

## Simple model of protein folding kinetics

ROBERT ZWANZIG

Laboratory of Chemical Physics, Building 5, Room 116, National Institutes of Health, Bethesda, MD 20892-0520

Contributed by Robert Zwanzig, July 14, 1995

**ABSTRACT** A simple model of the kinetics of protein folding is presented. The reaction coordinate is the “correctness” of a configuration compared with the native state. The model has a gap in the energy spectrum, a large configurational entropy, a free energy barrier between folded and partially folded states, and a good thermodynamic folding transition. Folding kinetics is described by a master equation. The folding time is estimated by means of a local thermodynamic equilibrium assumption and then is calculated both numerically and analytically by solving the master equation. The folding time has a maximum near the folding transition temperature and can have a minimum at a lower temperature.

The goal of this paper is to present a simple and easily solved generic model of the kinetics of protein folding. This model possesses many of the characteristic features of more realistic models.

Reviews of protein folding theory (1–3) focus on the importance of the energy landscape and its roughness. Indeed, this emphasis is necessary if one wants to understand why some amino acid sequences fold to their native structures easily and others get stuck in metastable states or fold to many different structures. Here we avoid this issue and assume from the beginning that we are dealing with a “good” sequence, one that folds easily and has a unique native structure. Then we can make some simplifying assumptions which allow an easy treatment of the kinetics. The model treated here is essentially the same as in an earlier discussion of Levinthal’s paradox (4), except that the completely folded state is treated as reversibly accessible and not merely as an absorbing sink.

The energy landscape of a protein is its potential energy as a function of many physical coordinates. Folding is a complex motion of the protein on this multidimensional potential surface. The folding process is not necessarily unique—there may be many trajectories or sequences of events that can lead to the native structure. In chemical kinetics one generally prefers, when possible, to deal with a free energy as a function of a single reaction coordinate instead of a potential energy as a function of many coordinates. For this reason, I choose not to describe the configuration of a protein properly by means of its physical coordinates. Instead, I specify its configuration in a cruder way by  $N$  discrete parameters. These parameters can also be regarded as coordinates, but they may have only a distant relation to physical coordinates. For example, each parameter could characterize the immediate environment of a particular amino acid, or a contact between nonneighboring amino acids that occurs in the native structure, or a small set of bond angles in a Ramachandran plot. The model treated here does not require an exact interpretation of these parameters. It is a weakness of the model that no direct connection is made with the structure of any specific protein. On the other hand, it is a strength of the model that the results may have some generality, at least for “good” proteins, independent of actual structures.

Each of the  $N$  parameters can take on any one of  $\nu + 1$  values. The total number of configurations is  $(\nu + 1)^N$ . Of the  $\nu + 1$  values, one is called correct, because it corresponds to the value that parameter has in the native protein or ground state—for example, the right environment of a particular amino acid, the right contact pair, or the right bond angles. The other  $\nu$  values are called incorrect. The number of parameters that have incorrect values in any particular configuration of the protein is denoted by  $S$ . This number is a measure of the distance of any state from the fully correct or native structure. When  $S$  reaches 0, the protein is correctly folded. The importance of correctness as a basic property of any protein configuration was recognized by Bryngelson and Wolynes (5, 6), and current theories and computer simulations of protein folding all make use of analogous quantities (7–12).

**Thermodynamics.** In this model, we assume for simplicity that the energy of a configuration, and the different ways of being incorrect, are determined solely by  $S$ . A simple choice for the energy of the protein, as a function of its distance  $S$  from the correctly folded configuration, is a “smooth funnel,”

$$E_S = SU - \varepsilon \delta_{S0}; S = 0, 1, \dots, N, \quad [1]$$

where both  $U$  and  $\varepsilon$  are assumed to be positive. Positive  $U$  avoids the “golf course landscape” that leads to the Levinthal paradox. Positive  $\varepsilon$ , or an energy gap, is needed to compensate for the low configurational entropy of the correctly folded state. The energy spacing between neighboring values of  $S$  is a constant  $U$  except for the larger energy gap  $U + \varepsilon$  between  $S = 0$  and  $S = 1$ . This is illustrated in Fig. 1, with the arbitrary choice  $N = 100$ ,  $U = 2$ , and  $\varepsilon = 24$ . The degeneracy of the state specified by  $S$  is the number of ways of choosing  $S$  incorrect values,

$$g_S = \nu^S \binom{N}{S}. \quad [2]$$

In this model the partition function is very easy to calculate; it is the sum of a binomial series,

$$Q = \sum_S g_S e^{-\beta E_S} = e^{\beta \varepsilon} + (1 + \nu e^{-\beta U})^N - 1. \quad [3]$$

Because it comes up frequently, we use the notation  $K = \nu e^{-\beta U}$ .

