# From Minimal Models to Real Proteins: Time Scales for Protein Folding Kinetics

D. Thirumalai

# From Minimal Models to Real Proteins: Time Scales for Protein Folding Kinetics

D. Thirumalai

Department of Chemistry and Biochemistry and Institute For Physical Science and Technology, University of Maryland, College Park, MD 20742, USA

**Abstract.** — The multipathway mechanism discovered using minimal protein models in conjunction with scaling arguments are used to obtain time scales for the various processes in the folding of real proteins. We consider pathways involving low energy native-like structures as well as direct pathways that proceed via a nucleation mechanism. The average activation barrier separating the low energy structures and the native state is predicted to scale as $\sqrt{N}$ where $N$ is the number of aminoacids in the proteins. In addition estimates of folding times for direct pathways in which collapse and folding are (almost) synchronous are given. It is argued folding sequences whose folding transition temperature is very close to the collapse transition temperature are likely to reach the native conformation rapidly.

## 1. Introduction

The study of minimal models for proteins has produced a very novel (and unconventional) framework for understanding general aspects of the kinetics of in vitro protein folding. The theoretical concepts underlying this framework has recently been referred to as the "new view" [1] of protein folding. A small but representative list of references espousing the "new view" is listed [2–8]. According to the conventional view point [9] protein folding is thought to occur via a sequential mechanism i.e. successive stages in the self-assembly process (characterized, perhaps, by discrete intermediates) are more fully folded and closer to the folded state. More importantly it is believed that the pathway leading the protein from a denatured state to the fully folded state is thought to be unique. In contrast the "new view" starts from the premise that the underlying energy landscape in proteins is rough. This means that there are many minima separated by barriers of differing heights. Thus one is forced to use statistical mechanical methods to understand the dynamics leading to the native state even when considering the folding of a single monomeric protein. In this scenario the protein folding problem amounts to searching for the native state (presumed to be unique) in the complex landscape. The description based on the rugged energy landscape implies that generically a heteropolymer formed randomly from the pool of twenty aminoacids would exhibit kinetic glassy behavior by becoming trapped in a deep non-native minimum. A possible way out of this conundrum is

to assume that in the process of evolution sequences have been designed that have an effective free energy bias toward the native state such that as the unfolding→folding reaction progresses these deep minima are effectively avoided. This general postulate led Bryngelson and Wolynes (BW) to suggest that natural proteins are "minimally frustrated" [2]. Onuchic *et al.* [10] have used these ideas further to characterize the energy landscape in real proteins in terms of the parameters introduced by BW.

The mechanism by which a protein reaches the native state has been studied using minimal models and variety of correlation functions [11]. These studies, which have correlated the folding kinetics with the underlying landscape for foldable sequences, have shown that the approach to the native state occurs in three distinct stages and involves multiple parallel pathways. The purpose of this note is to utilize the results of these model studies, which have only considered small system sizes, and propose scaling laws for extrapolating the time scales and mechanism for larger single domain proteins. This work builds on our earlier suggestion that such scaling laws can be formulated in analogy with the principles in polymer physics [11]. This together with simulations of minimal models can be used to analyze kinetic protein folding experiments.

We would like to emphasize that several aspects of this study are tentative. They have been obtained from the physical picture of folding which emerged from previous studies of minimal models. These models are caricatures of protein and differ from each other in many ways. Despite the variations in these models certain features of folding kinetics appear to be robust. For example both the thermodynamic and kinetic aspect for all these models are qualitatively the same and this is perhaps due to sequence heterogeneity inherent in proteins. It should also be emphasized that all these models and other realistic representations of proteins could lead to different predictions for certain aspects of the problem. For example formation of the precise ordering of secondary structures will crucially depend on the details of the models. One expects on general grounds that these features should not affect the overall tertiary structure formation to which this paper is most directly relevant. In other words the $N$ (number of amino acids) dependence of the kinetic laws obtained here should not significantly depend on these details.

