

Protein Folding

*In Vitro**

Biochemistry 412

February 20, 2007

[*Note: includes computational (*in silico*) studies]

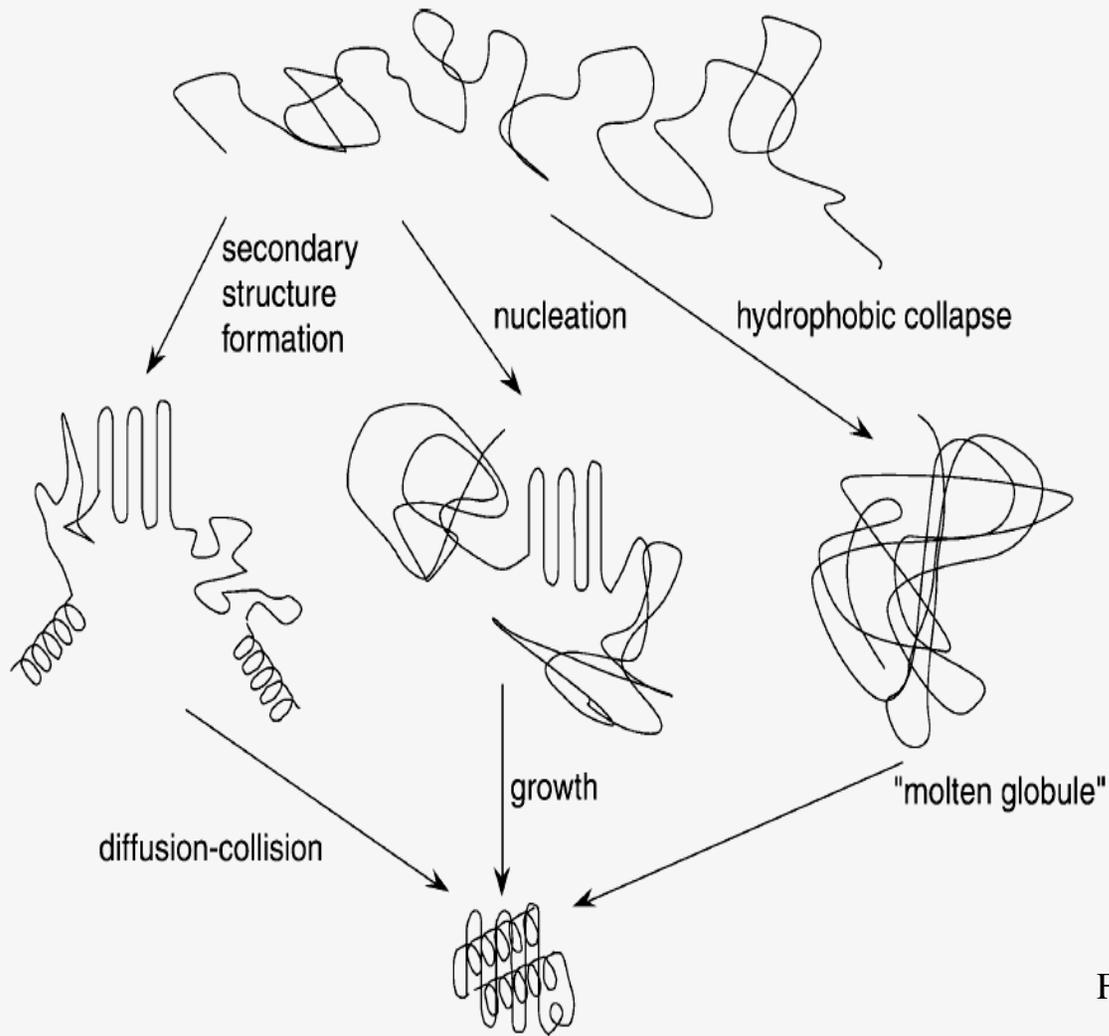


Figure 6. The Three Classical Mechanisms for Protein Folding

At left is pure framework or hierarchical, with secondary structure formed before tertiary. Center is nucleation. At right is tertiary interactions initiating secondary structure formation (modified from Fersht, 1999).

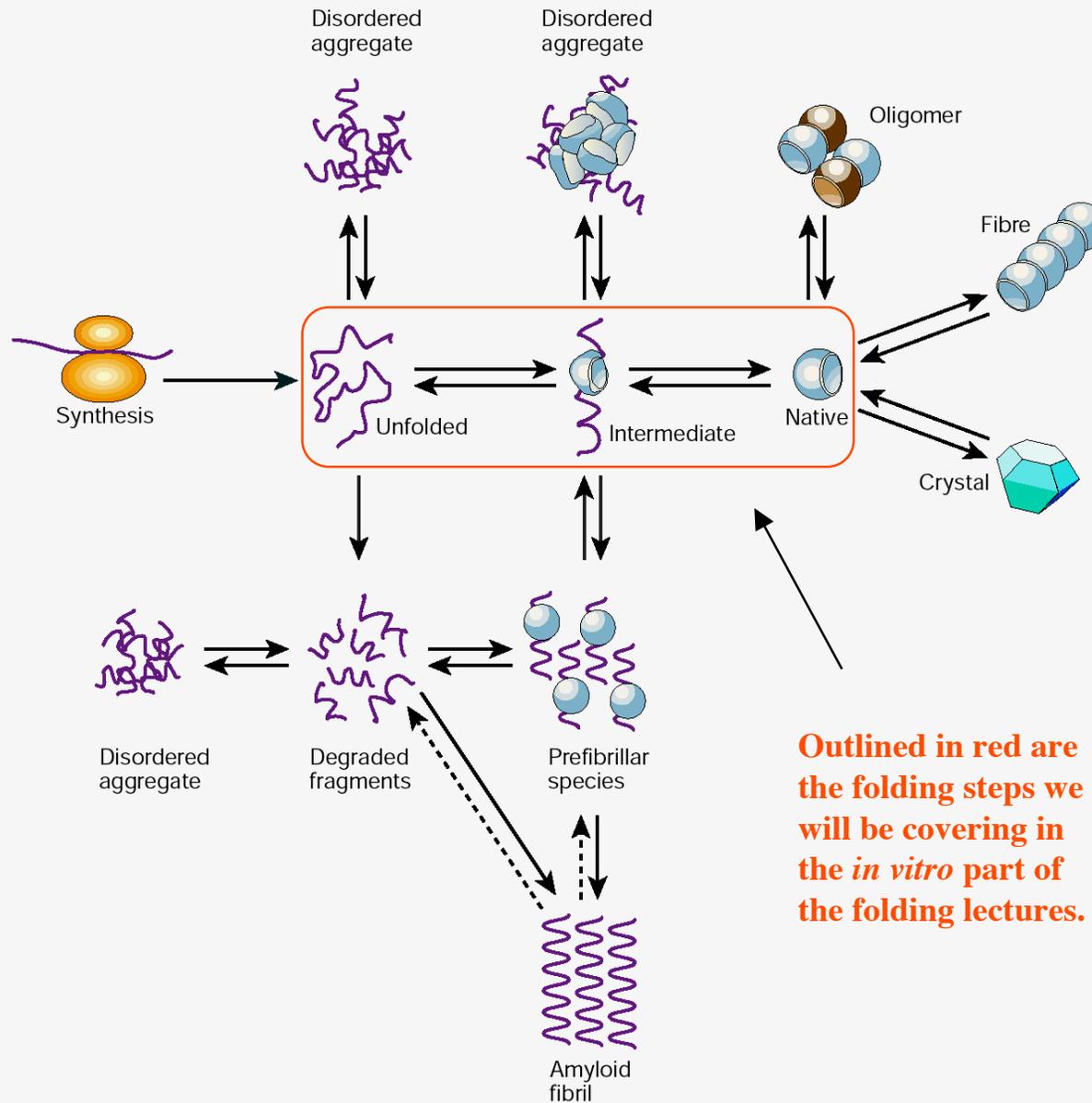
Fersht & Daggett (2002) *Cell* **108**, 573.

Some folding-related facts about proteins:

- **Many small, single domain proteins exhibit simple two-state folding behavior**
- **Most proteins are only marginally stable (5 - 15 kcal/mol) under physiological conditions**
- **Small proteins generally fold very rapidly, often in less than a second**
- **During folding, proteins sample only very few of the total number of possible conformations (see Levinthal's Paradox, below)**

And...

- *It is assumed that a protein's amino acid sequence uniquely determines its native 3D structure*



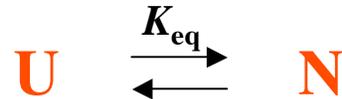
Outlined in red are the folding steps we will be covering in the *in vitro* part of the folding lectures.

Figure 4 A unified view of some of the types of structure that can be formed by polypeptide chains. An unstructured chain, for example newly synthesized on a ribosome, can fold to a monomeric native structure, often through one or more partly folded intermediates. It can, however, experience other fates such as degradation or aggregation. An amyloid fibril is just one form of aggregate, but it is unique in having a highly organized 'misfolded' structure, as shown in Fig. 3. Other assemblies, including functional oligomers, macromolecular complexes and natural protein fibres, contain natively folded molecules, as do the protein crystals produced *in vitro* for X-ray diffraction studies of their structures. The populations and interconversions of the various states are determined by their relative thermodynamic and kinetic stabilities under any given conditions. In living systems, however, transitions between the different states are highly regulated by the environment and by the presence of molecular chaperones, proteolytic enzymes and other factors. Failure of such regulatory mechanisms is likely to be a major factor in the onset and development of misfolding diseases. Adapted from ref. 54.

