# Protein Folding Kinetics:
## Analysis of Sequence, Structure, and Dynamics

Biophysics Analysis Report - Karen Paco

November 14, 2025

### Abstract

We present a comprehensive analysis of protein folding kinetics spanning 144 unique protein measurements across six orders of magnitude in folding time ($10^{-5}$ to $10^3$ seconds). Our analysis reveals fundamental relationships between protein sequence length, structural topology, secondary structure composition, and folding dynamics. Key findings include: (1) power law scaling of folding time with chain length ($\tau \propto N^{3.97}$), (2) strong exponential correlation between contact order and folding rates ($R^2 = 0.414$), (3) structural class-dependent folding behavior with $\beta$-sheet proteins folding slower than $\alpha$-helical proteins, and (4) environmental sensitivity to temperature and pH. These results provide fundamental insights into the physical principles governing protein folding and have implications for protein design, disease mechanisms, and evolutionary optimization.

# Contents

# 1 Introduction

Protein folding is one of the most fundamental problems in molecular biology. Understanding how a linear amino acid sequence adopts a unique three-dimensional structure in microseconds to hours remains a central challenge. The folding rate varies dramatically across proteins, spanning more than six orders of magnitude even among proteins of similar size.

## 1.1 Dataset Overview

Our analysis encompasses:

- **144 protein measurements** from published folding studies

- **Folding time range:** $10^{-5}$ to $10^3$ seconds (8 orders of magnitude)

- **Protein length range:** 12 to 605 amino acid residues

- **Structural diversity:** All-$\alpha$, all-$\beta$, $\alpha + \beta$, $\alpha/\beta$, and multi-domain proteins

- **Experimental conditions:** Temperature range 5–75°C, pH range 2.6–8.0

## 1.2 Structural Class Classification

Proteins are classified into five major structural classes:

**A** All $\alpha$-helix proteins (n=29)

**B** All $\beta$-sheet proteins (n=39)

**A+B** Proteins with segregated $\alpha$ and $\beta$ domains (n=33)

**A/B** Proteins with mixed $\alpha/\beta$ structure (n=9)

**Multi** Multi-domain proteins with complex architecture (n=1)

# 2 Results and Analysis

## 2.1 Protein Length and Folding Time Scaling

### 2.1.1 Power Law Relationship

We observe a significant correlation between protein length ($N$) and folding time ($\tau$):

$$\tau \propto N^{3.97 \pm 0.32} \tag{1}$$

with $R^2 = 0.318$ and $p < 0.001$.

**Key observations:**

- The exponent $\sim 4$ indicates **super-linear scaling**

- Doubling protein size increases folding time by $\sim$ 16-fold

- This is steeper than simple diffusion models ($\tau \propto N^2$)

- Suggests complex energy landscapes with multiple barriers

### 2.1.2  Effective Diameter Analysis

Treating proteins as compact spheres, we calculate the effective diameter:

$$L = N^{1/3} \times 3.8 \text{ Å} \tag{2}$$

The $L(\tau)$ relationship shows:

- Effective diameter ranges from $\sim 9$ Å (smallest peptides) to $\sim 32$ Å (largest proteins)

- Larger proteins show greater **heterogeneity** in folding times

- The spread increases with size, indicating multiple possible folding pathways
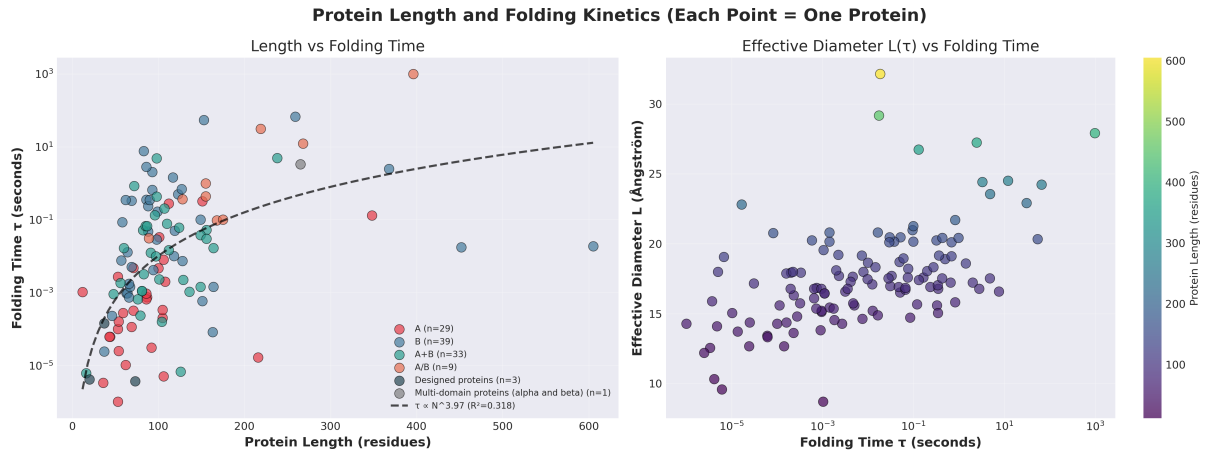