The thermal equilibrium probability of any configuration with  $S > 0$  is

$$P_S(\text{eq}) = \frac{1}{Q} K^S \binom{N}{S}, \quad [4]$$

and the occupancy of the correctly folded configuration is

$$P_0(\text{eq}) = \frac{e^{\beta \varepsilon}}{Q}. \quad [5]$$

This quantity can have a quite sharp “folding transition” as a function of temperature. This is illustrated in Fig. 2, using  $\nu = 2$ ,  $N = 100$ ,  $U = 2$ , and  $\varepsilon = 24$ . (For uniformity, these numerical values are used in all the following illustrations. Similar results

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked “advertisement” in accordance with 18 U.S.C. §1734 solely to indicate this fact.

Abbreviation: LTE, local thermodynamic equilibrium.

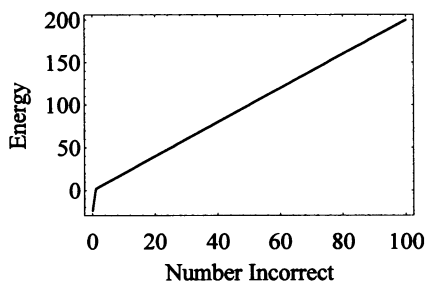


FIG. 1. Energy of a configuration as a function of the number of incorrect parameters. Discrete points are connected by lines to help the eye.

are found for different values of  $\nu$ ,  $N$ , and  $\epsilon$ .) The value of  $U$  is arbitrary and sets the temperature scale. Fig. 2 shows that the midpoint of the transition occurs at  $T = 1.00$ ; the occupancy changes from  $P_0 = 0.04$  at  $T = 1.05$  to  $P_0 = 0.97$  at  $T = 0.95$ .

The folding transition appears here as the result of a competition between entropy, favoring incorrect parameters, and energy, favoring correct parameters. Fig. 3 shows the free energy as a function of  $S$  at various temperatures. Note that the free energy has a local minimum at some nonzero value of  $S$ , another minimum at the correctly folded configuration  $S = 0$ , and a barrier between the two minima that is caused by the decrease of entropy as  $S$  decreases. The relative depth of the two minima is determined by the temperature; low temperatures favor the correctly folded protein. Higher temperatures favor configurations near the other minimum, at nonzero  $S$ , which are partially folded. It is tempting to suggest that the "heat-denatured" state is a compact structure, since it can have a substantial number of correct parameters, but no physical size occurs in this model. Similar free energy curves were seen by Sali *et al.* (12) in computer simulations of lattice models of proteins.

**Folding Kinetics.** So far, the model has been used to calculate equilibrium properties. To treat folding kinetics, we need some rules for moving around on the free energy surface that is determined by  $S$ . Whatever the rules, it is evident that the rate of folding is going to be affected by passage through the free energy barrier or bottleneck at  $S = 1$ . (This implies that the transition state must look very much like the native state, differing from it only in the incorrectness of a single parameter.) Folding is likely to involve a random motion back and forth between the free energy minimum of the partially folded protein and the free energy minimum of the correctly folded protein.

The kinetic model is based on the same general rules as in ref. 4. It appears to be the simplest model that leads to interesting results, and it is amenable to exact mathematical analysis. In this model, any individual configurational parameter can change from correct to incorrect or from incorrect to correct. (If one thinks in terms of actual molecular coordinates, this change could involve many atoms, as in the "crankshaft" moves of computer simulations, and it could involve the interactions of many atoms. We ignore these complexities here.) Then in any transition between configurations,  $S$  can change only by  $+1$  or  $-1$ . The protein does a biased nearest-neighbor random walk on the one-dimensional lattice  $S = 0, 1, 2, \dots, N$ . The probability that the protein is located at  $S$  at time  $t$  is denoted by  $P_S(t)$ . The random walk is described by a master equation for  $P_S(t)$  with transition rates constructed as follows.

