be $R_{G\Theta} = 0.22 \text{ nm } N^{1/2}$ (Eq. 3), the interaction energy ϵ , and the chain length N then fully determine the phase transition behavior of the unfolded chains within the framework of Sanchez theory (21). A comparison of the experimental data with a numerical evaluation of Eq. 1 in terms of the expansion factor $\alpha = R_G/R_{G\Theta}$ shows how the cooperativity of the collapse transition increases with increasing chain length (Fig. 3B). Strictly speaking, a second-order phase transition of the Landau type is only obtained in the limit of $N \rightarrow \infty$ (21). Hence, for the finite size of the proteins investigated here, with $33 \leq N \leq 163$, the transitions are pseudo-second-order, resulting in a rounding of the transition (21, 52).

Since the absolute value of ϵ depends on specific numerical factors in the theory, it is instructive to investigate the difference between the interaction energies in water, $\epsilon(0)$, and GdmCl solution $\epsilon(a_{\text{GdmCl}})$, respectively, $\Delta\epsilon = \epsilon(0) - \epsilon(a_{\text{GdmCl}})$. The values of $\Delta\epsilon$ determined for the different interdyer variants of length n_{DA} can then be rescaled to the full-length protein (n_{total}) according to $\Delta\epsilon_{\text{total}} = \Delta\epsilon(n_{DA}/n_{\text{total}})^{1/2}$ (*SI Appendix*). $\Delta\epsilon_{\text{total}}$ shows a pronounced dependence on the GdmCl activity for all six proteins (Fig. 5A). The effect of GdmCl on protein chains can be modeled as a preferential interaction of the denaturant with the polypeptide chain (49, 53). This weak-binding model describes the solvation free energy for the polypeptide chain as $\Delta g_{\text{sol}} = -\beta\gamma \log(1 + Ka_{\text{GdmCl}})$, where γ corresponds to the effective number of binding sites for GdmCl molecules, K is the apparent equilibrium constant for binding, and $\beta = (RT)^{-1}$, where R is the ideal gas constant and T is the temperature. Fits with this model provide a good description of the change in $\Delta\epsilon_{\text{total}}$ with GdmCl activity for all proteins investigated here (Fig. 5A). In addition, we find a remarkable agreement of the absolute values of $\Delta\epsilon_{\text{total}}$ with the transfer free energies (Δg_{sol}) of the polypeptide chains from water into GdmCl solutions (54) calculated based on