## 2. Three Stage Multipathway Kinetics

Before postulating the scaling laws for the time scales in the three stages of folding as a function of the number of aminoacids, $N$, comprising a protein we briefly discuss the origin of the multiphasic multipathway kinetics. This basic mechanism was originally discovered using Monte Carlo simulations of minimal models of the sort introduced by Dill and coworkers [12]. In these [11,13] and subsequent off-lattice studies [14] only the alpha carbon representation of the protein is utilized. These models are minimal representation of proteins because they contain some, but not all, of the features that play a role in imparting stability to globular proteins. Despite the seeming simplicity of the minimal models it has been shown that the underlying energy landscape is complex sharing many features in common with real proteins [11,13]. Simulations using minimal protein models have revealed two distinct mechanisms for protein to reach the native state. The first one is an indirect route that is described by a three stage multipathway mechanism and the second is a direct pathway which is characterized by a nucleation mechanism. It is crucial to appreciate that between the characteristic temperatures, $T_\theta$ and $T_f$, namely the collapse transition temperature and the folding transition temperature respectively both mechanisms are simultaneously operative. The amplitudes of the two mechanisms depend on temperature [14], chemical conditions [15], sequence, and will be very sensitive to mutations.

We first consider indirect pathways that invariably produce misfolded structures, and are described by a three stage multipathway mechanism (TSMM). The TSMM, thought to be suitable for the folding of single domain proteins, was first observed in a class of lattice [11] and off-lattice models [14] using the following technique. The model protein is quenched from a high temperature (infinite in lattice models) to a temperature close to the folding transition temperature. The approach to the native state is then monitored using several statistical mechanical correlation functions. The ensuing kinetics, averaged over an ensemble of distinct initial conditions, signalling the formation of the native state clearly occurs in three distinct stages. A brief description of each of the phases is given below.

i) Non-specific Collapse: in the first phase the model protein collapses to a compact conformation driven perhaps by the effective attractive interactions between the hydrophobic residues. The relaxation of the various correlation functions in this stage is quite complex and may in fact be consistent with a stretched exponential behavior. In contrast to homopolymers the non-specific collapse in proteins is not totally random. The structures obtained at this stage depend on various factors like loop formation probability, internal motions dominated by dihedral angle transitions etc. Thus in proteins even at this stage there is some specificity that is not expected in homopolymers.

ii) Kinetic Ordering: in the second phase the foldable chain effectively discriminates between the exponentially large number of compact conformations to attain a large fraction of native like contacts. The motion of the various segments in this regime are highly cooperative and in this phase the effect of the biasing forces inherent in the foldable sequences encapsulated in the principle of minimal frustration become operative. At the end of this stage the molecule finds one of the basins corresponding to the minimum energy structures. Although these structures have many native-like contacts they are also different from the native state in a significant way.

iii) All or None: the final stage of folding corresponds to activated transitions from one of the many native-like minimum energy structures to the native state. In this stage the few incorrect contacts are broken and the native contact established. This necessarily involves chain expansion and significant unravelling before the correct contacts can be established. A detailed analysis of several independent trajectories for both lattice and off-lattice simulations suggest that there are multiple pathways that lead to the structures found at the end of the second stage [11,14]. There are relatively few paths that connect the native state and the numerous native-like conformations located at the end of the second stage [11,16]. A summary of the three stage multipathway mechanism is given in Figure 1.