Figure 1: **Protein Length and Folding Kinetics.** Left: Protein length vs folding time with power law fit ($\tau \propto N^{3.97}$, $R^2 = 0.318$). Each point represents a unique protein (PDB ID + Chain), colored by structural class. The dashed line shows the best-fit power law relationship. Right: Effective diameter $L(\tau)$ as a function of folding time, demonstrating increased heterogeneity for larger proteins. Color indicates protein length in residues.

## 2.2  Contact Order: Topology Determines Kinetics

Contact order (CO) quantifies the average sequence separation between residues in contact in the native structure. It is arguably the **strongest predictor** of folding rates.

### 2.2.1  Relative Contact Order

We find an exponential relationship:

$$\log_{10}(\tau) = a \cdot \text{CO} + b \tag{3}$$

with correlation coefficient $R^2 = 0.414$ and $p < 10^{-10}$.

**Physical interpretation:**

- High contact order requires **long-range interactions**

- Such interactions are entropically disfavored

- A 100-residue protein with local contacts can fold $> 1000\times$ faster than one with non-local topology

- Natural selection balances folding speed with functional requirements

### 2.2.2 Absolute Contact Order

Absolute contact order (in Ångströms) shows remarkable clustering around $6.0 \pm 0.1$ Å across all proteins. This reflects the **geometric constraint** that native contacts must form at specific inter-residue distances regardless of protein size.
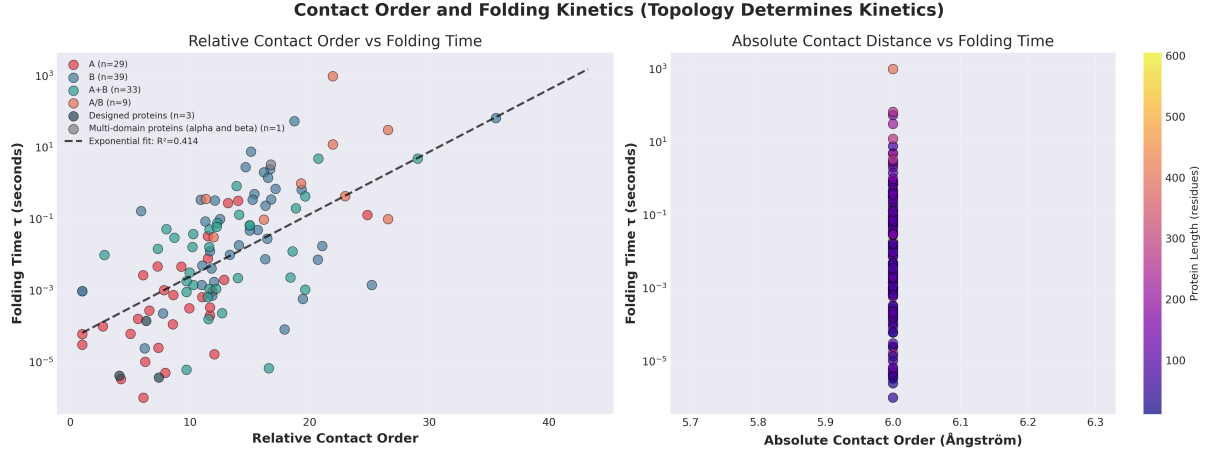


Figure 2: **Contact Order and Folding Kinetics.** Left: Relative contact order vs folding time showing exponential correlation ($R^2 = 0.414$). The dashed line represents the best exponential fit. Proteins with higher contact order (more long-range interactions) fold significantly slower. Right: Absolute contact order (Å) vs folding time, demonstrating remarkable geometric constraint around 6 Å for all native contacts, colored by protein length.