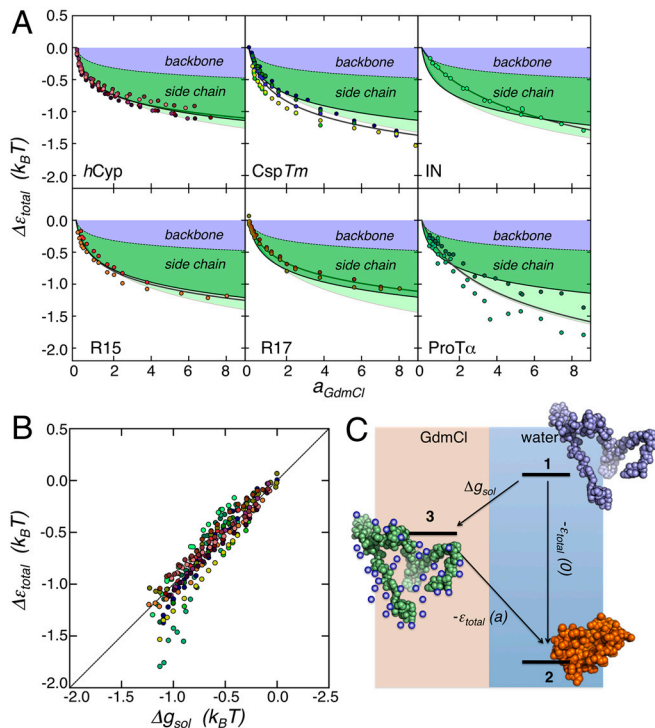


Fig. 5. Relative intrachain interaction energies, $\Delta\epsilon_{\text{total}}$, as a function of GdmCl activity, and comparison between $\Delta\epsilon_{\text{total}}$ and Δg_{sol} . (A) $\Delta\epsilon_{\text{total}}$ for the proteins of this study (circles, colors as in Fig. 3B) together with the fits according to the Schellman weak binding model (gray solid line), and, for comparison, the Tanford transfer free energies Δg_{sol} calculated for the full-length sequences (black line) according to ref. 54. Contributions from the backbone and side chains to Δg_{sol} are shaded in blue and green, respectively. The effect of the δg_{sol} -values estimated for Glu and Asp on Δg_{sol} is indicated as a light green shaded area. From the discrepancy between $\Delta\epsilon_{\text{total}}$ and Δg_{sol} for ProT α , we obtained δg_{sol} for Glu and Asp at 6 M GdmCl to be $-798 \text{ cal mol}^{-1}$ (*SI Appendix, Eq. S14 and Table S2*). (B) Correlation between $\Delta\epsilon_{\text{total}}$ and Δg_{sol} and thermodynamic cycle (C) illustrating the effect of GdmCl on the chain energy as explained in the main text. State 1 is a hypothetical expanded unfolded state in water and state 3 is the same state in the presence of GdmCl. State 2 is the collapsed unfolded state in water.

their amino acid sequences (Fig. 5A and B and *SI Appendix, Fig. S6*). This accordance suggests that the expansion of unfolded proteins, at least for the proteins investigated here, can be explained quantitatively by the change in free energy upon interaction of GdmCl molecules with the chain, implying $\Delta\epsilon_{\text{total}} = \Delta g_{\text{sol}}$. This finding strongly supports the use of this equality in a heteropolymer theory of protein folding (13) and in the molecular transfer model, where it was employed to predict the dimensions of denatured proteins at varying concentrations of GdmCl (14). A simple thermodynamic cycle, in which the total intrachain interaction energy, $-\epsilon_{\text{total}}(0)$, is reduced by the free energy of transferring the amino acid sequence from water to GdmCl (Δg_{sol}), illustrates the effect of GdmCl on the intrachain interaction energy, $-\epsilon_{\text{total}}(a)$, and R_G (Fig. 5C). Finally, these results directly support the correlation between the m -value for folding and the free energy change of collapse predicted by Alonso and Dill (13) and found experimentally by Ziv and Haran (11) (*SI Appendix*).

Effect of sequence composition on the scaling exponent. A detailed analysis of the effect of sequence composition on the scaling exponents of the six proteins in water reveals a pronounced positive correlation between ν and the net charge of the polypeptide (Fig. 6A), and a negative correlation between ν and sequence hydrophobicity (Fig. 6B). A similar correlation has recently been observed in molecular dynamics simulations of protamines,

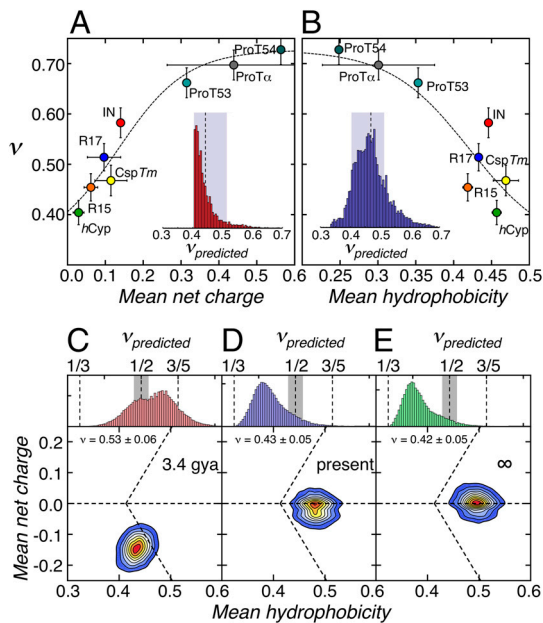


Fig. 6. Scaling exponents, sequence composition, and evolutionary trends. (A) Correlation between the scaling exponents of the proteins and the net charges of their sequences at pH 7. (B) Correlation between the scaling exponents of the six proteins and the mean hydrophobicity of their sequences. Horizontal error bars are the standard deviations as shown in Fig. 1A; vertical error bars reflect the changes in the scaling exponents upon variation of I_p^* by $\pm 10\%$. Dashed lines in A and B are global fits according to empirical equations chosen to give reasonable limits of ν (SI Appendix, Eq. S29). Insets: Frequency histograms of the predicted scaling exponents for the unfolded states of the proteins selected from the PDB shown in Fig. 1A and B based on the fits in A (red) and B (blue), respectively. The shaded areas indicate the regime of scaling exponents between $\nu = 0.40$ and $\nu = 0.51$, which encompass 93% of proteins in A and 71% of proteins in B. (C–E) Distributions of predicted scaling exponents (Top) and mean net charge versus hydrophobicity (Bottom) for 50,000 amino acid sequences drawn randomly from the amino acid frequency distribution of the last universal ancestor (C), current proteins (D), and predicted for the distant future (E). The mean scaling exponents are indicated. See SI Appendix, Eqs. S29–S31 for calculation of the scaling exponents. Amino acid frequencies were taken from table 3 in ref. 60.