## 3. Time Scales in Protein Folding

In order to make the above scenario applicable to proteins it is necessary to obtain the dependence of the time constants for the three stages as a function of $N$. Proteins are mesoscopic evolved heteropolymer systems with considerable branching. Here we adopt the point of view that from a coarse grained sense one can represent a protein with $N$ effective amino acid residues with $N$ up to roughly 200. Many single domain proteins fall in this category. Given this we wish to make plausible scaling arguments that give explicit expressions for the approximate time scales involved in the distinct stages of folding. The kinetic laws proposed in this paper are applicable when the refolding process is initiated by changing the solvent conditions.
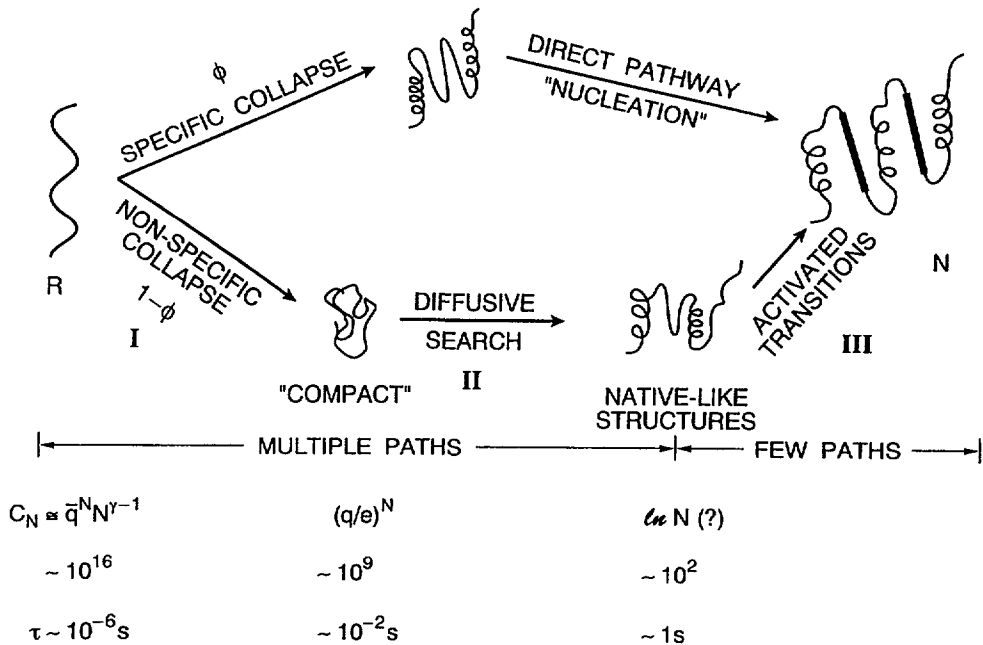
Fig. 1. — This figure summarizes the mechanisms by which proteins reach their native state starting from random coil conformations. This is based on simulation studies of several classes of minimal models. The mechanisms illustrated are applicable for sequences which fold in finite times. In the first process a fraction, $\phi$, of the initial population of molecules reaches the native state directly following collapse. This involves the formation of a critical nucleus. Clearly $\phi$ depends on temperature, ambient conditions, and is also sensitive to mutations. The remaining fraction $(1 - \phi)$ reach the native conformation via a three stage multipath- way mechanism (TSMM). It should be emphasized that this kinetic partitioning between nucleation pathways and TSMM is simultaneously present for a generic protein. The kinetic partition factor depends on intrinsic factors such as sequence as well as environment. By changing the condition one can change $\phi$ as has been done recently in Cytochrome c [15]. The TSMM involves a nonspecific collapse followed by a search among the compact conformations. The search among compact conformations leads the protein to one of the many native-like conformations. Those conformations typically have between (60-80)% of native tertiary contacts. The native-like conformations are often stabilized by incorrect tertiary interactions. The last stage consists of activated transitions from one of the native-like conformations to the native state. There are multiple paths leading to the native state and relatively few paths in the final stage of folding. The number of various conformations available in the three stages as a function of $N$ are given. Their value for $N = 27$ is also shown. The number of minimum energy (native-like) structures should be viewed as a conjecture. This conjecture is based on results reported elsewhere [17]. The last line in the figure gives the time scale for the three processes for $N = 100$ using estimates for physical quantities given in the text.

This is usually accomplished by diluting the concentration of the denaturant so that the folding process can begin. One can also initiate the folding by changing the temperature so that $T \approx T_{\mathrm{f}}$.

3.1. NON-SPECIFIC COLLAPSE. — This stage is analogous to the homopolymer collapse. An approximate time scale for this phase can be obtained using the arguments used by de Gennes

in his treatment of the kinetics of coil to globule transition [17]. It is possible that more recent models of homopolymer collapse could be used to get better estimates for the collapse process itself [18]. This seems worth pursuing. Here we restrict ourselves to the simpler description [17]. The time scale for non- specific collapse is largely determined by the average surface tension between the residues (hydrophobic) and water. In order to utilize de Gennes' arguments we imagine that the protein undergoes a temperature jump from slightly above the theta temperature $T_\theta$ to a value close to the folding transition temperature $T_f$. By generalizing slightly the arguments set forth in reference [17] it is easy to show that $\tau_c$ can be written as

$$\tau_c \simeq \left(\frac{\eta a}{\gamma}\right) \left(\frac{T_\theta - T_f}{T_\theta}\right)^3 N^2 \tag{1}$$

where $\eta$ is the solvent viscosity and a is roughly the persistence length of the protein. Notice that $\tau_c$ is not dependent on the average attractive interactions, $\epsilon_h$, between the hydrophobic residues which are thought to be responsible for the collapse. The value of the surface tension $\gamma$ clearly depends on $\epsilon_h$ and as a result $\tau_c$ depends implicitly on $\epsilon_h$. Using minimal lattice models of proteins we have found that equation (1) is roughly obeyed with the exponent lying between 2 and 2.2.