Two-State Behavior

Energetic and Kinetic Formalisms

Let U signify the unfolded state and N signify the native state:



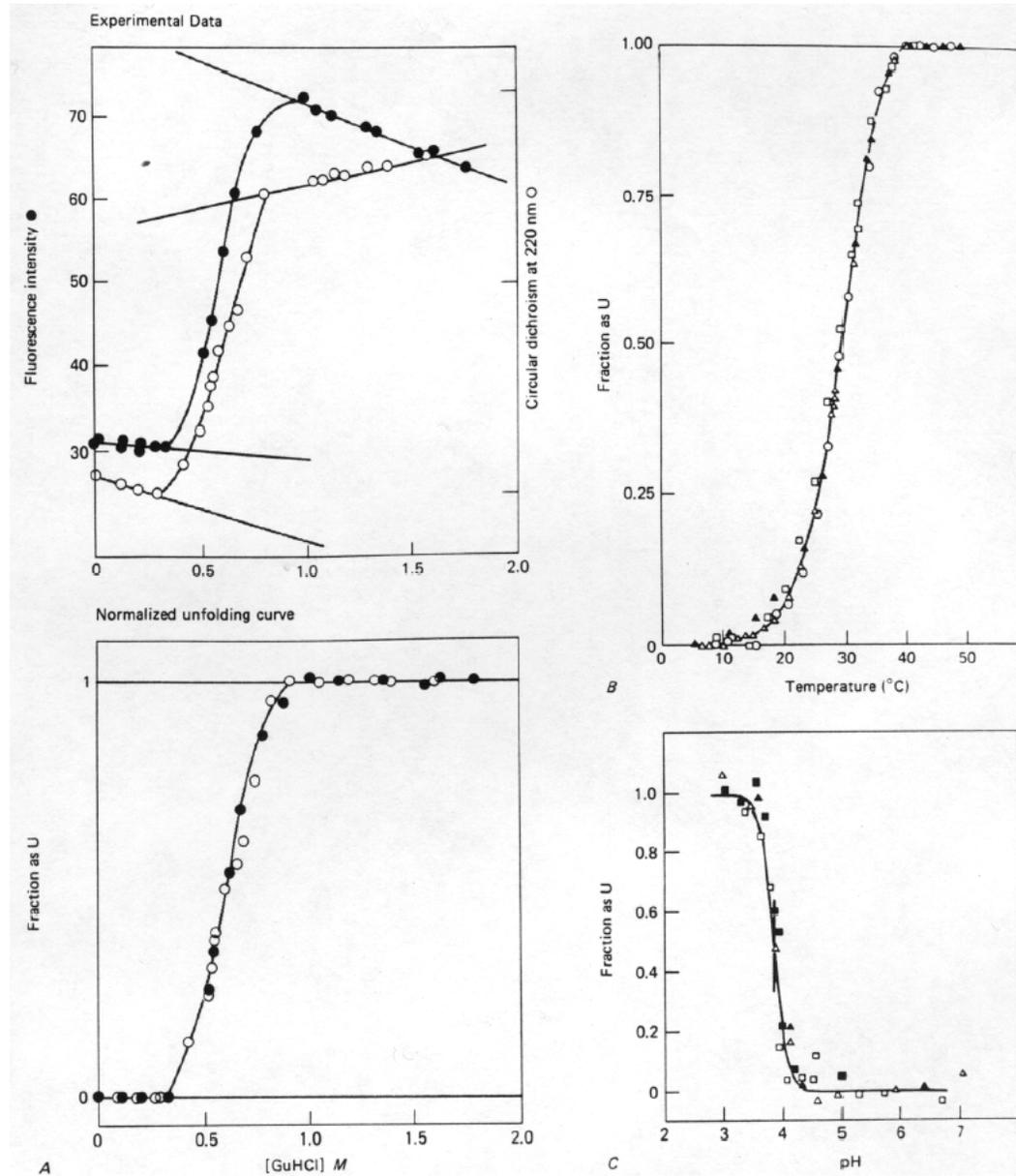
At equilibrium $k_{\text{fold}}[\text{U}] = k_{\text{unfold}}[\text{N}]$

So $K_{\text{eq}} = [\text{N}]/[\text{U}] = k_{\text{fold}}/k_{\text{unfold}}$

And likewise, the stabilization free energy can be expressed as

$$\Delta G^{\circ} = G_{\text{N}}^{\circ} - G_{\text{U}}^{\circ} = -RT \ln K_{\text{eq}}$$

Experimental (equilibrium) unfolding of proteins



Creighton
Proteins
W. H. Freeman
1984, p. 288

Note that K_{eq} is a function of the denaturant concentration, since denaturants *by definition* shift the equilibrium toward the unfolded state.

In fact, $\ln K_{\text{eq}}$ can be approximated as a linear function of the denaturant concentration, i. e.,

$$\ln K_{\text{eq}} = \ln K_{\text{eq}}^{\text{H}_2\text{O}} - c[\text{denaturant}]$$

Where c is a constant for a given protein and set of conditions.

The reason this is important is that linear plots enable you to accurately measure stabilization free energy differences between two different proteins (e. g., between a wild type protein and its mutant).

The Protein Folding Problem: Levinthal's Paradox

- Ribonuclease (124 residues) can potentially form about 10^{50} conformations. If it tries a different conformation every 10^{-13} seconds, it would take $10^{50}/10^{13} = 10^{37}$ seconds or $\sim 10^{30}$ years to try all conformations, yet ribonuclease folds in ~ 1 minute.
- There must be **pathways of folding** with sequential, dependent steps (intermediates), instead of a random “sampling” of all possible conformations.

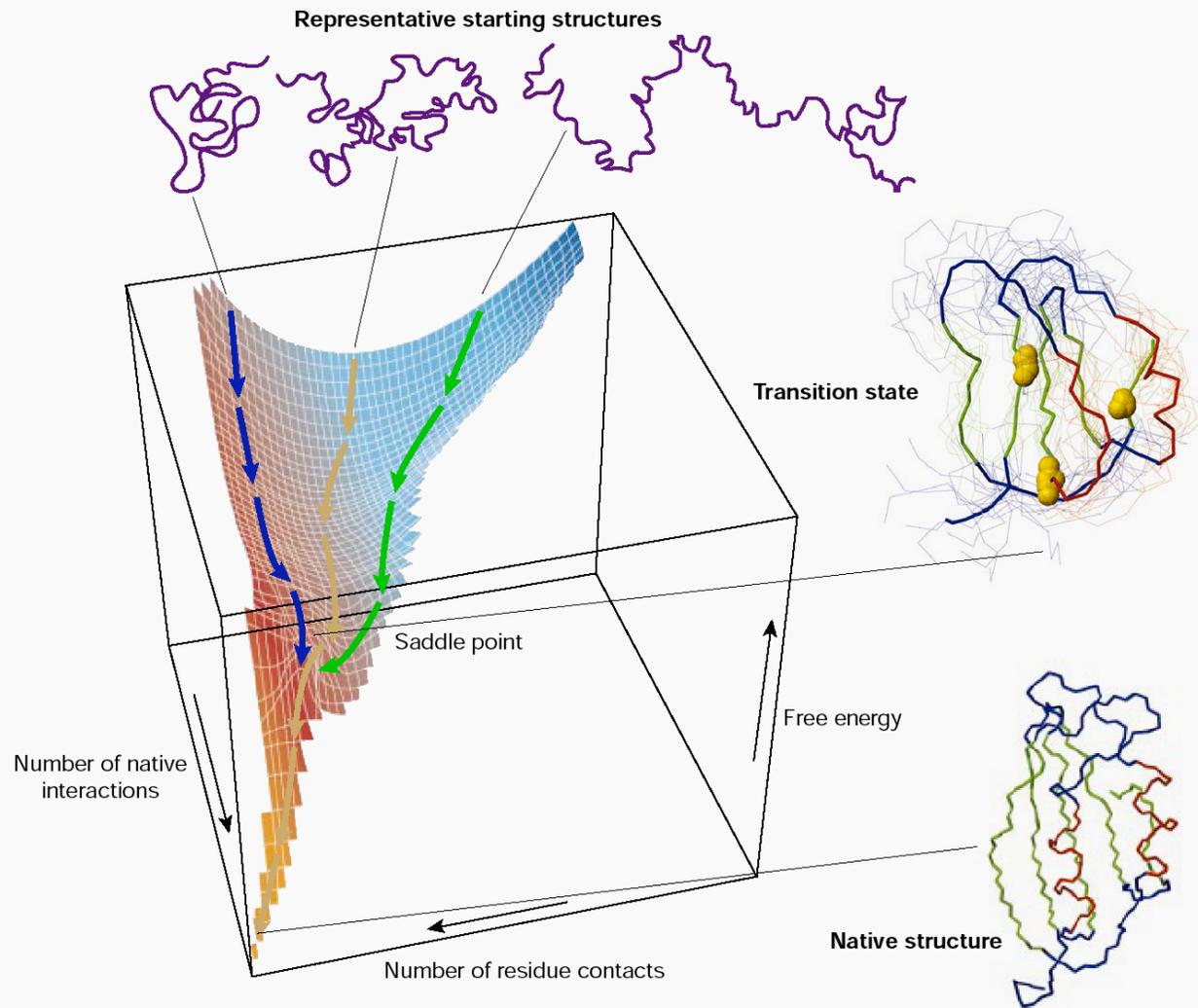
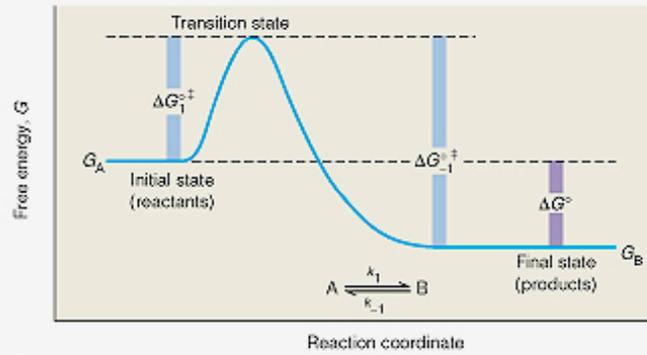
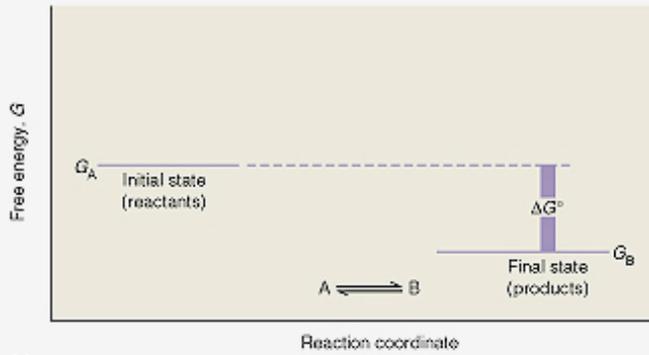


Figure 1 A schematic energy landscape for protein folding. The surface is derived from a computer simulation of the folding of a highly simplified model of a small protein. The surface 'funnels' the multitude of denatured conformations to the unique native structure. The critical region on a simple surface such as this one is the saddle point corresponding to the transition state, the barrier that all molecules must cross if they are to fold to the native state. Superimposed on this schematic surface are ensembles of structures corresponding to different stages of the folding process. The transition state ensemble was calculated by using computer simulations constrained by experimental data from mutational studies of acylphosphatase¹⁸. The yellow spheres in this ensemble represent the three 'key residues' in the structure; when these residues have formed their native-like contacts the overall topology of the native fold is established. The structure of the native state is shown at the bottom of the surface; at the top are indicated schematically some contributors to the distribution of unfolded species that represent the starting point for folding. Also indicated on the surface are highly simplified trajectories for the folding of individual molecules. Adapted from ref. 6.