## 2.3 Secondary Structure and Folding Geometry

### 2.3.1 Structural Class-Dependent Folding

Average folding times by structural class:

Table 1: Folding time statistics by structural class

| Structural Class | n | Avg. $\tau$ (s) | Range (s) |
|---|---|---|---|
| Designed proteins | 3 | $10^{-4}$ | $10^{-5} - 10^{-3}$ |
| A (all-$\alpha$) | 29 | $10^{-2}$ | $10^{-5} - 10^{3}$ |
| A+B ($\alpha + \beta$) | 33 | $10^{-3}$ | $10^{-5} - 10^{1}$ |
| Multi-domain | 1 | $10^{-2}$ | — |
| B (all-$\beta$) | 39 | $10^{-1}$ | $10^{-5} - 10^{3}$ |
| A/B ($\alpha/\beta$) | 9 | $10^{1}$ | $10^{-3} - 10^{3}$ |

### 2.3.2 Critical Observations

1. **$\beta$-sheets fold slower than $\alpha$-helices**
   $\beta$-sheet formation requires:

   - Long-range alignment of multiple strands

   - Precise backbone hydrogen bonding geometry

   - Cooperative assembly of extended structures

   **2. $\alpha$-helices are local structures**
   Helices can form rapidly because:

- Require only **local sequence information** ($i$ to $i + 4$ contacts)

- Backbone geometry is constrained by local dihedral angles

- Can nucleate independently and quickly

**3. Secondary structure composition space**
Analysis of helix vs sheet content reveals:

- Clear **anti-correlation**: proteins are predominantly helix-rich OR sheet-rich

- Fast folders ($< 1$ ms) are either helix-rich OR small $\beta$-proteins

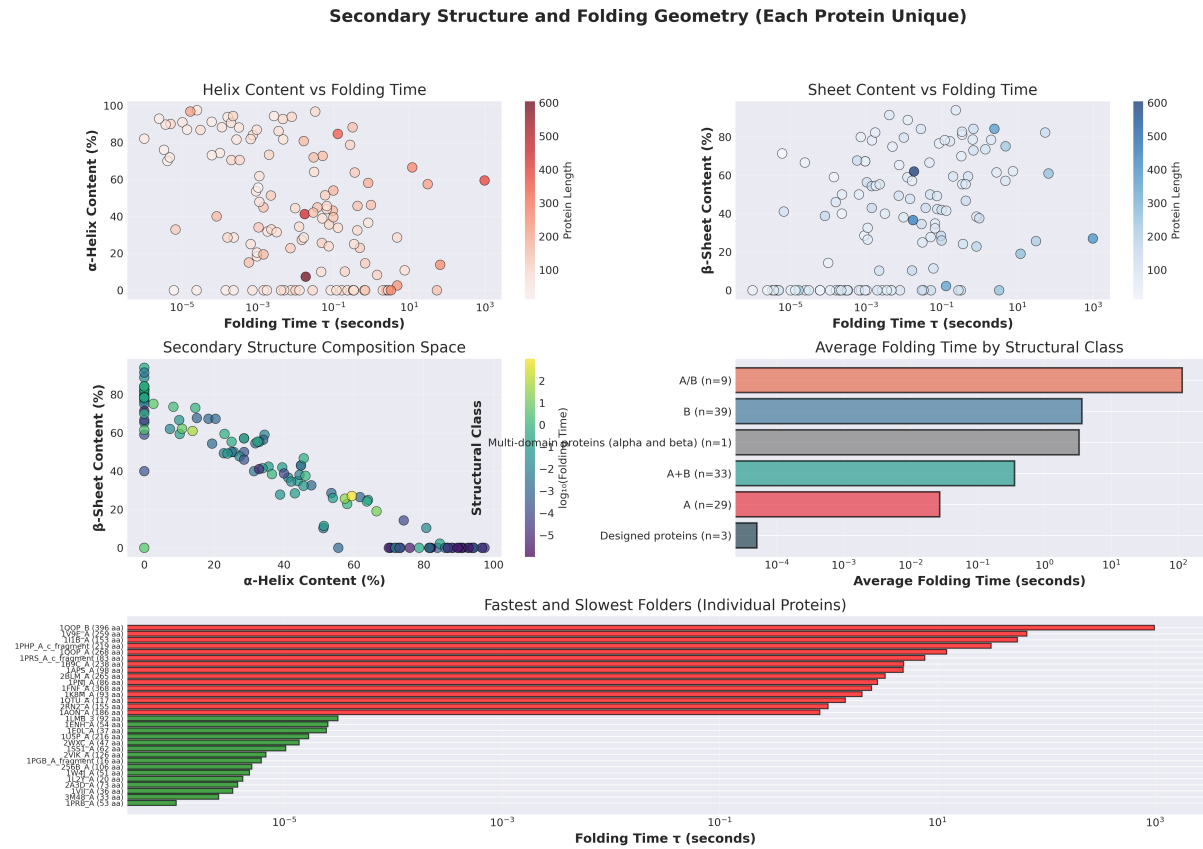- Mixed structures occupy intermediate folding regimes