3.2. KINETIC ORDERING. — As mentioned before in this stage the protein molecule searches among the space of compact structures in such a manner that as the forward reaction proceeds biases in the sequence directs it to structures that are native-like. In our previous study we had shown that these low energy structures could correspond to one of the many minimum energy structures [11, 19]. These structures have considerable native character as measured by the degree of overlap introduced elsewhere [11]. The degree of overlap between a given structure and the native conformation is assessed quantitatively by comparing the distances between non-bonded residues (residues separated by at least three beads) and the corresponding ones in the native conformation. If the distance is less than the bond length then it is presumed that the contact between the specified residues is a native contact. If the fraction of native contacts in a given structure is greater than about 50% these structures are classified as native-like. In fact we typically find that the minimum energy structures have (60 - 80)% of native contacts [11]. The search among the compact structures leading to one of the minimum energy structures (see Fig. 1) has been previously argued to proceed by a diffusive process in the rugged energy landscape [19]. Based on extensive numerical studies we had suggested that the formation of native-like structures from an ensemble of native structures takes place via a series of (largely) local cooperative motions. Furthermore this process was argued to be analogous to reptation in polymer melts [19]. We should emphasize that the local cooperative motions in a monomeric protein is reminiscent of reptation and consequently one may expect similar $N$ dependence to be obtained. With this analogy the approximate time scale for diffusion in subspace of compact structures can be written as [11]

$$\tau_{KO} \approx \tau_D N^\zeta \tag{2}$$

where $\zeta$ is a dynamical folding exponent, and $\tau_D$ is a an undetermined time constant. We have estimated from limited numerical data that $\zeta$ should be approximately 3. This value is consistent with the notion that the diffusive search among a manifold of compact structures leading ultimately to one of the numerous minimum energy structures takes place via a reptation like mechanism. This estimate for $\zeta$ should be regarded as a conjecture and thus it would be desirable to have better estimates of this using minimal models. The time scale $\tau_{KO}$ is somewhat reminiscent of the arguments given by Grosberg et al. [20b] in their description of the second

stage of kinetics of approach to the compact phase of homopolymers following a temperature quench. In our case the physical process that takes place is very different from that of homopolymers. In proteins the second stage leads the system to native-like structures driven mostly by formation of favorable free energy contacts while the compact globule formation in the second stage is dictated by geometrical factors and entanglements [20b].

3.3. ALL OR NONE. — The last stage of the folding process corresponds to an activated transition from one of the many minimum energy native-like structures to the native state. Since there are many native-like structures it follows that the transition state, which occurs late in the folding process for those molecules that reach the native state via TSMM, the transition state is not unique. Because of the number of precursors to the native state it follows that the major driving force even in the late stages for folding is entropic in origin. This concept has been explicitly illustrated using lattice models of disulfide bonded proteins [21]. From a theoretical per- spective this is in harmony with similar notions introduced to describe activated transitions in glasses [22] based on a detailed study of Potts glass models [22]. This analogy to Potts glasses can in fact be utilized to obtain an estimate of the average barrier height separating the minimum energy structures and the native state as a function of $N$. Here we give the arguments based our pervious study of the temperature dependence of relaxation times in systems with rough energy landscape [21, 23]. The diffusive folding time for the last stage may be written as

$$\tau_F \approx \tau_0 e^{\Delta F^{\ddagger}/k_B T} \tag{3}$$

where $\Delta F^{\ddagger}$ is the average free energy separating the native state and the precursors, and $\tau_0$ is an appropriate time constant. The estimate of $\Delta F^{\ddagger}$ for a 46-mer minimal model with a $\beta$-barrel structure as the native state is estimated to be $6k_B T_f$ where $T_f$ is the folding transition temperature [24]. Apparently a smaller estimate ($\Delta F^{\ddagger} \sim 2.4k_B T_f$) has been made for 27-mer three letter code lattice model of proteins [9].

The scaling of $\Delta F^{\ddagger}$ with $N$ can be anticipated by using the following argument. We assume that the free energy distribution of the low energy native-like structures is given by a Gaussian distribution. Since there are an ensemble of independent transition states connecting these native like conformations and the native state it is natural to assume that the barrier height distribution is also roughly Gaussian with the dispersion $\langle \Delta F^2 \rangle$ that scales as $N$. It should be emphasized that this distribution results only after averaging over the free energies of the native-like structures. Since the barrier height distribution is essentially a Gaussian it follows that $\Delta F^{\ddagger} \sim \langle \Delta F^2 \rangle^{1/2} \sim \sqrt{N}$.