A transition from  $S$  to the more correct  $S - 1$  occurs with a rate  $w(S \rightarrow S - 1)$ . This rate is the number  $S$  of incorrect parameters that can be changed, times the rate  $k_1$  of changing any incorrect one to a correct one,

$$w(S \rightarrow S - 1) = Sk_1, \quad [6]$$

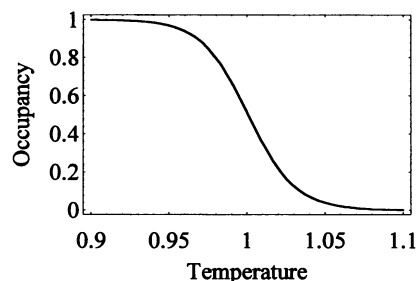


FIG. 2. Occupancy of the correctly folded state as a function of temperature.

for  $S > 0$ . Generally,  $k_1$  depends on temperature; in a kinetically rough landscape it may contain an Arrhenius factor  $\exp(-\Delta E/kT)$ . Transitions from  $S$  to the less correct  $S + 1$  occur at the rate  $w(S \rightarrow S + 1)$  and are determined by the requirement of detailed balance,

$$w(S \rightarrow S + 1)P_S(\text{eq}) = w(S + 1 \rightarrow S)P_{S+1}(\text{eq}). \quad [7]$$

Then for  $S > 0$ , the rate of going from  $S$  to  $S + 1$  is

$$w(S \rightarrow S + 1) = (N - S)Kk_1. \quad [8]$$

In this model, the state  $S = 0$  is treated differently from the way it was done in the earlier treatment (4),

$$w(0 \rightarrow 1) = NKk_1 e^{-\beta\epsilon}. \quad [9]$$

There, to calculate a first passage time, the state was assumed to be fully absorbing. Now we allow a nonzero rate of escape,  $\exp(-\beta\epsilon)NKk_1$ . The limit of infinite  $\epsilon$  leads to the earlier model.

With these rules, the master equation is

$$\begin{aligned} \frac{dP_S}{dt} = & w(S - 1 \rightarrow S)P_{S-1} - w(S \rightarrow S - 1)P_S \\ & + w(S + 1 \rightarrow S)P_{S+1} - w(S \rightarrow S + 1)P_S. \end{aligned} \quad [10]$$

It is easy to verify that this master equation is solved by the thermodynamic equilibrium distribution  $P_S(\text{eq})$  that was given earlier.

**Estimating the Folding Time.** We can estimate the folding time as follows. Because of the large energy gap, equilibration between the states  $S = 0$  and  $S > 0$  is expected to be the rate-determining step. We guess that all states  $S > 0$  come rapidly to local thermodynamic equilibrium (LTE), conditional on the current value of  $P_0(t)$ ,

$$P_S(t) \rightarrow \text{constant} \cdot \binom{N}{S} K^S, \quad [11]$$

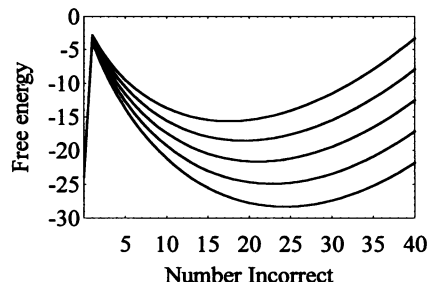


FIG. 3. Free energy as a function of the degree of incorrectness, for various temperatures. From top to bottom, the temperatures are  $T = 0.90, 0.95, 1.00, 1.05$ , and  $1.10$ .