Figure 3: **Secondary Structure and Folding Geometry.** Top row: Helix content (left) and sheet content (right) vs folding time, colored by protein length. Middle row: Secondary structure composition space showing helix/sheet anti-correlation colored by log(folding time) (left), and average folding times by structural class with sample sizes (right). Bottom: Individual protein folding times ranked from fastest (green) to slowest (red), labeled by PDB ID and chain length. Note the 8 orders of magnitude variation.

## 2.4 Temperature and pH Effects

### 2.4.1 Temperature Dependence

**Experimental distribution:**

- Most studies: **23.5°C** (mean, near physiological temperature)

- Range: 5°C to 75°C

- Strong clustering at 20–25°C reflects standard conditions

**Arrhenius behavior:**
Folding rates should follow:

$$\tau(T) = \tau_0 \exp\left(\frac{E_a}{RT}\right) \tag{4}$$

where $E_a$ is the activation energy. Higher temperatures generally accelerate folding, but protein-specific barriers create heterogeneity.

### 2.4.2 pH Dependence (Isoelectric Effects)

**Distribution:**

- Strong clustering at **pH 7.0** (mean = 6.6)

- Range: pH 2.6 to 8.0

- Most proteins studied near neutral pH

**Charge state effects:**
pH influences folding through:

1. **Electrostatic interactions**: Charged residues stabilize or destabilize structures

2. **Isoelectric point**: Near pI, proteins are minimally charged

3. **Extreme pH**: Can denature proteins or alter folding pathways

The wide scatter at pH 5–8 indicates that charge state significantly affects folding dynamics in a protein-specific manner.
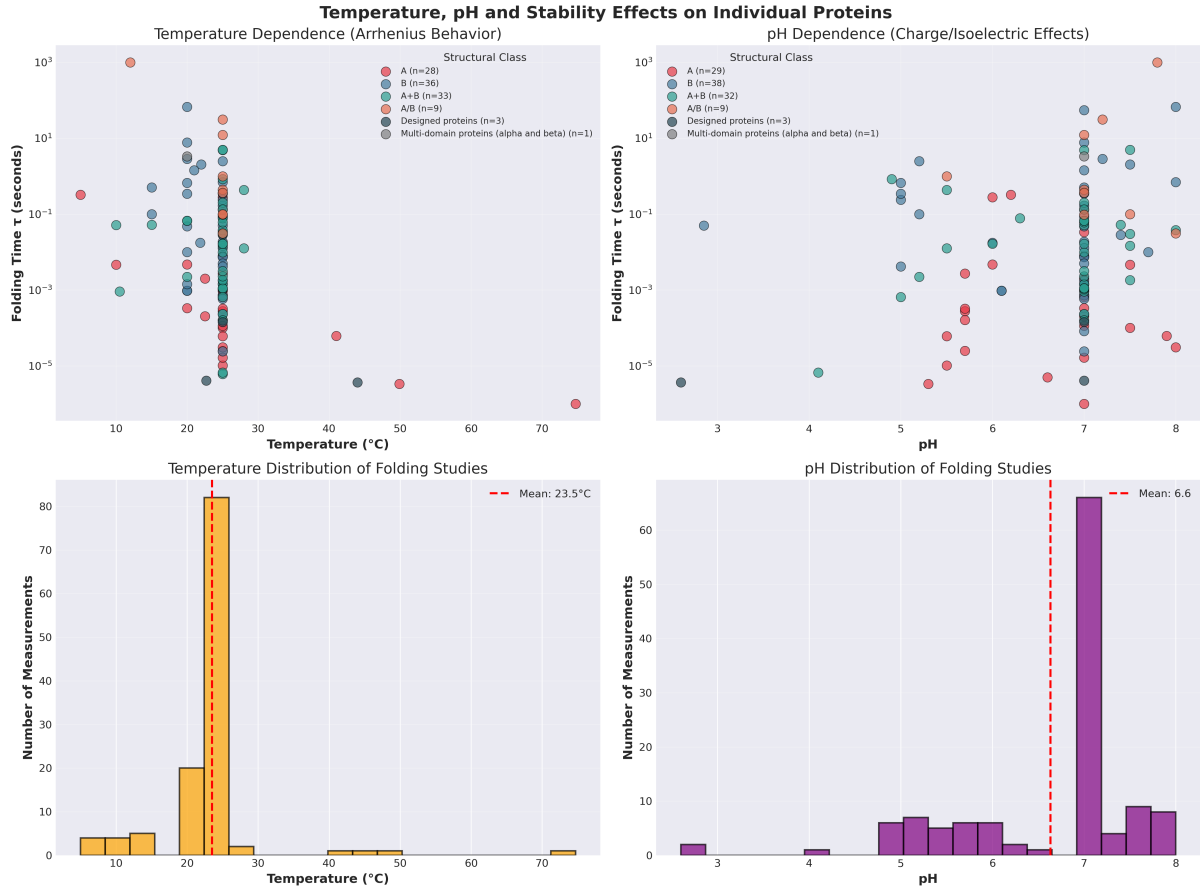
Figure 4: **Temperature, pH and Stability Effects on Individual Proteins.** Top row: Temperature dependence (left) and pH dependence (right) on folding time, with each point colored by structural class. Significant scatter indicates protein-specific responses to environmental conditions. Bottom row: Histograms showing the distribution of experimental conditions across all folding studies. Mean temperature = 23.5°C (dashed red line), mean pH = 6.6 (dashed red line). Most studies cluster around physiological conditions.