The physical picture given above and the analogy with relaxation processes in proteins and glasses can be further used to obtain a more precise scaling behavior of $\Delta F^{\ddagger}$. The relaxation time for glasses is often written as [22]

$$\tau \sim e^{\Delta F^{\ddagger}/k_B(T-T_k)} \tag{4}$$

where $T_k$ is the Kauzmann temperature signalling the vanishing of the configurational entropy. (It should be pointed out that $T_k$ is the same as $T_g$ in the random energy model in that both are computed from the vanishing of entropy, and hence are equilibrium glass transition temperature [25]. From here on we will use $T_k$ and $T_g$ interchangeably.) The formula in equation (4) was derived using scaling arguments based on the physics of droplet excitations in Potts glass models [22, 23]. In reference [20a] it was shown that sufficiently close to $T_k$ the activation barrier scales as

$$\Delta F^{\ddagger}/k_B T_k \sim t^{-1} \tag{5}$$

where $t = (T - T_k)/T_k$. We will assume that this behavior persists at temperatures close to $T_f$. In order to obtain the $N$ dependence of $\Delta F^\ddagger$ we note that the activation physically corresponds to cooperative structural rearrangement on a scale which becomes increasingly large as $T \to T_k$. For proteins this implies that as $T \to T_g$ (or $T_k$) such a length scale $\xi$ would essentially be the whole protein molecule and hence the folding times near $T_g$ would essentially be infinite or more precisely much grater than the biologically relevant time scales. Using the relation $t^{-1} \sim \xi^{1/\nu} \sim N^{1/\nu d}$ we get

$$\Delta F^\ddagger / k_B T_k \sim N^{1/\nu d} \tag{6}$$

where $\nu$ is a correlation length exponent. Because of the possible mapping between the minimal protein models and Potts glasses and REM we take $\nu = 2/d$ [22a]. This yields

$$\Delta F^\ddagger \sim k_B T_k \sqrt{N}. \tag{7}$$

By letting $\alpha = T_f/T_g$ the barrier height can be written as

$$\Delta F^\ddagger \sim \alpha^{-1} k_B T_f \sqrt{N}. \tag{8}$$

There are a few comments about the scaling laws for the three stage multipathway kinetics that are worth making.

a) The estimate of the time scale for non-specific collapse in equation (1) can be made using realistic parameters. Typical ranges of $\eta$, $a$, and $\gamma$ are (0.01-0.1) Poise, (5-10) Å, and (40-60) cal/Å$^2$ mole. This yields $\eta a/\gamma$ in the range $(10^{-10} - 10^{-11})$ s. Thus at $T_f = T_\theta/2 \approx \tau_c$ is found to be between (0.1-10) $\mu$s for $N = 200$. This time is very short and it therefore is reasonable that non-specific collapse would, at subsequent times, invariably lead to misfolded structures.

b) The time constant $\tau_D$ in equation (2) correspondly roughly to the time scale for local dihedral angle transitions [26]. Our previous simulation studies suggest that to a large extent a series of local dihedral angle transitions are responsible for the diffusive search [24]. Typically $\tau_D \approx 10^{-8}$ s. This is also consistent with similar estimates based on lattice MC simulations [26]. Thus $\tau_{KO} \approx 10$ ms. for $N = 100$.

c) Since there are a manifold of different kinds of states in a protein it is worthwhile being specific about the nature of states for which equation (8) is expected to be valid. We expect that the barrier height separating any one of the manifolded of low energy native-like states and the native state to be given by equation (1). The barrier to go from the native state to one of the many low energy structures is given by $\Delta F^\ddagger + \Delta_s$ where $\Delta_s$ is the stability gap. If the analogy with Potts glasses is further utilized then it follows that the relative free energy difference between the distinct low energy states can differ at most by $\sqrt{N}$. Thus for native state to be stable it follows that the stability gap should obey the inequality

$$\Delta_s / k_B T \geq \Omega \sqrt{N} \tag{9}$$

where $\Omega$ is an unspecified constant. The relative free energy stability of the native conformation compared to that of the manifold of higher free energy states at $T = 37$ °C (physiological temperature) for $N = 150$ is roughly 10 kcal/mole according to equation (9) assuming $\Omega \sim O(1)$. It is known that native conformations of single domain proteins ($N \leq 200$) are stable with respect to the other states by only (5-15) kcal/mole which is in approximate agreement with equation (9).