and we find the constant from the time-dependent normalization

$$\sum_{S=1}^N P_S(t) = 1 - P_0(t). \quad [12]$$

Then the guessed-at solution is

$$P_S(t) \cong \frac{1}{Q_0} \binom{N}{S} K^S (1 - P_0(t)), \quad [13]$$

where  $Q_0 = (1 + K)^N - 1$ . In particular this gives a relation between  $P_0(t)$  and  $P_1(t)$ ,

$$P_1(t) \cong \frac{NK}{Q_0} (1 - P_0(t)). \quad [14]$$

When this is substituted in the equation for the rate of change of  $P_0(t)$ , we get

$$\begin{aligned} \frac{dP_0}{dt} &= -k_1 NK e^{-\beta\epsilon} P_0 + k_1 P_1 \\ &\cong k_1 \left[ \frac{NK}{Q_0} (1 - P_0) - NK e^{-\beta\epsilon} P_0 \right]. \end{aligned} \quad [15]$$

This shows a competition between the rate of gain  $k_1 NK/Q_0$ , which decreases as  $T$  increases, and the rate of loss  $k_1 NK \exp(-\beta\epsilon)$ , which increases as  $T$  increases. These two quantities are equal at the folding temperature. Since the overall folding rate is the sum of the gain and loss rates, it has a minimum at the folding temperature (13).

The rate equation can be rewritten as

$$\frac{dP_0}{dt} \cong -\frac{1}{\tau_f} (P_0 - P_0(\text{eq})), \quad [16]$$

where the folding time  $\tau_f$  (i.e., the observed relaxation time) is the reciprocal of the folding rate,

$$\tau_f \cong \frac{(1 + K)^N - 1}{NKK_1} P_0(\text{eq}). \quad [17]$$

This folding time is smaller than the mean first passage time to  $S = 0$ , except in the limit of infinite  $\epsilon$ . The dashed line in Fig. 4 shows the logarithm (base 10) of the folding time as a function of temperature. The same set of parameters ( $\nu = 2$ ,  $N = 100$ ,  $U = 2$ , and  $\epsilon = 24$ ) is used in this example, along with the arbitrarily chosen rate  $k_1 = 10^9 \text{ s}^{-1}$ . A rate  $k_1$  that is independent of temperature corresponds to a hypothetical kinetically smooth funnel—no barriers are encountered when  $S$  decreases. (A kinetically rough funnel will be illustrated later.) Note that the folding time has the expected maximum near the transition temperature.

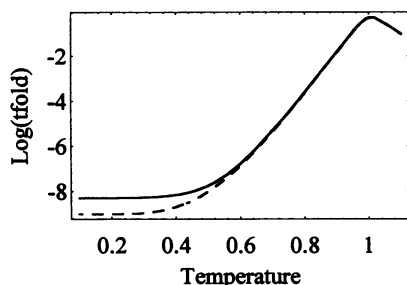


FIG. 4. Logarithm (base 10) of the folding time, in seconds, as a function of temperature. The solid line comes from the exact solution of the master equation, the dashed line from the LTE approximation. Both are based on the rate constant  $k_1 = 10^9 \text{ s}^{-1}$ .

**Analytic Solution.** The argument based on the LTE assumption is open to question. As long as there is a flow into the native state  $S = 0$ , the LTE assumption cannot be strictly true. However, the result agrees fairly well with what one gets from an exact treatment of the master equation. We define a relaxation function  $\varphi(t)$  by

$$P_0(t) = P_0(\text{eq}) + (P_0(0) - P_0(\text{eq}))\varphi(t). \quad [18]$$

In simple first-order kinetics, this function is  $\exp(-t/\tau_f)$ . Whether or not it actually decays exponentially, a good estimate of the folding time is its time integral,

$$\tau_f = \int_0^\infty dt \varphi(t). \quad [19]$$

An exact expression for this quantity, in terms of definite integrals, is derived in *Appendix*. The folding time generally depends on the initial state of the system; for illustrative purposes, we take the special initial condition in which the most incorrect state  $S = N$  is fully occupied.

To check the LTE approximation, we computed the definite integrals numerically. The solid line in Fig. 4 shows the temperature dependence of the logarithm (base 10) of the folding time, for the same choice of parameters as before. The dashed line shows the corresponding LTE estimate. There is some deviation at low temperatures, but on the whole, the LTE approximation appears to work quite well.

**Numerical Solution.** As a further test, the master equation was solved numerically for the same choice of parameters (and specifically for  $T = 0.96$ ). The exact folding time is about 0.2 s. In a very short time period, of the order of  $10^{-8}$  s, the initial condition relaxes to a distribution resembling the prediction of LTE. Over a very much longer time scale, of the order of 1 s, the occupancy of the native state increases steadily to its equilibrium value. The result of numerical integration could not be distinguished from the solution of Eq. 16 using the exact folding time.