positively charged intrinsically disordered peptides (55). These correlations allow us to estimate the scaling exponents also for other proteins. Values of the scaling exponents predicted for the unfolded states of 10,905 monomeric proteins from the Protein Data Bank, based on the correlation between ν and net charge (Fig. 6A, Inset), and ν and hydrophobicity (Fig. 6B, Inset) indicate that the majority of these proteins fall into the range of the scaling exponents observed with the foldable proteins in this study. A value of 0.45 ± 0.03 is obtained as a mean value of the two distributions, remarkably close to the value expected for the Θ -state ($\nu = 1/2$).

Discussion

In order to quantify the thermodynamics of unfolded proteins with polymer theory, information about the Θ -point of the unfolded protein is indispensable (11, 21). Using smFRET, we determined the effective Θ -point of unfolded polypeptide chains by extracting the scaling exponents for four foldable proteins (CspTm, hCyp, R15, R17) and two intrinsically disordered proteins (ProT α and IN). The R_G -values and scaling exponents obtained at high GdmCl are in quantitative agreement with values from SAXS (40) (Fig. 2) and SANS (41), indicating that smFRET is not only a precise but also an accurate method to determine the chain dimensions of unfolded proteins. With the ability to resolve subpopulations, smFRET allows us additionally to obtain the full range of scaling exponents down to physiological solvent conditions.

The higher net charge of the two intrinsically disordered proteins IN and ProT α (Fig. 1A) affects the scaling exponents and leads to an increase of ν at very low GdmCl concentrations (Fig. 3A). The resulting expanded conformations under physiological conditions might reflect an optimization of the sequences for the interaction with their cellular ligands, in keeping with suggestions from theory and simulations that binding kinetics can be accelerated in extended unfolded conformer ensembles (5). In contrast to the IDPs, the scaling exponents of the four foldable proteins decrease monotonically with decreasing solvent quality (Fig. 3A). However, with a mean scaling exponent of 0.46 ± 0.05 in water, they are still much more expanded than a dense globule, which would obey a scaling exponent of $1/3$, as observed for folded globular proteins. Note that the scaling exponents of the two coexisting regimes, folded and unfolded, in water are significantly different ($\nu_{\text{folded}} = 0.34$, $\nu_{\text{unfolded}} \approx 0.46$). Although theory for homopolymers predicts a phase separation into compact globules ($\nu = 1/3$) and expanded chains ($\nu = 1/2$) in poor solvent at high concentrations of the polymer (23), these theories are insufficient to reconcile the two coexisting scaling regimes under our experimental conditions of almost infinite dilution.

In heteropolymer theory, the effective intrachain interaction energy can be approximated by the sum of two mean-field terms, one for backbone interactions (ϵ_{bb}) and one for side-chain interactions (ϵ_{sc}), $\epsilon = \epsilon_{\text{bb}} + \epsilon_{\text{sc}}$. Simulations (29) and experiments (33, 56) suggest that backbone interactions of polypeptide chains are attractive in water, implying that water is a poor solvent for the polypeptide chain backbone with $\epsilon_{\text{bb}} > 1$. Our mean scaling exponent of 0.46 ± 0.05 of unfolded proteins in water (i.e. $\epsilon \approx 1$) (Fig. 3A and B) would then imply that ϵ_{sc} is on average repulsive, i.e. $\epsilon_{\text{sc}} < 0$. Hence, backbone and side-chain interactions nearly compensate in water, leading to a chain close to its critical point. In case the cooperative formation of specific interactions in folded proteins exceeds the mean-field energy term ϵ , compact folded proteins with $\nu = 1/3$ and expanded unfolded proteins with $\nu > 1/3$ can coexist. This scenario is in accord with lattice simulations that suggest that the folding of proteins can occur without populating a dense unstructured globule (57).