d) It is amusing that the estimate given in equation (8) gives $\Delta F^\ddagger \sim 3.25 k_B T_f$ for $N = 27$ and 4.24 $k_B T_f$ for $N = 46$ using a value of $\alpha \sim 1.6$. These values are in rough accord with the

numerically estimated barrier heights given previously. This agreement should be considered very good because equation (8) is not expected to be accurate for small values of $N$. We believe that equation (8) would provide good estimates for real proteins.

e) It is interesting to calculate the folding times $\tau_F$ in equation (3) for values of $N$ corresponding to typical values in a single domain protein. These roughly range from $N = 50 - 200$. In order to obtain the approximate time scale $\tau_F$ an estimate for $\tau_0$ is needed. It is quite difficult $\tau_0$ obtain an estimate of $\tau_0$ because in general $\tau_0$ is a function of solvent viscosity and internal viscosity of proteins. Here we follow the suggestions of Onuchic et al. [10] and set $\tau_0 = 2\pi\tau_{corr}$ where $\tau_{corr}$ is the correlation time for harmonic fluctuations in a suitably defined reaction coordinate. Onuchic et al. estimate $\tau_{corr} \sim 20,000$ monte carlo steps (MCS) for a three letter 27-mer model [10] which yields $\tau_{corr} \sim 1.26 \times 10^{-4}$ s assuming 1 MCS= $10^{-9}$ s. With this estimate one obtains an expression for $\tau_F$ at $T = T_f$ which is given by

$$\tau_F = 0.126\ e^{0.6\sqrt{N}}\ \text{ms} \tag{10}$$

using $\alpha^{-1} = 0.6$. Equation (10) yields $\tau_F$ in the range (0.05-4) s for $N$ between 100 and 300. These values are consistent with typical experimental findings which estimate that the time constants for the slow phase of the folding is roughly on the order of seconds. In obtaining equation (10) we have assumed that $\tau_{corr}$ is independent of $N$.

f) It is interesting to compare $\tau_{KO}$ (cf. Eq. (2)) and $\tau_F$ (cf. Eq. (10)) for the much studied case of $N = 27$. In this case $\tau_{KO} \sim 2$ ms i.e., the time scales for the second and third stages are practically the same. Folding of this model peptide will appear to be roughly two stage like. In order to clearly distinguish between the various stages simulations for a range of values of $N$ should be performed.

g) The value of $\alpha = T_f/T_g$ has been set to 1.6 in getting equations (8) and (9). These estimates should really be viewed as tentative even though they seem to be supported by simulations of minimal models. In general for foldable sequences it is necessary that $\alpha > 1$. This shows that as $\alpha$ becomes large the effective activation barrier decreases. This would necessarily lead to smaller values of the diffusing folding times $\tau_F$ (cf Eq. (3)). Large values of $\alpha$ imply $T_f/T_g \gg 1$ and the above arguments indicate that fast folding sequences maximize $\alpha$ as anticipated by Bryngelson and Wolynes [27]. The upper bound for $T_f$ is clearly $T_\theta$ the collapse transition temperature because above $T_\theta$ the protein behaves as a denatured species with large values of the radius of gyration. (Here it should be stressed that even near strongly denaturing conditions many protein molecules could retain some element of the secondary structure, and hence the radius of gyration may not correspond to the Flory estimate for self avoiding walk.) Thus fast folding sequences should be characterized by having $T_f$ as close to $T_\theta$ as possible. From the bound

$$\max\ (T_f/T_g) = T_\theta/T_g \tag{11}$$

it follows that there should be correlations between folding rates and

$$t_{CT} = (T_\theta - T_f)/T_\theta. \tag{12}$$

In fact in our earlier simulation studies [11] on lattice models we had discovered that fast folding sequences are characterized by small values of $t_{CT}$. In other words smaller the value of $t_{CT}$ the faster is the folding rate. Although *minimizing* $t_{CT}$ has been derived from the crucial observation of BW that alpha should be *maximized* for practical applications such as de novo design of proteins the criterion of having small $t_{CT}$ may be easier to implement because both $T_f$ and $T_\theta$ can be calculated using equilibrium statistical mechanics.

A note about $T_\theta$ is in order. It might appear that one may not easily apply the concept of $\theta$ point for a purely random heteropolymer by using the vanishing of the effective two

body interaction. The reasons for this argument is that even if the value of the effective two body interaction is obtained by suitable renormalization of the energy scales its value may be positive implying that the system is in a good solvent. However, proteins are not random heteropolymers. The fraction of hydrophobic residues is in slight excess of others so that an effective two body interaction is only positive under strongly denaturing conditions. Minimal model studies have clearly shown that for a given foldable sequence $T_\theta$ may be obtained by conventional methods (by monitoring the temperature dependence of energy fluctuations) and the collapse transition is a finite sized second order phase transition.