**Kinetically Rough Landscape.** In this model, kinetic roughness occurs in two ways. First, it is connected with the free energy penalty for increasing  $S$ ; this has already been taken into account. But decreasing  $S$  can also involve climbing an energy barrier. This can be represented by a rate  $k_1$  containing an extra Arrhenius factor, for example  $k_1 = 10^9 \exp(-\Delta E/kT)$ . Fig. 5 shows the exact folding time, for the same set of parameters as before. The solid line is the earlier result, for  $\Delta E = 0$ , and the dashed line is the folding time when  $\Delta E = 4$ . The folding time passes through a minimum and then increases again as the temperature decreases. In a crude sense, this is like "glassy" behavior. Similar behavior was seen in simulations of lattice models of proteins, in connection with first passage times rather than folding times (14, 15).

**Summary.** A generic picture of protein folding kinetics has been presented. It is based on the idea of the degree of

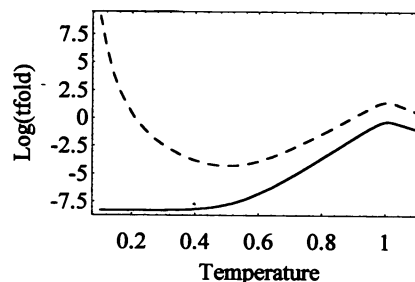


FIG. 5. Logarithm (base 10) of the folding time, in seconds, as a function of temperature. The rate constant is  $k_1 = 10^9 \text{ s}^{-1} \exp(-\Delta E/kT)$ . The solid line uses  $\Delta E = 0$ , and the dashed line uses  $\Delta E = 4$ .

correctness of a protein configuration, but not on any particular definition of correctness. It is based on a simple set of thermodynamic rules, which lead to a free energy that has two minima, one corresponding to the native state and the other to an ensemble of partially folded states. The picture is based on a simple set of kinetic rules, equivalent to a random walk in correctness. After a short induction time, during which the protein arrives at a LTE distribution, the overall kinetic behavior is like that of a two-state system, and the picture provides an estimate of the folding rate. It also allows an analytic treatment of the kinetics, which supports the LTE approximation. What this picture does not do is deal directly with any individual protein; the thermodynamic and kinetic rules may apply only approximately or qualitatively to real proteins. The picture shows only what one might expect in a general way for proteins that fold easily.

### Appendix

**Master Equation.** The master equation can be solved by standard methods of classical analysis (16). (Because  $t$  occurs in the combination  $k_1 t$ , I set  $k_1 = 1$ . The correct  $k_1$  is added later.) First I introduce the generating function

$$G(t, x) = \sum_{S=0}^N x^S P_S(t). \quad [\text{A1}]$$

Then  $G$  satisfies a first-order partial differential equation with an inhomogeneous term. Its solution has two parts, one coming from an initial condition and the other (which vanishes if  $\varepsilon = 0$ ) from the inhomogeneous term,

$$G(t, x) = G^{(0)}(t, x) - (1 - e^{-\beta\varepsilon}) \int_0^t dt' \frac{d\Theta(t', x)}{dt'} P_0(t - t'), \quad [\text{A2}]$$

where the contribution from the initial condition is

$$G^{(0)}(t, x) = \Theta(t, x) G\left(0, \frac{1 + Kx - (1 - x)e^{-(1+K)t}}{1 + Kx + K(1 - x)e^{-(1+K)t}}\right) \quad [\text{A3}]$$

and the function  $\Theta(t, x)$  is

$$\Theta(t, x) = \left[ \frac{1 + Kx + K(1 - x)e^{-(1+K)t}}{1 + K} \right]^N. \quad [\text{A4}]$$