Dobson (2003) *Nature* **426**, 884.

Generalized Free Energy Diagrams

[for folding, let N (native state) = "B" = "P"
and U (unfolded state) = "A" = "S"]

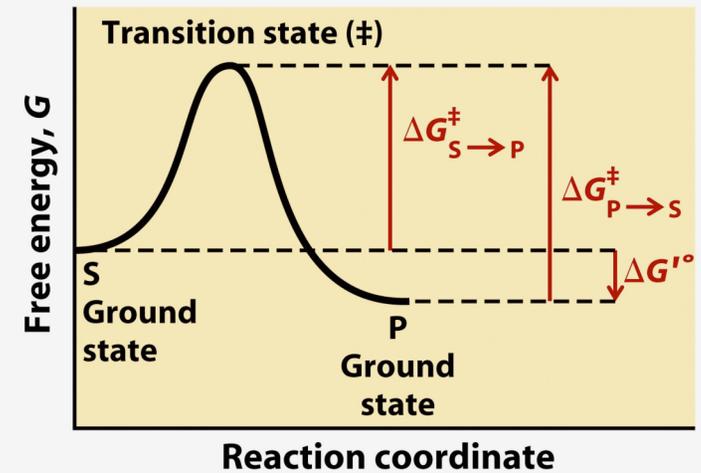
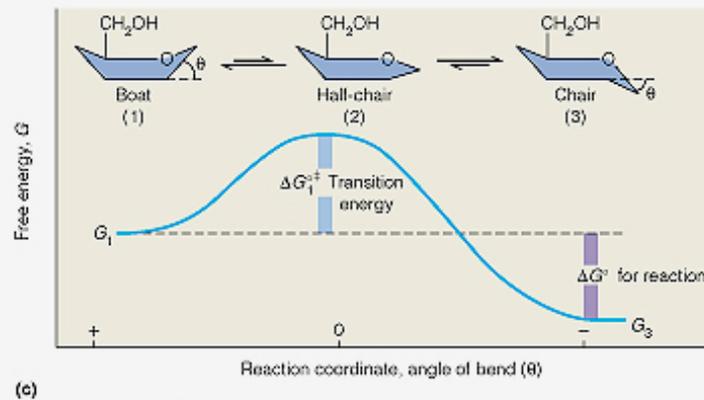


$$K = e^{-\Delta G^\circ/RT}$$

For $A \rightleftharpoons A^\ddagger$

$$[A]^\ddagger/[A]_0 = e^{-\Delta G^\ddagger/RT}$$

$$[A]^\ddagger = [A]_0 e^{-\Delta G^\ddagger/RT}$$



Copyright © 2000 Benjamin/Cummings, an imprint of Addison Wesley Longman, Inc.

K = equilibrium constant

\ddagger = transition state

$[A]^\ddagger$ = concentration of molecules having the activation energy

$[A]_0$ = total concentration

$-\Delta G^\circ$ = standard free energy change of activation (activation energy)

Note that the transition state (TS) energy, G^\ddagger , can be indirectly measured based on its difference with the unfolded and native state free energies.

Thus,

$$\Delta G_{\text{TS-U}} = G^\ddagger - G_{\text{U}}^\circ \quad \text{and} \quad \Delta G_{\text{N-TS}} = G_{\text{N}}^\circ - G^\ddagger$$

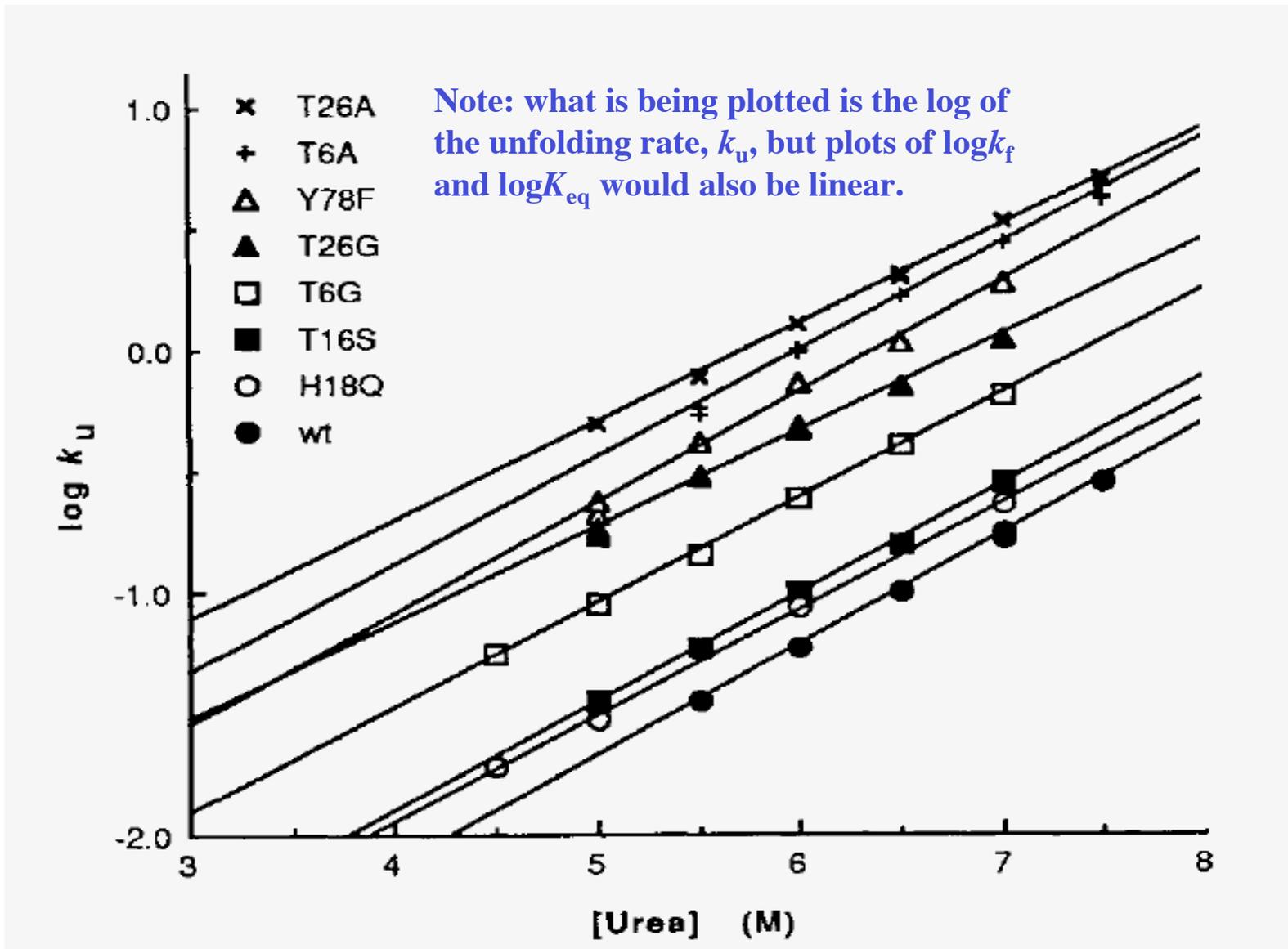
And

$$\Delta G_{\text{TS-U}} = -RT \ln k_{\text{fold}} \quad \text{and} \quad \Delta G_{\text{N-TS}} = RT \ln k_{\text{unfold}}$$

Now....

Protein engineering rears its head!

Denaturation Data for Barnase Mutants



Matouschek et al (1989) *Nature* **340**, 122.

TABLE 1 Changes in activation ($\Delta\Delta G_U^\ddagger$) and free energies of unfolding ($\Delta\Delta G_U$) on mutation of barnase

Mutation	Function of position	$\Delta\Delta G_U$ kcal mol ⁻¹	$\Delta\Delta G_U^\ddagger$ (at 4 M urea) kcal mol ⁻¹	$\Delta\Delta G_U^\ddagger$ (in H ₂ O) kcal mol ⁻¹	$\Delta\Delta G_U^\ddagger/\Delta\Delta G_U$ (at 4 M urea)	$\Delta\Delta G_U^\ddagger/\Delta\Delta G_U$ (in H ₂ O)
Thr → Gly 6	N cap*	1.34	0.86	1.01	0.64	0.76
Thr → Ala 6	N cap*	2.23	1.66	1.76	0.74	0.79
Thr → Gly 26	N cap*	1.58	1.33	1.67	0.84	1.06
Thr → Ala 26	N cap*	2.14	1.91	2.19	0.89	1.02
His → Gln 18	C cap—charge/helix dipole	1.60	0.23	0.36	0.14	0.22
Thr → Ser 16	Hydrophobic on helix surface	1.87	0.28	0.37	0.15	0.20
Tyr → Phe 78	Bridges loop	1.50	1.39	1.41	0.92	0.94
Leu → Ala 14	Hydrophobic core	4.80	1.91	1.86	0.40	0.39
Ile → Val 88	Hydrophobic core	1.49	0.30	0.28	0.20	0.19
Ile → Ala 88	Hydrophobic core	4.46	0.68	0.33	0.15	0.07
Ile → Val 96	Hydrophobic core	0.98	0.48	0.55	0.49	0.56
Ile → Ala 96	Hydrophobic core	3.52	0.74	0.60	0.21	0.17

All data at 25 °C and pH 6.3. Values of $\Delta\Delta G_U$ are taken from refs 11, 12 and unpublished data (L.S. and A.R.F; manuscript in preparation), apart from those for Thr → Ser 16 and Tyr → Phe 78 which were measured as described in ref. 12.* The nomenclature of ref. 21 for residue that occurs at the N terminus of α -helix and can make hydrogen bonds with the backbone NH groups not involved in intrahelical backbone hydrogen bonding²⁰.