## 4. Specific Collapse, Direct Pathways, Nucleation Mechanism

If in fact $t_{CT} \ll 1$ it follows that the process of folding would occur simultaneously with collapse itself or in other words it would be very difficult to differentiate between collapse and folding in typical experiments. In these situations specific collapse is almost synchronous with folding. Let us first emphasize that for proteins which are characterized by several energy scales $t_{CT}$ cannot be identically zero. The physical reason is that protein at $T = T_\theta$ is not compact and at $T_\theta$ the chain behaves like an ideal one. It is known that proteins are compact with the radius of gyration satisfying the inequality $aN^{1/2} < R \leq aN^{1/3}$. It follows that even though $t_{CT}$ cannot be identically zero one can engineer sequences (in principle) such that $t_{CT} \ll 1$. In this case, as mentioned above, at least experimentally collapse and folding cannot be easily separated. Two recent studies in fact suggest that in general there are pathways that exist that directly reach the native state without forming any detectable intermediates. Let us briefly discuss both these studies.

1) The first study probed the kinetics of approach to the native state in a model $\beta$-barrel structure. It [14, 23] was shown that depending on the temperature there exists a fraction of initial trajectories in which the native state is reached rapidly following collapse. This process was shown to occur via the formation of a critical nucleus [14, 24, 28]. In Figure 1 the direct pathway proceeding by a nucleation mechanism is shown as dashed lines.

2) More importantly experiments on Cytochrome c [15] clearly show that these direct pathways exist such that the native state is reached almost simultaneously with the collapse itself. Furthermore these authors have shown that by suitably modifying the chemical conditions indirect pathways leading to the formation of low energy misfolded structures can often be eliminated. This would imply that under suitable chemical and temperature conditions folding would occur in essentially a single step. If this occurs, as the authors of reference [24] have shown for Cytochrome c, then one would infer using our scenario that $T_f \sim T_\theta$. Clearly more experiments would be needed to establish the general validity of these findings.

For the specific collapse, which ensures the correct formation of the native structures and loops, the folding kinetics would be exponential. An estimate for the time constant for specific collapse may be obtained by generalizing de Gennes arguments to include large loop formation probability. This changes equation (1) to

$$\tau_{SC} \simeq \left(\frac{\eta a}{\gamma}\right) \left(\frac{T_\theta - T_f}{T_\theta}\right)^3 N^{2+\theta} \tag{13}$$

where $\theta$ approximately accounts for the probability of long range (approximately the size of protein) tertiary interactions. In three dimensions $1.8 \leq \theta \leq 2.2$ and using the estimate of $10^{-10}$ s for $(\eta a/\gamma)\tau_{sc} \simeq 3$ ms for $N = 100$ at $T_f = T_\theta/2$.

## Conclusions

Analogies between the folding of proteins and other problems in polymer physics and other disordered systems have been used to obtain expressions for the time constants for the various stages in the pathways that inevitably (for entropic reason) lead the protein to one of the low energy misfolded structures. The scaling laws presented here can be verified experimentally. It has been shown using pulsed exchange NMR methods on a number of proteins that the protection kinetics exhibits a fast phase and a slow phase after a crucial transient time of about 5 ms [29]. This has been shown to arise from a combination of nucleation pathway and pathways leading to misfolded structures [14] (see Fig. 1). The fraction $(1 - \phi)$ of molecules follows the TSMM enroute to the native state. For sufficiently large $N$ the third stage would be the rate determining step whose time scale is determined by the estimated barrier height given in equation (8). If the scenario presented in Figure 1 is valid then time scale for the slow phase in refolding experiments would follow the Arrhenius behavior with time being given by $\tau_F \approx \tau_0 e^{\sqrt{N}}$. If these experiments are performed as a function of temperature then the activation barrier heights can be inferred from the temperature dependence of the slow phase [14]. Since the various proteins studied span a reasonable range of $N$ the kinetic laws proposed here can be tested.

A very important theme in the mechanism given in Figure 1 is that natural proteins can fold extremely rapidly on the time scale of collapse itself. In fact suitable modification of the environment can (almost) eliminate the indirect pathways involving kinetic traps. The time scale for this process is on the order of several msecs which roughly coincides with the "burst phase" in pulsed hydrogen exchange labeling experiments. Thus in order to fully understand the detailed mechanism in the direct pathway one needs experiments that probe initial events, namely loop formation between distant residues, which happen on the order of tens of microseconds [30,31]. It appears that experiments ranging from $10^{-6}$ s to s are needed to verify the full consequences of the kinetics of folding as predicted by the "new view".