From the definition of the generating function,  $P_0$  is obtained from  $G$  by setting  $x = 0$ ,  $P_0(t) = G(t, 0)$ . Then  $P_0$  satisfies an inhomogeneous integral equation,

$$P_0(t) = \rho_0(t) - (1 - e^{-\beta\varepsilon}) \int_0^t dt' \frac{d\Theta(t', 0)}{dt'} P_0(t - t'), \quad [\text{A5}]$$

where, for brevity, the solution when  $\varepsilon = 0$  is denoted by  $\rho_0$ ,

$$\rho_0(t) = \Theta(t, 0) G\left(0, \frac{1 - e^{-(1+K)t}}{1 + Ke^{-(1+K)t}}\right). \quad [\text{A6}]$$

This equation can be solved by Laplace transforms. The result (using  $z$  as the Laplace transform variable, and denoting transforms by  $\hat{\phantom{x}}$ ) is

$$\hat{P}_0 = \frac{1}{e^{-\beta\varepsilon} + (1 - e^{-\beta\varepsilon})z\hat{\Theta}(z, 0)} \hat{\rho}_0. \quad [\text{A7}]$$

**Folding Time.** To find the complete solution, one must invert a Laplace transform, and this is a tedious undertaking. However, there is a convenient shortcut to the folding time which does not require a Laplace inversion. We use the relaxation function  $\varphi(t)$  that was defined earlier. The folding time is the integral (restoring the factor  $k_1$ )

$$k_1 \tau_f = \int_0^\infty dt \varphi(t) = \lim_{z \rightarrow 0} \hat{\varphi}(z), \quad [\text{A8}]$$

which can be found from the small  $z$  behavior of  $\hat{P}_0(z)$ . For example, if the initial state is where the most incorrect state is fully occupied, or  $P_S(0) = \delta_{SN}$ , then

$$\tau_f = \frac{1}{k_1} \left[ \frac{e^{\beta\varepsilon} - 1}{Q} a_2 - a_1 \right], \quad [\text{A9}]$$

where  $a_2$  and  $a_1$  are given by

$$a_2 = \int_0^\infty dt \{ [1 + Ke^{-(1+K)t}]^N - 1 \}, \quad [\text{A10}]$$

$$a_1 = \int_0^\infty dt \{ [1 - Ke^{-(1+K)t}]^N - 1 \}. \quad [\text{A11}]$$

These are the integrals that were referred to in the section on the analytic solution.

I thank William A. Eaton for helpful remarks.

1. Wolynes, P. G., Onuchic, J. N. & Thirumalai, D. (1995) *Science* **276**, 1619–1620.
2. Bryngelson, J. D., Onuchic, J. N., Socci, N. D. & Wolynes, P. G. (1995) *Proteins Struct. Funct. Genet.* **21**, 167–195.
3. Dill, K. A., Bromberg, S., Yue, K., Fiebig, K. M., Yee, D. P., Thomas, P. D. & Chan, H. S. (1995) *Protein Sci.* **4**, 561–602.
4. Zwanzig, R., Szabo, A. & Bagchi, B. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 20–22.
5. Bryngelson, J. D. & Wolynes, P. G. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 7524–7528.
6. Bryngelson, J. D. & Wolynes, P. G. (1989) *J. Phys. Chem.* **93**, 6902–6915.
7. Honeycutt, J. D. & Thirumalai, D. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 3526–3529.
8. Honeycutt, J. D. & Thirumalai, D. (1992) *Biopolymers* **32**, 695–709.
9. Guo, Z., Thirumalai, D. & Honeycutt, J. D. (1992) *J. Chem. Phys.* **97**, 525–535.
10. Shakhnovich, E. I. & Gutin, A. M. (1990) *Nature (London)* **346**, 773–775.
11. Sali, A., Shakhnovich, E. I. & Karplus, M. (1994) *J. Mol. Biol.* **235**, 1614–1636.
12. Sali, A., Shakhnovich, E. I. & Karplus, M. (1994) *Nature (London)* **369**, 248–251.
13. Creighton, T. E. (1994) in *Mechanisms of Protein Folding*, ed. Pain, R. H. (Oxford Univ. Press, Oxford) pp. 1–25.
14. Chan, H. S. & Dill, K. A. (1994) *J. Chem. Phys.* **100**, 9238–9257.
15. Socci, N. D. & Onuchic, J. N. (1994) *J. Chem. Phys.* **101**, 1519–1528.
16. Montroll, E. W. & Shuler, K. E. (1957) *J. Chem. Phys.* **26**, 454–464.