What do our results imply for protein folding? Although a collapse to a very dense state ($\nu = 1/3$ and $R_G/R_H = 0.77$) favors folding by reducing the conformational entropy, it could drastically slow down the dynamics of the chain (57) by processes such as internal friction, which have been shown to increase with increasing compaction of unfolded proteins (16, 17, 33, 58). However, especially during the early stages of the folding process, many interactions have to be sampled to find the correct contacts that incrementally decrease the energy of the protein. Simulations based on simple models predict that unfolded chains close to the Θ -regime can accomplish this optimization process more efficiently than chains that are in the completely collapsed globule regime (9, 25, 26). Our results for hCyp, CspTm, R15, and R17 (Figs. 2 and 3), and a comparison of their hydrophobicity and net charge with those of a large number of foldable protein sequences (Fig. 6) implies that natural sequences are indeed close to this regime, and only very few proteins are expected to reach the maximally compact regime with $\nu = 1/3$ in their unfolded state (Fig. 6). However, not only extreme compaction, but also expansion caused by a high net charge of the polypeptide (4, 55) can impede folding, as exemplified by IDPs that are folding incompetent without their biological ligands (48). An intermediate regime of compaction as prevalent in current sequences (Fig. 6) therefore indeed seems most favorable for folding. Within this regime, however, topology-specific effects such as contact order (59) appear to play the dominant role in determining the folding rates of current foldable proteins.

The correlations among net charge, hydrophobicity, and scaling exponents (Fig. 6) finally also allow us to assess the change in average chain dimensions during protein evolution. Based on

bioinformatics analyses (60), ancestral proteins are assumed to have consisted of only eight to ten different amino acids with high average hydrophilicity (Fig. 6 C–E). The resulting scaling exponent of 0.53 ± 0.06 for these ancestral proteins (SI Appendix, Eqs. S29–S31) is close to what we observe for current IDPs, implying that IDPs may be remnants of ancestral protein sequences, whereas foldable sequences with more compact unfolded states are a more recent result of protein evolution (Fig. 6 C–E).