## Acknowledgments

## References

[1] Baldwin R.L., *Nature* **369** (1994) 183.

[2] Bryngelson J.D. and Wolynes P.G., *Proc. Natl. Acad. Sci. (USA)* **84** (1987) 7524-7528.

[3] Garel T. and Orland H., *Europhys. Lett.* **6** (1988) 307.

[4] Honeycutt J.D. and Thirumalai D., *Proc. Natl. Acad. Sci. (USA)* **87** (1990) 3526-3529.

[5] Shakhnovich E., Farztdinov G., Gutin A.M. and Karplus M., *Phys. Rev. Lett.* **67** (1991) 1665-1668.

[6] Miller R., Danks C.A., Fasolka M.J., Balazs A.C., Chan H.S. and Dill K.A., *J. Chem. Phys.* **96** (1992) 768-780.

[7] Leopold P.E., Montal M. and Onuchic J.N., *Proc. Natl. Acad. Sci. USA* **89** (1992) 8721-8725.

[8] Sali A., Shakhnovich E. and Karplus M.K., *J. Mol. Biol.* **235** (1994) 1614.

[9] Kim P.S. and Baldwin R.L., *Ann. Rev. Biochem.* **51** (1982) 459-489.

[10] Onuchic J.N., Wolynes P.G., Luthey-Schulten Z. and Socci N.D., (Preprint, Jan., 1995).

[11] Camacho C.J. and Thirumalai D., *Proc. Natl. Acad. Sci. USA* **90** (1993) 6369-6372.

[12] Lau K. and Dill K.D., *Macromolecules* **22** (1989) 3986-3997.

[13] Chan H.S. and Dill K.D., *J. Chem. Phys.* **100** (1994) 9238-9257.

[14] Thirumalai D. and Guo Z., *Biopolymers Research Communications* **35** (1995) 137-140.

[15] Sosnick T., Mayne L., Hiller R. and Englander S.W., *Nature* **1** (1994) 149-156.

[16] Bryngelson J.D., Onuchic J.N., Socci N.D. and Wolynes P.G., *Proteins: Structure, Function and Genetics* (1995) in press.

[17] de Gennes P.G., *J. Phys. Lett.* **46** (1985) L639-L642.

[18] Ostrovsky B. and Bar-Yan Y., *Europhys. Lett.* **25** (1994) 409.

[19] Camacho C.J. and Thirumalai D., *Phys. Rev. Lett.* **71** (1993) 2505-2508.

[20] a) de Gennes P.G., *J. Chem. Phys.* **55** (1971) 572; b) Grosberg A., Nechaev S. and Shakhnovich E., *J. Phys.* **49** (1988) 2095.

[21] Camacho C.J. and Thirumalai D., *Proteins: Structure, Function, and Genetics* **22** (1995) 27.

[22] a) Kirkpatrick T.R., Thirumalai D. and Wolynes P.G., *Phys. Rev. A* **40** (1989) 1045-1054; b) Kirkpatrick T.R. and Thirumalai D., *Phys. Rev. B* **37** (1988) 5342-5350.

[23] Kirkpatrick T.R. and Thirumalai D., *J. Phys. I France* **5** (1995) 777.

[24] Guo Z. and Thirumalai D., *Biopolymers* **36** (1995) 83.

[25] Derrida B., *Phys. Rev. B* **24** (1981) 2613-2626.

[26] Thirumalai D., in Statistical Mechanics, Protein Structure, and Protein Substrate Interactions, S. Doniach, Ed. (Plenum Press, 1994) 115-134.

[27] Bryngelson J.D. and Wolynes P.G., *J. Phys. Chem.* **93** (1989) 6902-6915.

[28] Abkevich V.I., Gutin A.M. and Shakhnovich E.I., *Biochemistry* **33** (1994) 10026-10036.

[29] Radford S.E., Dobson C.M. and Evans P.A., *Nature* **358** (1992) 302-307.

[30] Jones C.M., Henry E.R., Hu Y., Chan C.K., Luck S.D., Bhuyama A., Roder H., Hofrichter J. and Eaton W.A., *Proc. Natl. Acad. Sci. USA* **90** (1993) 11860-11864.

[31] Camacho C.J. and Thirumalai D., *Proc. Natl. Acad. Sci. USA* **92** (1995) 1277.