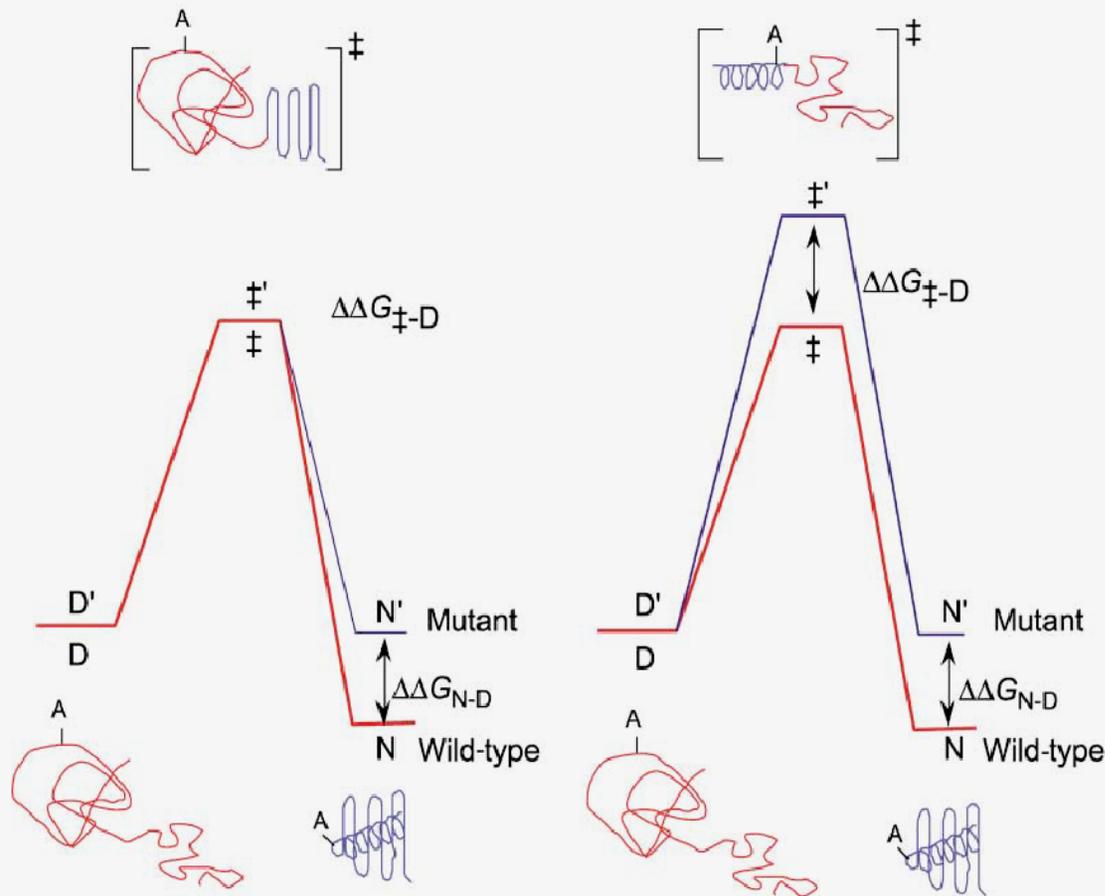


Figure 1. Reaction Profile and ϕ Value Analysis

Schematic profiles are sketched in red for a protein that has an alanine residue (A) in a helix, and in blue for a mutant in which the alanine is mutated to a glycine.

(Left) The transition state (\ddagger), at the top of the energy profile, has the helical region as denatured as in the denatured state D. The energy of the transition state is affected by A to G by the same energy as in D, and so the change in energy of \ddagger relative to that of D, $\Delta\Delta G_{\ddagger-D}$, is 0. Thus, $\phi = \Delta\Delta G_{\ddagger-D}/\Delta\Delta G_{N-D} = 0$.

(Right) The opposite case when the helix is fully structured in the transition state has $\Delta\Delta G_{\ddagger-D} = \Delta\Delta G_{N-D}$, and so $\phi = 1$ (modified from Fersht, 1999). The value of $\Delta\Delta G_{\ddagger-D}$ is calculated from the ratio of rate constants for folding of wild-type ($k_{f(wt)}$) and mutant ($k_{f(mut)}$) proteins [$\Delta\Delta G_{\ddagger-D} = RT \ln(k_{f(wt)}/k_{f(mut)})$]. The value of $\Delta\Delta G_{N-D}$ is calculated by subtracting the free energy of folding of wild-type protein ($\Delta G_{N-D(wt)}$) from that of mutant ($\Delta G_{N-D(mut)}$). The free energies of folding are usually measured from urea-, guanidinium chloride-, or thermal-denaturation curves. ϕ value analysis requires measuring rate and equilibrium constants.

Φ -Value Analysis

$$\Phi = (\Delta G_{\text{TS-D}} - \Delta G'_{\text{TS-D}}) / (\Delta G_{\text{N-D}} - \Delta G'_{\text{N-D}}) = \Delta\Delta G_{\text{TS-D}} / \Delta\Delta G_{\text{N-D}} \quad (\text{Eqn1})$$

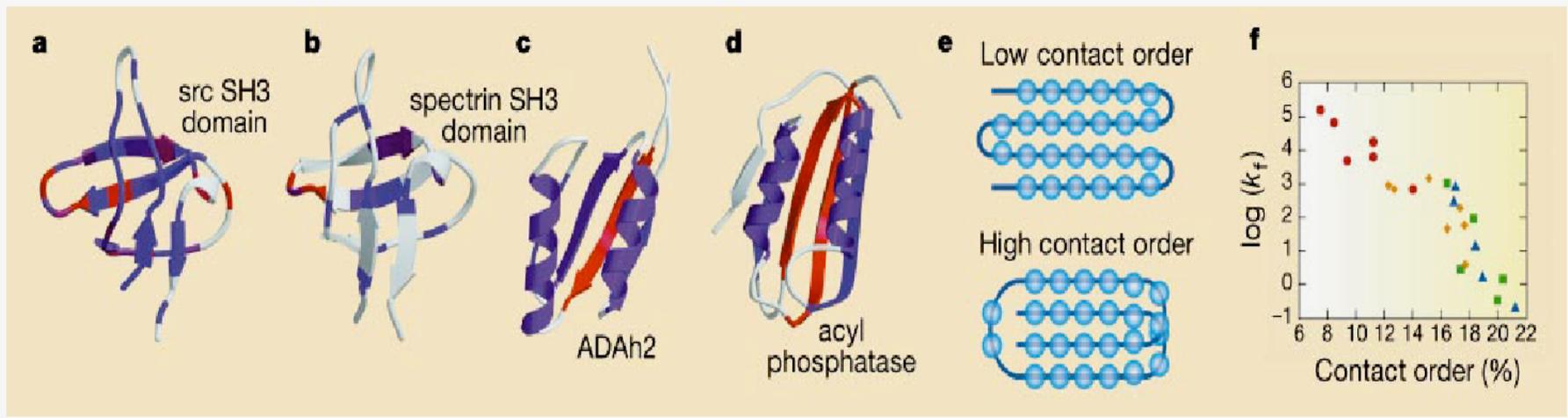
where $\Delta G_{\text{TS-D}}$ and $\Delta G_{\text{N-D}}$ are the free energies of the transition state (this could also be an intermediate state) and the native state, respectively, relative to the denatured state for the wild-type protein (Fig. 1).

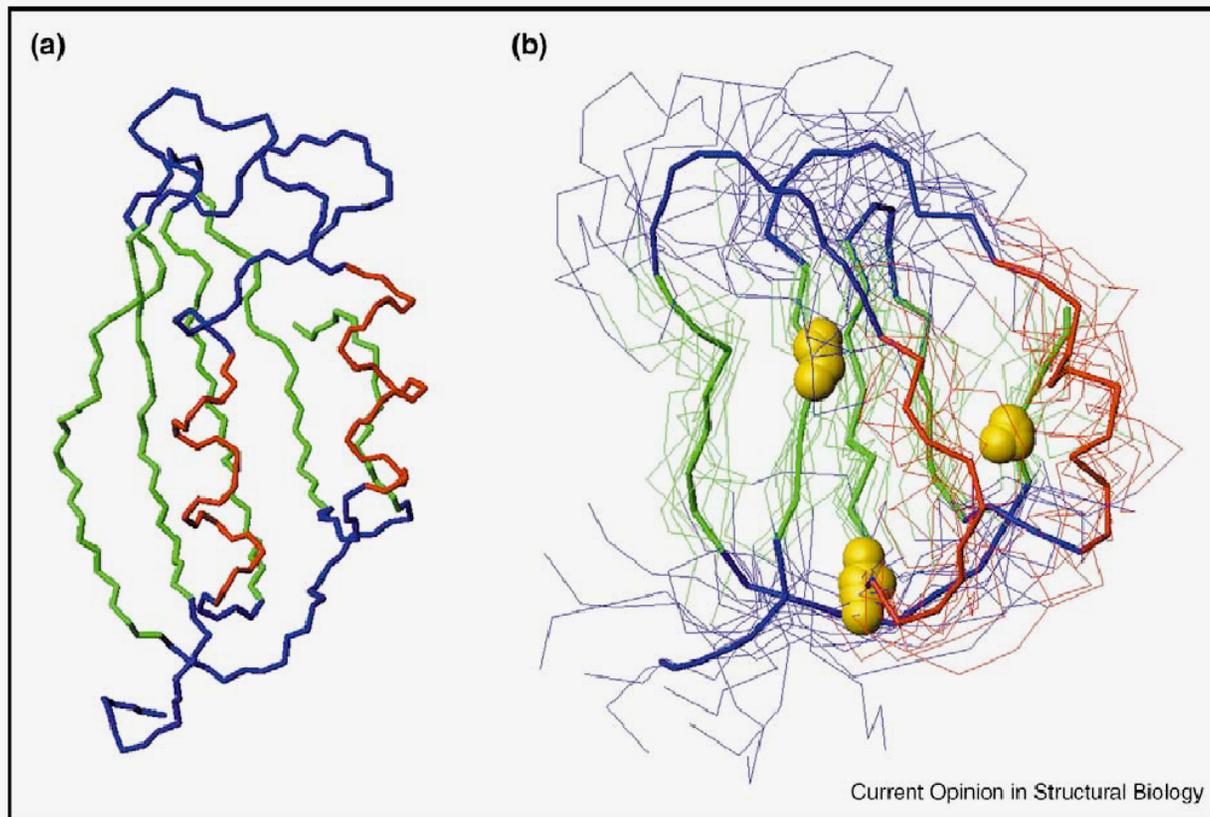
Corresponding terms for the mutant are indicated by a prime. $\Delta\Delta G_{\text{N-D}}$ and $\Delta\Delta G_{\text{TS-D}}$ are the destabilization energies of the native state and transition states of interest, respectively, caused by the mutation.