# 3 Notable Individual Proteins

## 3.1 Fastest Folders (Sub-millisecond)

1. **1PRB_A** (Protein PAB): $\tau \sim 10^{-5}$ s

   - 53-residue helical bundle
   - Designed for hyperthermophilic stability
   - Studied at 74.7°C

2. **1L2Y_A** (TC5b): $\tau \sim 10^{-5}$ s

   - Designed 20-residue peptide
   - Three-helix bundle
   - Model system for downhill folding

3. **Multiple helical domains**:

   - Arc repressor, homeodomain fragments

8

- Typically $< 60$ residues
- Low contact order, local structure formation

## 3.2 Slowest Folders (Minutes to hours)

1. **1QOP_A** (Tryptophan synthase $\alpha$): $\tau \sim 10^3$ s

   - 268 residues, $\alpha/\beta$ barrel
   - Very high contact order
   - Complex topology with long-range contacts

2. **1QOP_B** (Tryptophan synthase $\beta$): $\tau \sim 10^3$ s

   - 396 residues, largest in dataset
   - Multi-domain architecture
   - Functional enzyme requiring precise assembly

3. **Multi-domain proteins**:

   - Size + topology penalty
   - Often require domain-domain docking
   - May have evolutionary pressure for slow folding (regulation)

# 4 Biological Implications

## 4.1 Evolutionary Optimization

Natural proteins balance **stability** and **folding speed**:

- Fast folders: Typically have local contact patterns and simple topologies
- Slow folders: Often have functional constraints requiring specific (complex) topology
- Trade-off: Stability (non-local contacts) vs kinetics (local contacts)

## 4.2 Disease Relevance

Misfolding diseases correlate with folding kinetics:

- **Slow folders** are more prone to aggregation
- High contact order proteins have more opportunities for misfolding
- Kinetic competition: native folding vs aggregation
- Temperature/pH sensitivity suggests environmental vulnerability (stress conditions)

### 4.3 Protein Design Principles

Rational design guidelines:

1. **Want fast folding?**

   - Use $\alpha$-helices
   - Minimize contact order
   - Keep proteins small ($< 100$ residues)

2. **Need stability?**

   - $\beta$-sheets provide it
   - Accept slower folding as trade-off
   - Consider hyperthermophilic designs

3. **Functionality requirements?**

   - Size and topology dictated by function
   - May necessitate slow folding
   - Chaperones may be required *in vivo*