1. Hagen SJ, Hofrichter J, Szabo A, Eaton WA (1996) Diffusion-limited contact formation in unfolded cytochrome c: Estimating the maximum rate of protein folding. *Proc Natl Acad Sci USA* 93:11615–11617.
2. Bieri O, et al. (1999) The speed limit for protein folding measured by triplet–triplet energy transfer. *Proc Natl Acad Sci USA* 96:9597–9601.
3. Schuler B, Lipman E, Eaton W (2002) Probing the free-energy surface for protein folding with single-molecule fluorescence spectroscopy. *Nature* 419:743–747.
4. Müller-Spätth S, et al. (2010) Charge interactions can dominate the dimensions of intrinsically disordered proteins. *Proc Natl Acad Sci USA* 107:14609–14614.
5. Shoemaker B, Portman J, Wolynes P (2000) Speeding molecular recognition by using the folding funnel: The fly-casting mechanism. *Proc Natl Acad Sci USA* 97:8868–8873.
6. Sugase K, Dyson H, Wright PE (2007) Mechanism of coupled folding and binding of an intrinsically disordered protein. *Nature* 447:1021–1025.
7. Chan HS, Dill KA (1991) Polymer principles in protein structure and stability. *Annu Rev Biophys Chem* 20:447–490.
8. Onuchic JN, Luthey-Schulten Z, Wolynes PG (1997) Theory of protein folding: The energy landscape perspective. *Annu Rev Phys Chem* 48:545–600.
9. Thirumalai D, O'Brien E, Morrison G, Hyeon C (2010) Theoretical perspectives on protein folding. *Annu Rev Biophys* 39:159–183.
10. Sherman E, Haran G (2006) Coil-globule transition in the denatured state of a small protein. *Proc Natl Acad Sci USA* 103:11539–11543.
11. Ziv G, Haran G (2009) Protein folding, protein collapse, and Tanford's transfer model: Lessons from single-molecule FRET. *J Am Chem Soc* 131:2942–2947.
12. Bryngelson J, Wolynes P (1990) A simple statistical field-theory of heteropolymer collapse with application to protein folding. *Biopolymers* 30:177–188.
13. Alonso DO, Dill KA (1991) Solvent denaturation and stabilization of globular proteins. *Biochemistry* 30:5974–5985.
14. O'Brien E, Ziv G, Haran G, Brooks B, Thirumalai D (2008) Effects of denaturants and osmolytes on proteins are accurately predicted by the molecular transfer model. *Proc Natl Acad Sci USA* 105:13403–13408.
15. Haran G (2012) How, when, and why proteins collapse: The relation to folding. *Curr Opin Struct Biol* 22:14–20.
16. Waldauer S, Bakajin O, Lapidus L (2010) Extremely slow intramolecular diffusion in unfolded protein L. *Proc Natl Acad Sci USA* 107:13713–13717.
17. Nettels D, Gopich I, Hoffmann A, Schuler B (2007) Ultrafast dynamics of protein collapse from single-molecule photon statistics. *Proc Natl Acad Sci USA* 104:2655–2660.
18. Nettels D, et al. (2009) Single-molecule spectroscopy of the temperature-induced collapse of unfolded proteins. *Proc Natl Acad Sci USA* 106:20740–20745.
19. Schuler B, Eaton W (2008) Protein folding studied by single-molecule FRET. *Curr Opin Struct Biol* 18:16–26.
20. Grosberg A, Kuznetsov D (1992) Quantitative theory of the globule-to-coil transition. 4. Comparison of theoretical results with experimental data. *Macromolecules* 25:1996–2003.
21. Sanchez I (1979) Phase transition behavior of the isolated polymer chain. *Macromolecules* 12:980–988.
22. Flory P (1949) The configuration of real polymer chains. *J Chem Phys* 17:303–310.
23. de Gennes P-G (1979) *Scaling Concepts in Polymer Physics* (Cornell Univ Press, Ithaca, NY and London), pp 113–123.
24. Ha B-Y, Thirumalai D (1992) Conformations of a polyelectrolyte chain. *Phys Rev A* 46:R3012–R3015.
25. Camacho C, Thirumalai D (1993) Kinetics and thermodynamics of folding in model proteins. *Proc Natl Acad Sci USA* 90:6369–6372.
26. Thirumalai D (1995) From minimal models to real proteins: Time scales for protein-folding kinetics. *J Phys (Paris)* 5:1457–1467.
27. Uversky VN (2002) Natively unfolded proteins: A point where biology waits for physics. *Protein Sci* 11:739–756.
28. Crick SL, Jayaraman M, Frieden C, Wetzel R, Pappu RV (2006) Fluorescence correlation spectroscopy shows that monomeric polyglutamine molecules form collapsed structures in aqueous solutions. *Proc Natl Acad Sci USA* 103:16764–16769.
29. Tran HT, Mao A, Pappu RV (2008) Role of backbone-solvent interactions in determining conformational equilibria of intrinsically disordered proteins. *J Am Chem Soc* 130:7380–7392.
30. Uzawa T, et al. (2006) Time-resolved small-angle X-ray scattering investigation of the folding dynamics of heme oxygenase: Implication of the scaling relationship for the submillisecond intermediates of protein folding. *J Mol Biol* 357:997–1008.
31. Kallen J, et al. (1991) Structure of human cyclophilin and its binding site for cyclosporin A determined by X-ray crystallography and NMR spectroscopy. *Nature* 353:276–279.

Materials and Methods

Details of the expression, purification, and labeling of the protein variants and single-molecule measurements are described in detail in the SI Appendix.

ACKNOWLEDGMENTS. We thank Robert Best, Gilad Haran, Rohit Pappu, and Devarajan Thirumalai for helpful discussions. This work was supported by the Swiss National Science Foundation, the Swiss National Center of Competence in Research for Structural Biology, and by a Starting Investigator Grant of the European Research Council.