Dependence of folding mechanisms on topology

The structures of folding transition states are similar in proteins with similar native structures. The distribution of structure in the transition state ensemble can be probed by mutations at different sites in the chain; mutations in regions that make stabilizing interactions in the transition state ensemble slow the folding rate, whereas mutations in regions that are disordered in the transition state ensemble have little effect⁴. For example, in the structures of the SH3 domains of src¹⁸ (**a**) and spectrin¹⁷ (**b**), and the structurally related proteins Adah2 (ref. 37; **c**) and acyl phosphatase¹⁶ (**d**), the colours code for the effects of mutations on the folding rate. Red, large effect; magenta, moderate effect; and blue, little effect. In the two SH3 domains, the turn coloured in red at the left of the structures appears to be largely formed, and the beginning and end of the protein largely disrupted, in the transition state ensemble. (To facilitate

the comparison in **c** and **d**, the average effect of the mutations in each secondary structure element is shown.) This dependence of folding rate on topology has been quantified by comparing folding rates and the relative contact order of the native structures. The relative contact order is the average separation along the sequence of residues in physical contact in a folded protein, divided by the length of the protein. **e**, A low- and high-contact-order structure for a four-strand sheet. In **f**, black circles represent all-helical proteins, green squares sheet proteins and red diamonds proteins comprising both helix and sheet structures. The correlation between contact order and folding rate (k_f) is striking, occurring both within each structural subclass and within sets of proteins with similar overall folds (proteins structurally similar to the α/β protein acyl phosphatase¹⁶ are indicated by blue triangles).





Comparison of **(a)** the native state structure and **(b)** the most representative structures of the transition state ensemble of acylphosphatase, determined by all-atom molecular dynamics simulations [26]. Native secondary structure elements are shown in colour (the two α helices are plotted in red and the β sheet is in green). The three key residues for folding are shown as gold spheres [5*,26].

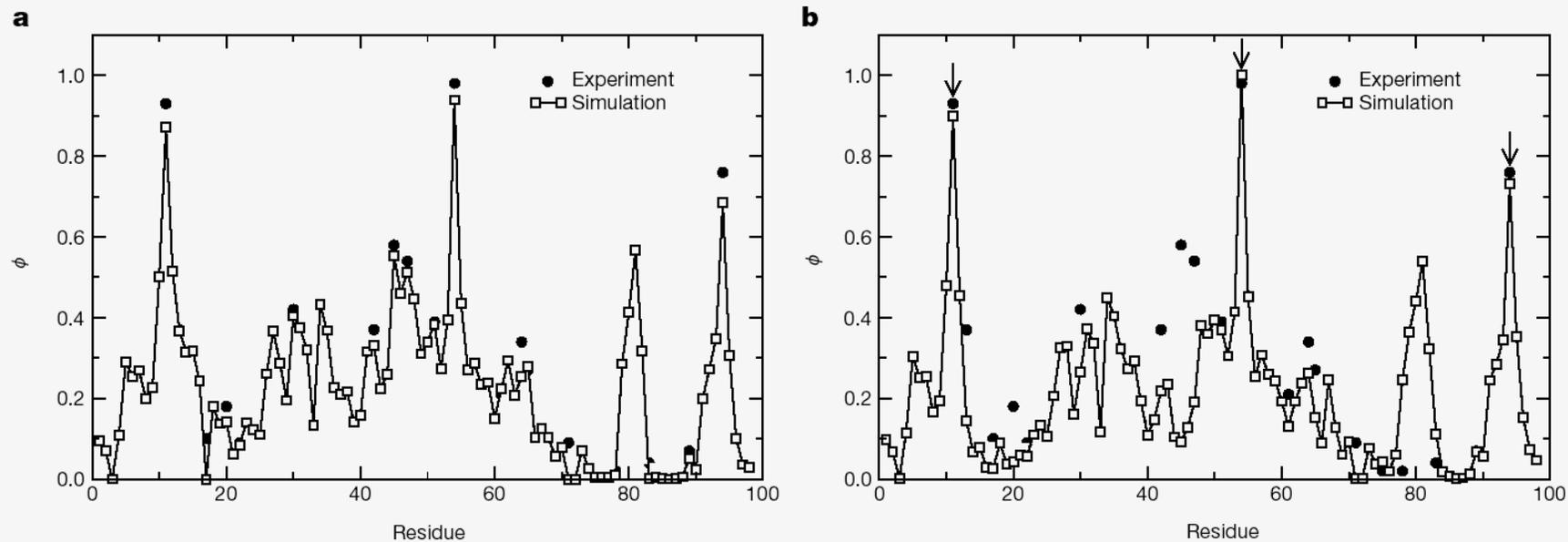


Figure 2 Comparison of ϕ_i^{calc} with ϕ_i^{exp} . **a**, When all 24 of the latter are used as restraints (filled circles); **b**, when only Y11, P54 and F94, are used as restraints (open squares, see text). Only four ϕ_i^{exp} are significantly underestimated; they are V13, T42, G45 and V47. The reason for the underestimation is that this subset of residues forms long-range

contacts mainly to each other (for example, V13 with G45) and are, therefore, almost totally unrestrained in the primitive model when Y11, P54 and F94 are used as restraints; the only common contact is between residues Y11 and V47. Except for V13, which is close to Y11, they are not part of the folding core.

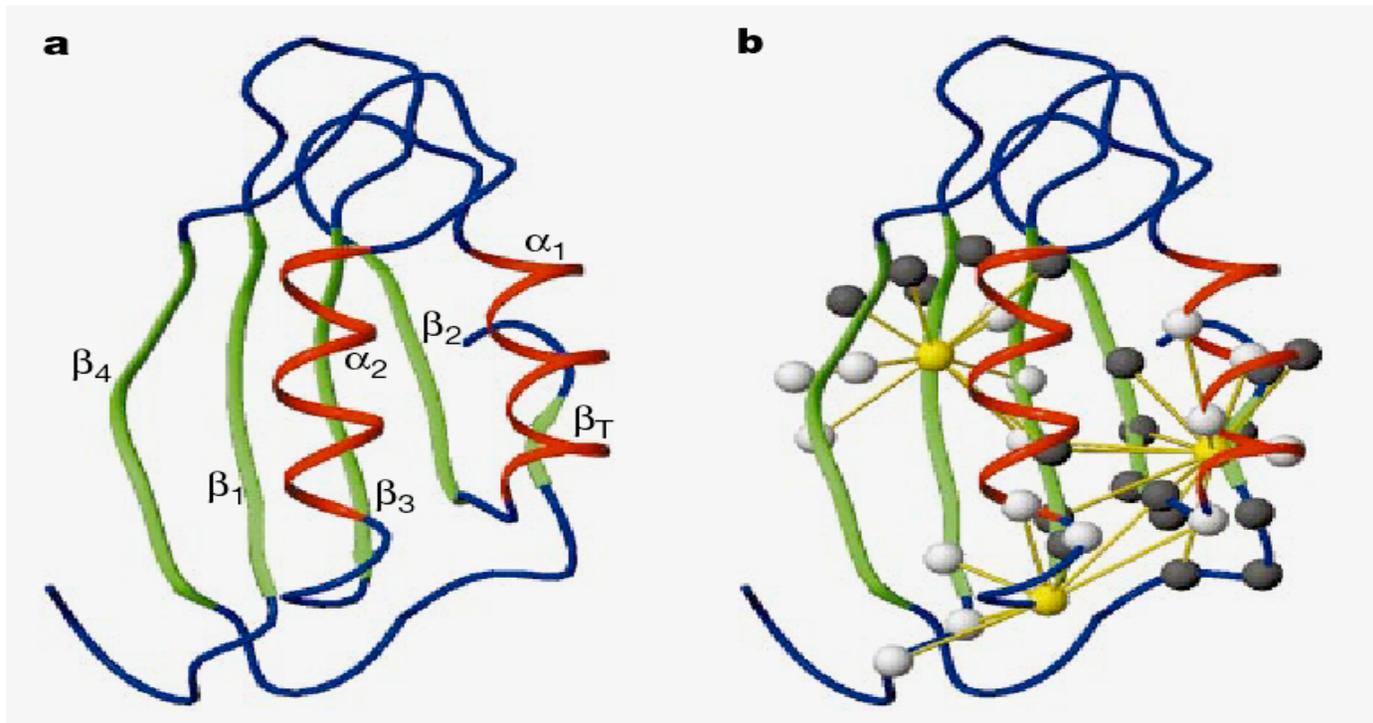


Figure 1 Native structure of acylphosphatase, AcP. **a**, The secondary structure elements are: β_1 (residues 7–13), α_1 (residues 22–33), β_2 (residues 36–42), β_3 (residues 46–53), α_2 (residues 55–56), β_4 (residues 77–85), β_T (residues 93–97); residues between these regions are parts of loops. **b**, the key residues 11, 54 and 94 found from the transition-state analysis are shown as gold spheres on the native structure (see text). They are connected by gold-coloured bonds to the residues (white and black spheres) forming a native contact with them (Y11 has long-range contacts with 47–52 and 78–81, P54 with 5–7 and 34–35, and F94 with 26–31, 36–39 and 50–52). Residues forming the transition-state core identified in the present work are shown as black spheres. They fall into two groups in spatial proximity; there is a larger group (residues 28, 35–39, 51–54 and 90–95) comprising parts of α_1 , β_2 , α_2 and β_T , and a smaller group (residues 11–13, 47 and 78, 79) comprising parts of the β_1 , β_3 and β_4 strands in the native structure.

Vendruscolo et al (2001) *Nature* **409**, 641.