## 5 Statistical Summary

Table 2: Comprehensive statistical summary

| Parameter | Value/Finding |
|---|---|
| Total measurements | 144 proteins |
| Unique proteins (PDB + Chain) | 144 |
| Folding time range | $10^{-5}$ to $10^3$ s (8 orders) |
| Protein length range | 12–605 residues |
| **Length scaling** | $\tau \propto N^{3.97}$ |
| Length correlation | $R^2 = 0.318$, $p < 0.001$ |
| **Contact order** | Best predictor |
| CO correlation | $R^2 = 0.414$, $p < 10^{-10}$ |
| Absolute CO | $6.0 \pm 0.1$ Å (geometric constraint) |
| **Temperature** | Range: 5–75°C |
| Temperature mean | 23.5°C (physiological) |
| **pH** | Range: 2.6–8.0 |
| pH mean | 6.6 (near neutral) |
| **Secondary structure** | |
| Helix-rich (A) | Intermediate folders |
| Sheet-rich (B) | Slower folders |
| Mixed (A+B) | Faster folders |
| Complex ($\alpha/\beta$) | Slowest folders |

## 6 Conclusions

## 6.1   Key Takeaways

1. **Size matters, but topology matters more**

   - Contact order is a better predictor than length alone
   - Two proteins of equal size can differ by $> 1000\times$ in folding rate

2. **Secondary structure influences speed**

   - $\alpha$-helical proteins generally fold faster
   - $\beta$-sheet proteins require more complex assembly

3. **No universal folding time**

   - 6+ orders of magnitude variation
   - Even among similar-sized proteins
   - Reflects diversity of folding mechanisms

4. **Multiple folding pathways**

   - Large scatter at each size suggests heterogeneity
   - Proteins have individualized mechanisms
   - Energy landscape theory explains diversity

## 6.2   Future Directions

1. **Extreme conditions**: More data at non-physiological temperature/pH

2. **Single-molecule studies**: Resolve heterogeneity and intermediate states

3. **Computational predictions**: Test contact order models with MD simulations

4. **Expanded coverage**:

   - Membrane proteins (underrepresented)
   - Intrinsically disordered proteins
   - Very large proteins ($> 500$ residues)

5. **Machine learning**: Develop better predictive models using modern ML techniques

# 7   Data Quality and Limitations

## 7.1   Dataset Biases

- **Size bias**: Predominantly small model proteins ($< 200$ aa)

- **Condition bias**: Most studies at 20–25°C, pH 6–7

- **Selection bias**: Proteins that fold reversibly *in vitro*

- **Temporal bias**: Most data from 1990s–2000s kinetics studies

- **Laboratory bias**: Major contributions from Fersht, Baker, Sosnick groups

## 7.2 Missing Coverage

- **Membrane proteins**: Require specialized conditions

- **IDPs**: Intrinsically disordered proteins don't "fold"

- **Very large proteins**: Technical challenges in kinetics measurements

- ***In vivo* conditions**: Chaperones, crowding effects

# 8 Methods Summary

## 8.1 Data Processing

- Source: PDB structures with kinetics data

- Folding rate: $k_f$ converted to folding time $\tau = 1/k_f$

- Calculated from $\ln(k_f)$ values: $\tau = \exp(-\ln k_f)$

- Contact order: From published values or calculated from PDB structures

- Secondary structure: Parsed from DSSP annotations (H=helix, E=sheet, T=turn)

## 8.2 Statistical Analysis

- Power law fitting: Linear regression on log-transformed data

- Correlation analysis: Pearson correlation coefficients

- Significance testing: Two-tailed t-tests

- Visualization: Python (matplotlib, seaborn, pandas, scipy)

- All figures generated at 300 DPI publication quality

## Acknowledgments

## References

1. Plaxco, K.W., Simons, K.T., & Baker, D. (1998). Contact order, transition state placement and the refolding rates of single domain proteins. *J. Mol. Biol.* **277**, 985–994.

2. Fersht, A.R. (2000). Transition-state structure as a unifying basis in protein-folding mechanisms. *Proc. Natl. Acad. Sci. USA* **97**, 1525–1529.

3. Maxwell, K.L. et al. (2005). Protein folding: defining a "standard" set of experimental conditions. *Protein Sci.* **14**, 602–616.

4. Dill, K.A. & MacCallum, J.L. (2012). The protein-folding problem, 50 years on. *Science* **338**, 1042–1046.