32. Wensley B, et al. (2010) Experimental evidence for a frustrated energy landscape in a three-helix-bundle protein family. *Nature* 463:685–688.
33. Möglich A, Joder K, Kiefhaber T (2006) End-to-end distance distributions and intrachain diffusion constants in unfolded polypeptide chains indicate intramolecular hydrogen bond formation. *Proc Natl Acad Sci USA* 103:12394–12399.
34. Flory P (1989) *Statistical Mechanics of Chain Molecules* (Carl Hanser Verlag, Munich, Vienna, and New York).
35. Dima R, Thirumalai D (2004) Asymmetry in the shapes of folded and denatured states of proteins. *J Phys Chem B* 108:6564–6570.
36. Theodorou DN, Suter UW (1985) Shape of unperturbed linear-polymers: Polypropylene. *Macromolecules* 18:1206–1214.
37. Tran HT, Pappu RV (2006) Toward an accurate theoretical framework for describing ensembles for proteins under strongly denaturing conditions. *Biophys J* 91:1868–1886.
38. Hammouda B (1993) SANS from homogeneous polymer mixtures: A unified overview. *Adv Polymer Sci* 106:87–133.
39. Zhou H (2004) Polymer models of protein stability, folding, and interactions. *Biochemistry* 43:2141–2154.
40. Kohn J, et al. (2004) Random-coil behavior and the dimensions of chemically unfolded proteins. *Proc Natl Acad Sci USA* 101:12491–12496.
41. Wilkins D, et al. (1999) Hydrodynamic radii of native and denatured proteins measured by pulse field gradient NMR techniques. *Biochemistry* 38:16424–16431.
42. Goldenberg D (2003) Computational simulation of the statistical properties of unfolded proteins. *J Mol Biol* 326:1615–1633.
43. Vitalis A, Wang X, Pappu R (2007) Quantitative characterization of intrinsic disorder in polyglutamine: Insights from analysis based on polymer theories. *Biophys J* 93:1923–1937.
44. Fitzkee N, Rose G (2004) Reassessing random-coil statistics in unfolded proteins. *Proc Natl Acad Sci USA* 101:12497–12502.
45. Zhou H (2002) Dimensions of denatured protein chains from hydrodynamic data. *J Phys Chem B* 106:5769–5775.
46. Damaschun G, Damaschun H, Gast K, Zirwer D (1998) Denatured states of yeast phosphoglycerate kinase. *Biochemistry (Moscow)* 63:259–275.
47. Tanford C, Kawahara K, Lapanje S (1966) Proteins in 6M guanidine hydrochloride: Demonstration of random coil behavior. *J Biol Chem* 241:1921–1923.
48. Uversky V, Gillespie J, Fink A (2000) Why are “natively unfolded” proteins unstructured under physiologic conditions. *Proteins* 41:415–427.
49. O'Brien E, Dima R, Brooks B, Thirumalai D (2007) Interactions between hydrophobic and ionic solutes in aqueous guanidinium chloride and urea solutions: Lessons for protein denaturation mechanism. *J Am Chem Soc* 129:7346–7353.
50. Wu C, Zhou S (1996) First observation of the molten globule state of a single homopolymer chain. *Phys Rev Lett* 77:3053–3055.
51. Dertinger T, et al. (2007) Two-focus fluorescence correlation spectroscopy: A new tool for accurate and absolute diffusion measurements. *Chemphyschem* 8:433–443.
52. Steinhilber MO (2005) A molecular dynamics study on universal properties of polymer chains in different solvent qualities. Part I. A review of linear chain properties. *J Chem Phys* 122:94901–94913.
53. Schellman J (2002) Fifty years of solvent denaturation. *Biophys Chem* 96:91–101.
54. Nozaki Y, Tanford C (1970) The solubility of amino acids, diglycine, and triglycine in aqueous guanidine hydrochloride solutions. *J Biol Chem* 245:1648–1652.
55. Mao AH, Crick SL, Vitalis A, Chicoine CL, Pappu RV (2010) Net charge per residue modulates conformational ensembles of intrinsically disordered proteins. *Proc Natl Acad Sci USA* 107:8183–8188.
56. Teufel DP, Johnson CM, Lum JK, Neuweiler H (2011) Backbone-driven collapse in unfolded protein chains. *J Mol Biol* 409:250–262.
57. Gutin A, Abkevich V (1995) Is burst hydrophobic collapse necessary for protein folding? *Biochemistry* 34:3066–3076.
58. Soranno A, et al. (2012) Quantifying internal friction in unfolded and intrinsically disordered proteins with single molecule spectroscopy. *Proc Natl Acad Sci USA*, doi:10.1073/pnas.1117368109.
59. Plaxco K, Simons K, Baker D (1998) Contact order, transition state placement and the refolding rates of single domain proteins. *J Mol Biol* 277:985–994.
60. Jordan IK, et al. (2005) A universal trend of amino acid gain and loss in protein evolution. *Nature* 433:633–638.
61. Hoffmann A, et al. (2007) Mapping protein collapse with single-molecule fluorescence and kinetic synchrotron radiation circular dichroism spectroscopy. *Proc Natl Acad Sci USA* 104:105–110.