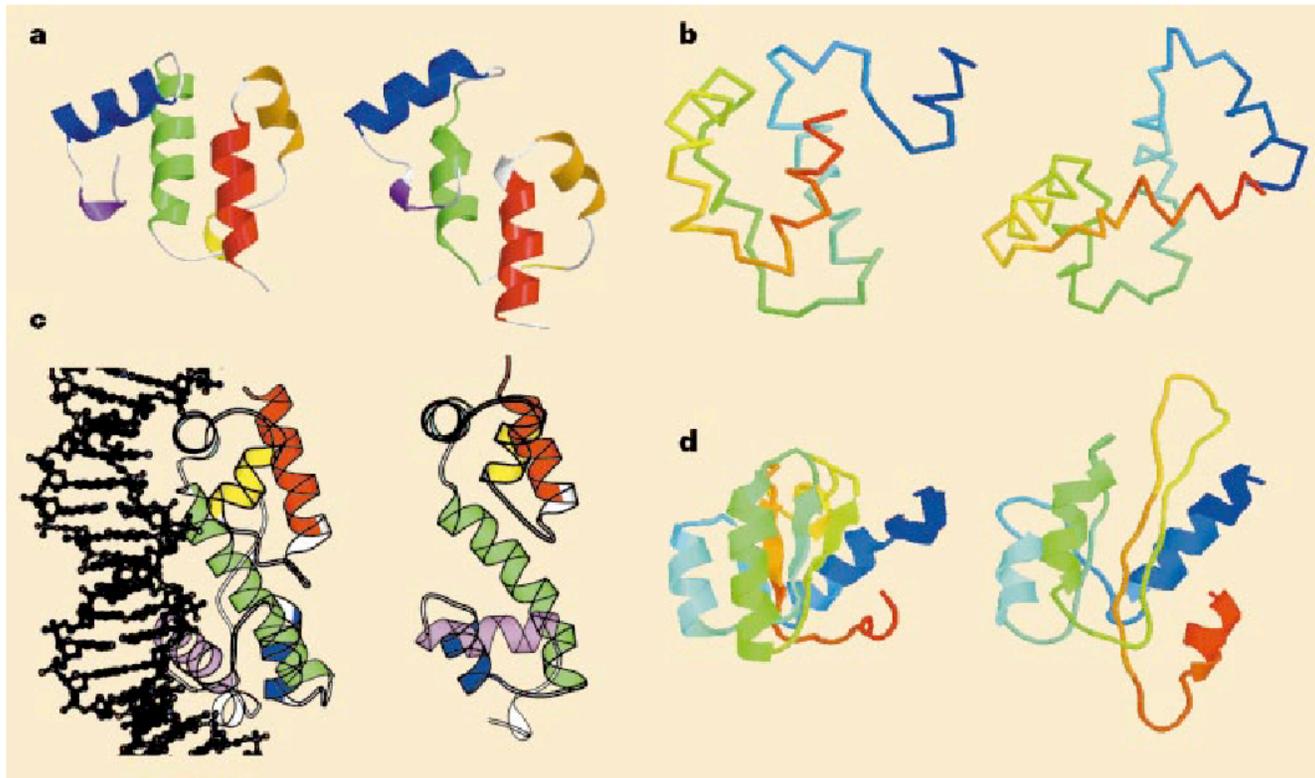
More Computational Protein Folding

**What else have the theorists
been doing lately?**

Ab initio structure predictions

Blind *ab initio* structure predictions for the CASP3 protein structure prediction experiment. For each target, the native structure is shown on the left with a good prediction on the right (predictions by Baker³⁹(**a**, **c**), Levitt⁴⁰(**b**) and Skolnick⁴¹(**d**) and colleagues; for more information see <http://predictioncentre.llnl.gov/> and *Proteins* Suppl. **3**, 1999). Segments are colour coded according to their position in the sequence (from blue (amino terminus) to red (carboxy terminus)). **a**, DNA B helicase⁴¹. This protein had a novel fold and thus could not be predicted using standard fold-recognition methods. Not shown are N- and

C-terminal helices which were positioned incorrectly in the predicted structure. **b**, Ets-1 (ref. 43). **c**, MarA⁴⁴. This prediction had potential for functional insights; the predicted two-lobed structure suggests the mechanism of DNA binding (left, X-ray structure of the protein–DNA complex). **d**, L30. A large portion of this structure was similar to a protein in the protein databank but the best *ab initio* predictions were competitive with those using fold-recognition methods. The three approaches that produced these predictions used reduced-complexity models for all or almost all of the conformational search process.



Progress in *de novo* protein structure prediction & design

A Prediction and design are inverse problems

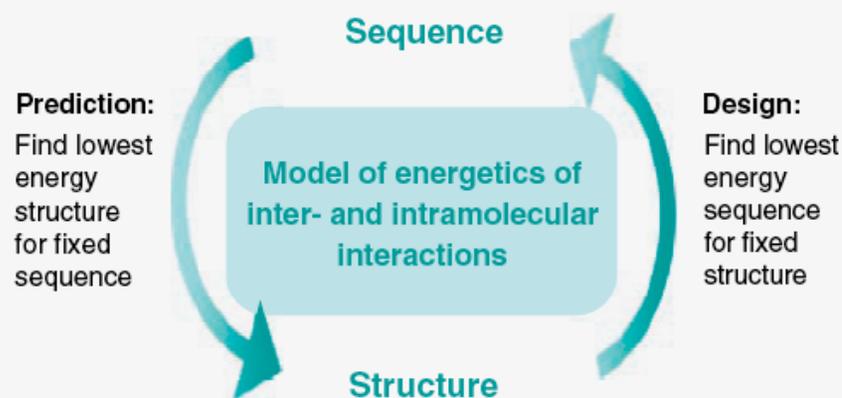
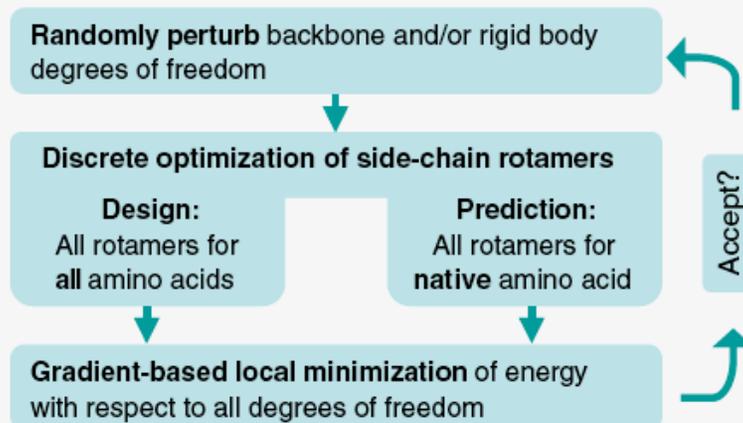


Fig. 1. Prediction and design. (A) Structure prediction and fixed backbone design are inverse problems. Completing the cycle corresponds to flexible backbone design, which requires optimization of both sequence and struc-

B Similarity of flexible backbone design and structure prediction



ture. (B) Algorithmic similarity of structure prediction, protein-protein docking, and flexible backbone design illustrated by the Monte Carlo minimization (MCM) high-resolution refinement protocol.

Present protein folding theory suggests that proteins must find the right polypeptide chain topology (“topomer”) first, then they can form 2° structure and snap into the correct 3D conformation relatively quickly.

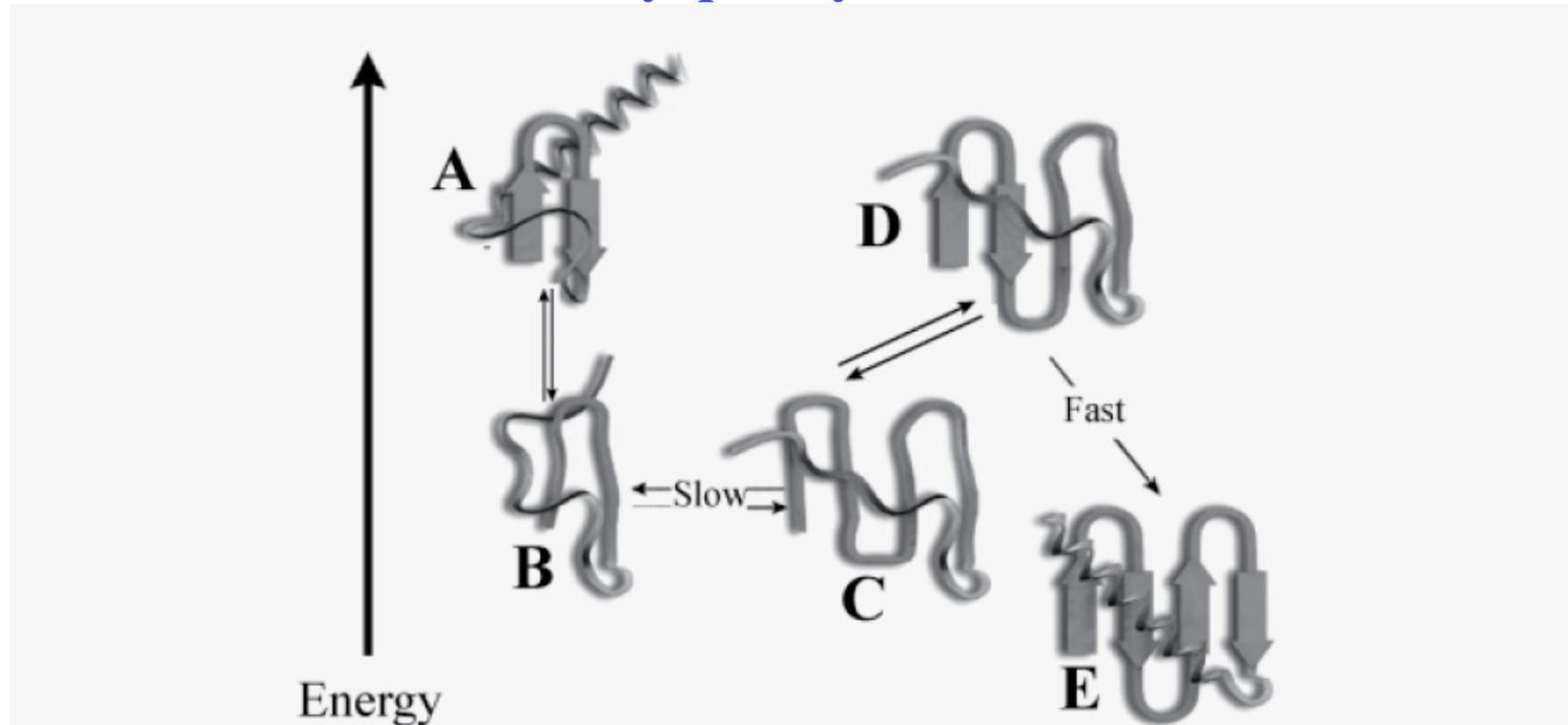
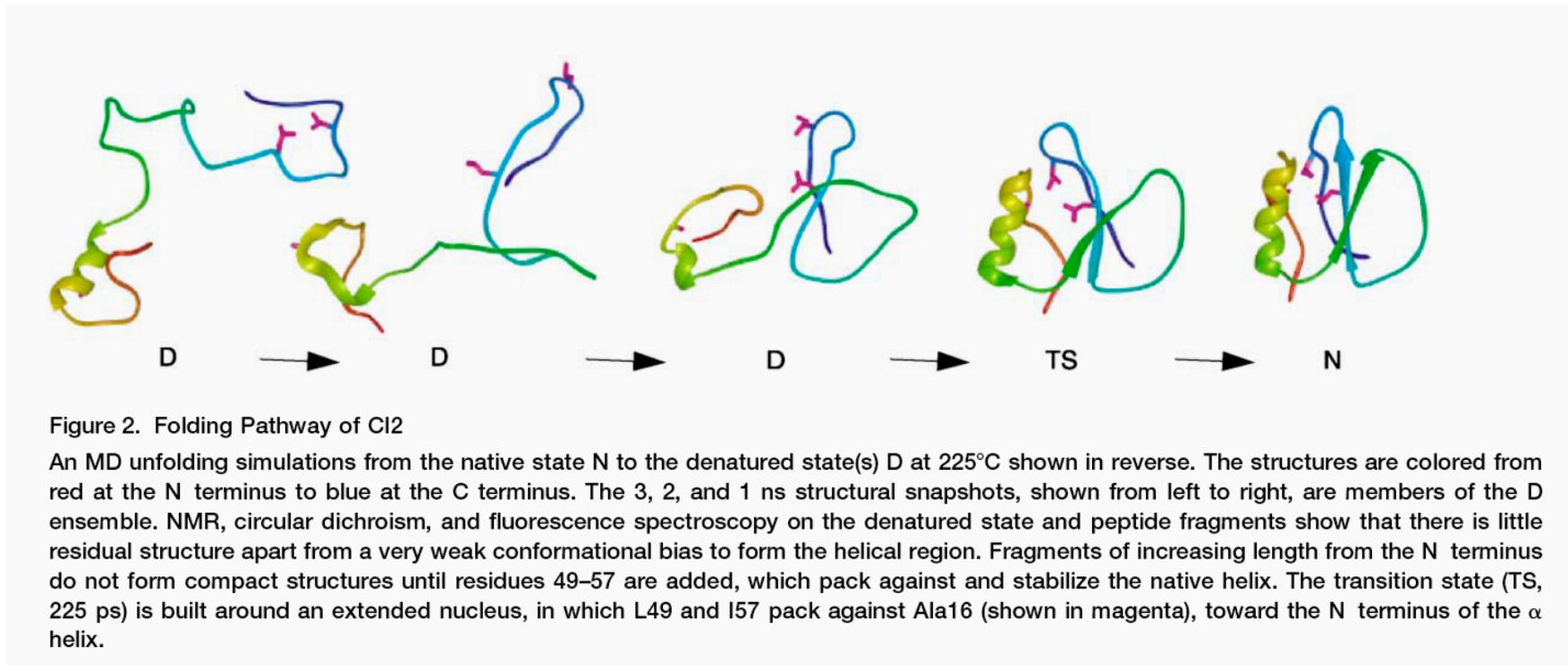


Figure 10 The essence of the topomer search model is that the rate with which an unfolded polymer diffuses between distinct topologies is much slower than the rate with which local structural elements zipper (and, critically, unzip) [reviewed in (59)].

Structures of hypothetical folding intermediates*, including a putative transition state, obtained by a molecular dynamics (MD) simulation of protein *unfolding*



***Note: structures shown are actually only representative structures of ensembles of intermediates.**

Fersht & Daggett (2002) *Cell* **108**, 573.

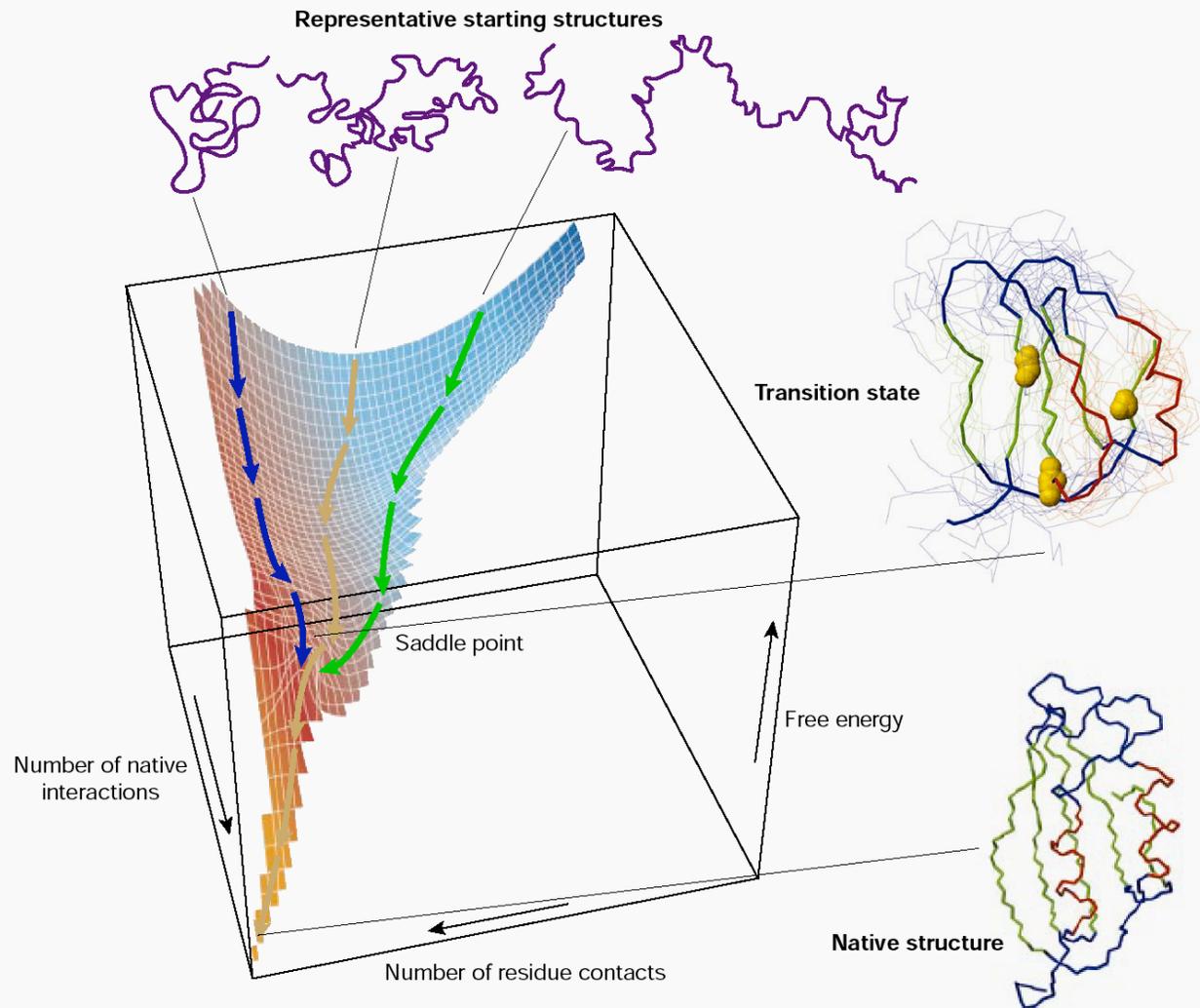


Figure 1 A schematic energy landscape for protein folding. The surface is derived from a computer simulation of the folding of a highly simplified model of a small protein. The surface 'funnels' the multitude of denatured conformations to the unique native structure. The critical region on a simple surface such as this one is the saddle point corresponding to the transition state, the barrier that all molecules must cross if they are to fold to the native state. Superimposed on this schematic surface are ensembles of structures corresponding to different stages of the folding process. The transition state ensemble was calculated by using computer simulations constrained by experimental data from mutational studies of acylphosphatase¹⁸. The yellow spheres in this ensemble represent the three 'key residues' in the structure; when these residues have formed their native-like contacts the overall topology of the native fold is established. The structure of the native state is shown at the bottom of the surface; at the top are indicated schematically some contributors to the distribution of unfolded species that represent the starting point for folding. Also indicated on the surface are highly simplified trajectories for the folding of individual molecules. Adapted from ref. 6.

Dobson (2003) *Nature* **426**, 884.

Calculating structures with their associated energies: Deep & narrow energy wells are a hallmark of near-correct structures

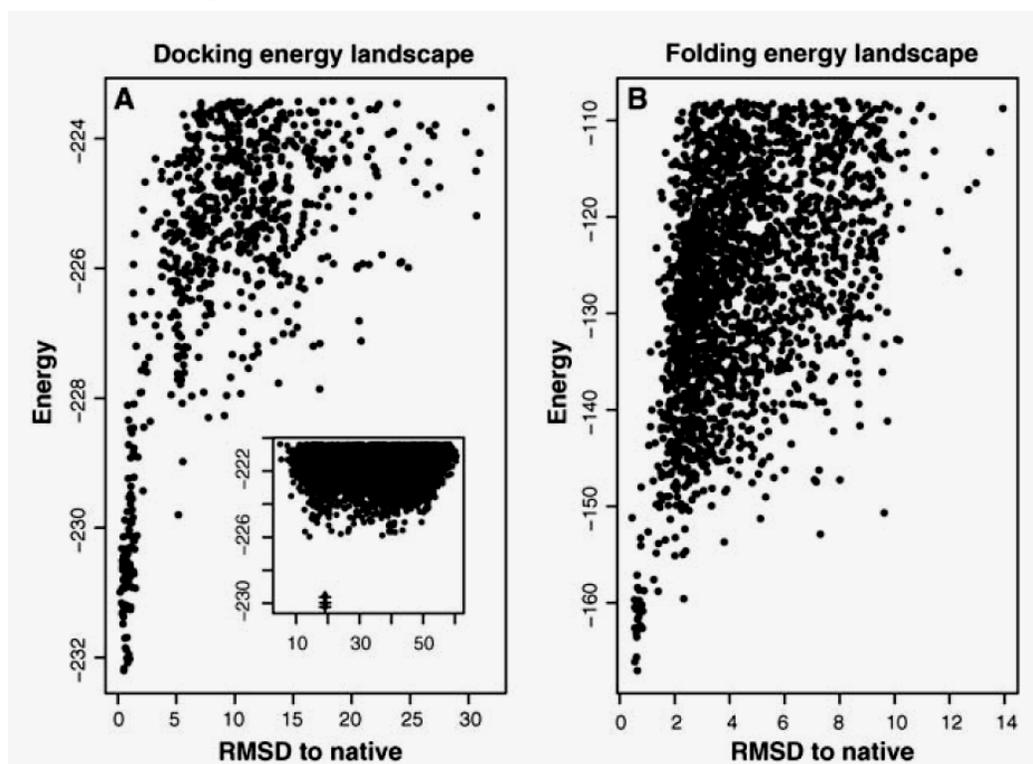
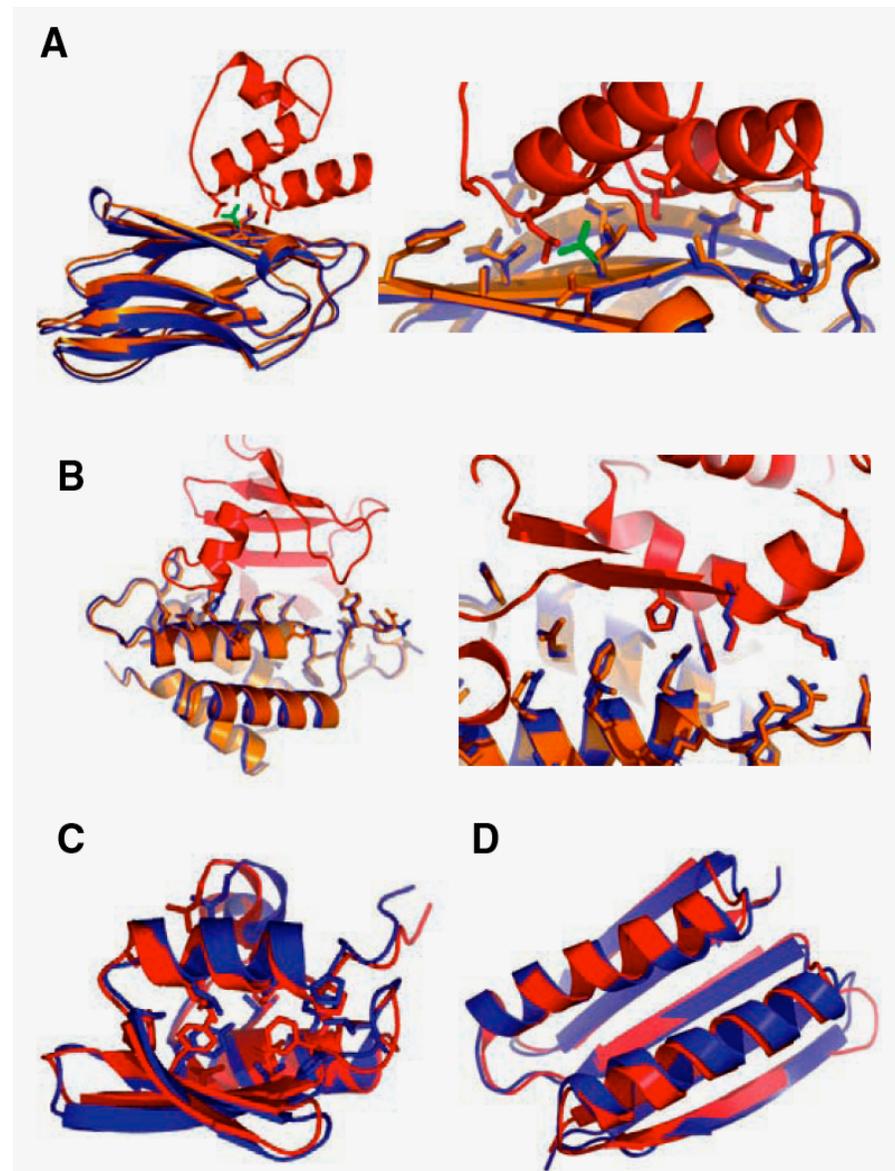


Fig. 2. Energy landscapes. Each point represents the lowest energy structure sampled in a single MCM trajectory. (A) Docking energy landscape for Capri Target 12 [cohesin-dockerin complex; PDB (Protein Data Bank) ID 1ohz (66)]. (Inset) In a large collection of trajectories starting from different random orientations that were carried out for the CAPRI experiment, a small number of structures (+) are distinguished from the background population by a significant energy gap. The x axis is the RMSD to an arbitrary reference orientation. (Main panel) Trajectories starting from these low-energy structures map out a narrow energy funnel. The x axis is the RMSD to the native structure. A deep energy funnel, as in this example, is a strong indicator that a prediction is correct. (B) Folding landscape for double-stranded RNA binding protein [PDB ID 1di2 (67)]. The backbone RMSD is to the native structure. The energy function (units are in kcal/mol) includes entropic contributions from solvation effects, but not the configurational entropy associated with protein vibrational and side-chain degrees of freedom, and hence is not the true free energy.

For small, single domain proteins, some of the predictions are getting very good!

Fig. 3. Examples of high-resolution prediction and design. (A) CAPRI Target 12 [dockerin-cohesin (66); interface residue backbone RMSD = 0.27 Å]. The lowest energy structure in Fig. 2A, main panel, is shown here. The side chain of Leu-83 (green in the free monomer) changes conformation upon binding. Side-chain conformations in red were provided; those in blue were predicted. (B) CAPRI Target 15 [ColicinD-Immunity protein D (68); interface residue backbone RMSD = 0.23 Å]. No side-chain information was provided for either partner. (C) CASP6 de novo structure prediction Target 0281 [hypothetical protein from *Thermus thermophilus* Hb8, PDB ID 1whz (69); backbone RMSD = 1.59 Å]. (D) TOP7 (RMSD = 1.2 Å) (42). (A) and (B) are adapted from figure 1 of (70). Blue: models; red and orange: x-ray structures.

Schueler-Furman et al (2005) *Science* **310**, 638.



Rosetta@home needs your help to determine the 3-dimensional shapes of proteins in research that may ultimately lead to finding cures for some major human diseases. By running the Rosetta program on your computer while you don't need it you will help us speed up and extend our research in ways we couldn't possibly attempt without your help. You will also be helping our efforts at designing new proteins to fight diseases such as HIV, Malaria, Cancer, and Alzheimer's (See our [Disease Related Research](#) for more information). Please [join us](#) in our efforts!

Site search

Join Rosetta@home

1. [Rules and policies](#)
2. [System requirements](#)
3. [Create account](#)
4. [Download BOINC](#)
5. [A welcome from David Baker](#)

About

- [Quick Guide to Rosetta and Its Graphics](#)
- [Rosetta@home Science FAQ](#)
- [Disease Related Research](#)
- [Research Overview](#)
- [News & Articles about Rosetta](#)
- [Science News](#)
- [Technical News](#)
- [BOINCWiki](#) - more about BOINC

Returning participants

- [Your account](#) - view stats, modify preferences
- [Results](#) - view your results
- [Teams](#) - create or join a team
- [Applications](#)
- [Server status](#)
- [Add-ons](#)

Community

- [Message boards](#)
- [Questions and answers](#)
- [Participant profiles](#)
- [Images](#)

Statistics

- [Top participants](#)
- [Top computers](#)
- [Top teams](#)
- [Top predictions](#)

Powered by



<http://boinc.bakerlab.org/rosetta.cgi/cgi>

User of the day



[Idle Layabout](#)

Server Status as of 6 Mar 2006 20:31:55 UTC

[Scheduler running] Queued: 24,860
 In progress: 164,278
 Successes last 24h: 86,951
 Users [L](#) (last day [L](#)) : 43,199 (+164)
 Hosts [L](#) (last day [L](#)) : 91,678 (+330)
 Credits last 24h [L](#) : 1,899,309
 Total credits [L](#) : 193,996,138
 TeraFLOPS estimate: 18.993

News

March 2, 2006
 The default cpu run time is now set at 2 hours instead of 8. This change will affect new work units only.

February 18, 2006
Rosetta application update! Graphics are now available for Mac OSX platforms.

February 17, 2006
 Work is flowing again. Today we upgraded our database server. Unfortunately, we will be delaying the application update for a day or two to work out a few minor issues. See [Technical News](#) for details.

February 17, 2006
Outage Notice: The project will be down starting today at 3pm PST for maintenance. The server should be back online later in the evening.

We will be updating the rosetta application today. There are a number of new features:

- Work units will have a default cpu run time of 8 hours, and users will have the option to change the cpu run time as a project specific preference. The length of work units will no longer depend on the number of predicted structures. This option was added to allow participants to reduce bandwidth usage per work unit and maintain consistent run times.
- Users will also have the option to change the frame rate and cpu use for graphics.
- A new graphics version will be available for Mac OSX users.

February 15, 2006
Volunteers needed!!

We are seeking volunteers for our new alpha test project, [RALPH](http://ralph.bakerlab.org) (<http://ralph.bakerlab.org>).

There are a number of recent improvements to the rosetta application but we need volunteers to speed up the process of testing to make the updated application available for production as soon as possible.

If you are interested in helping to improve Rosetta@home and can spare a few extra cycles for testing, please [join RALPH@home](#).

[...more](#)

News is available as an [RSS feed](#).

You, too, can do
 theoretical protein folding
 with David Baker!!

See <http://boinc.bakerlab.org/rosetta>
 But make sure your computer doesn't overheat!

Folding@home

<http://folding.stanford.edu/>

(Vijay Pande lab)