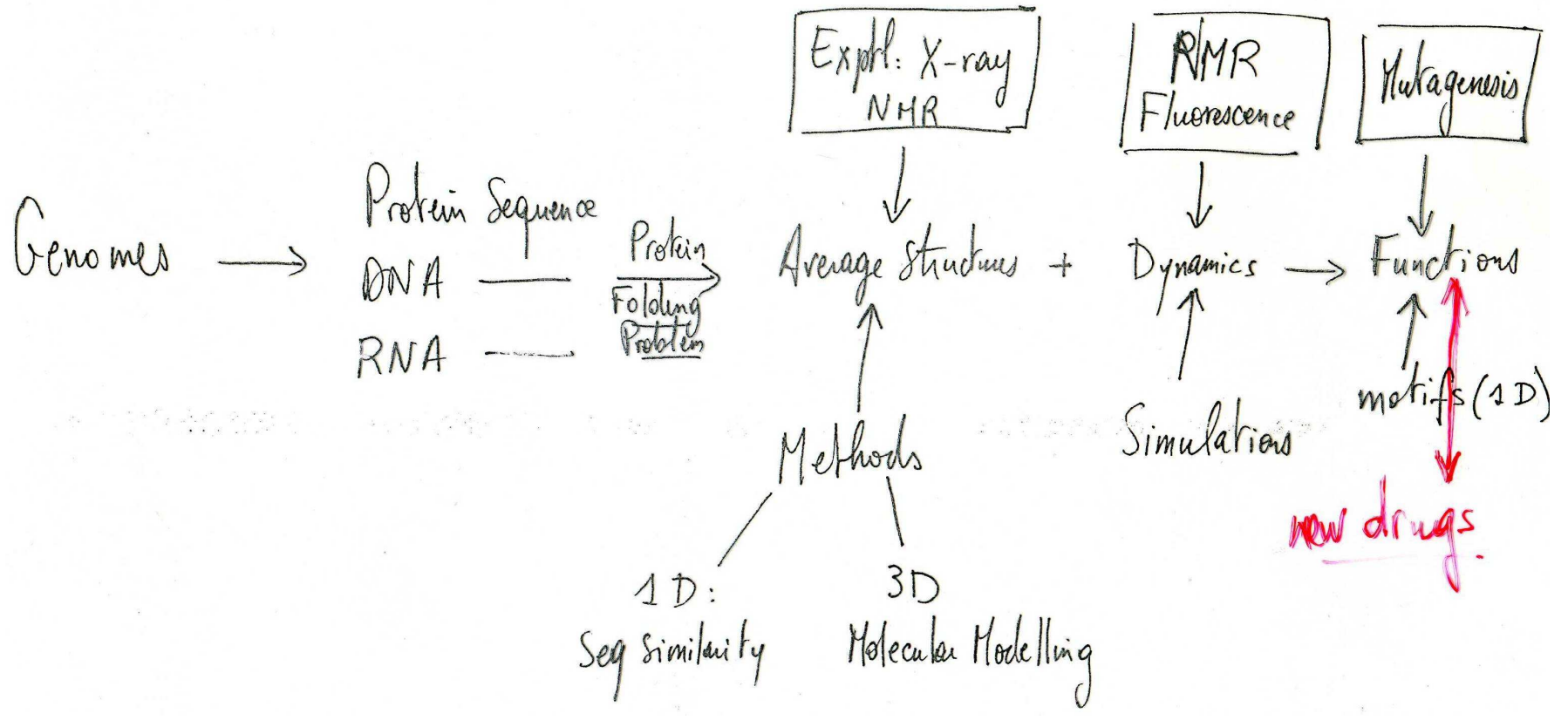


M1 Spécialité Bioinformatique Lecture 1A

P. Derreumaux

**The protein folding problem:
Generalities**

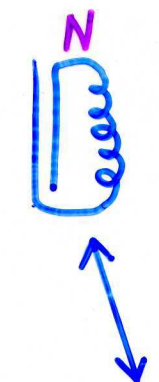
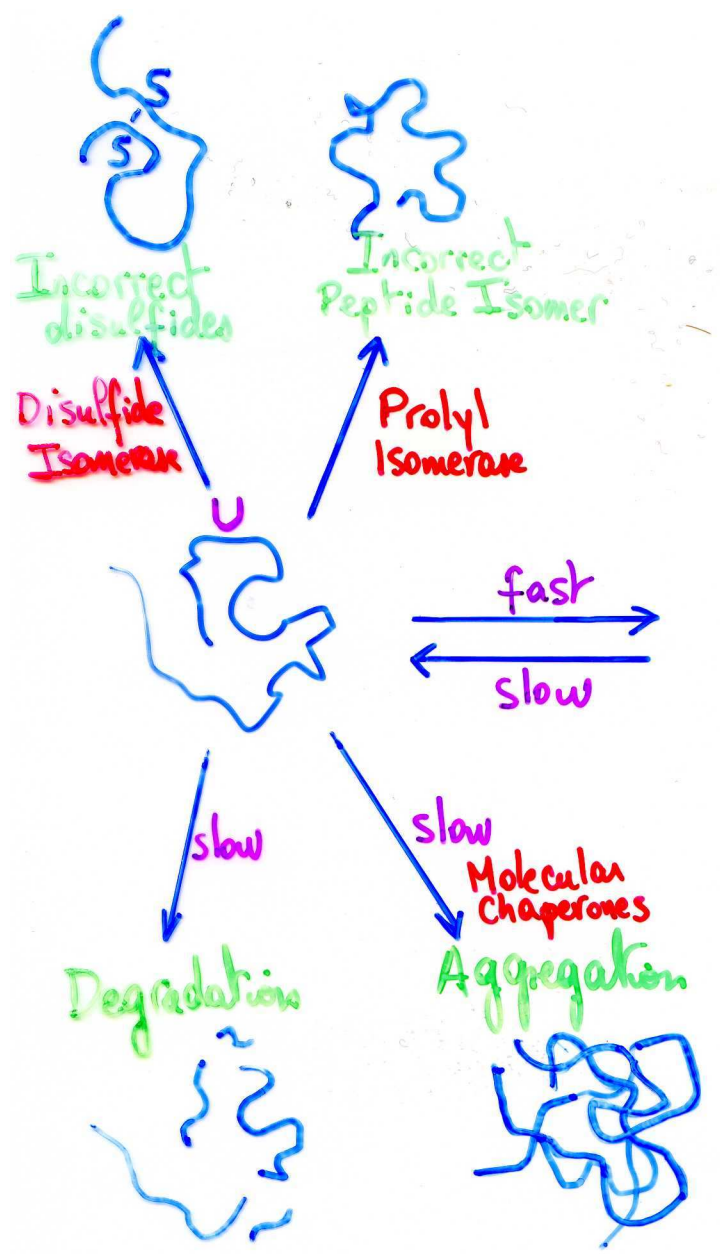


The Protein Folding problem

2 aspects

- structures from sequences.

- sequence events from disordered to native structures (characterization of TSE) →
sequence elements coding for protein topology



- stable (pH, T, mutations)
- bioactive: (flexibility \rightleftharpoons dysfunction)
- variety of functions

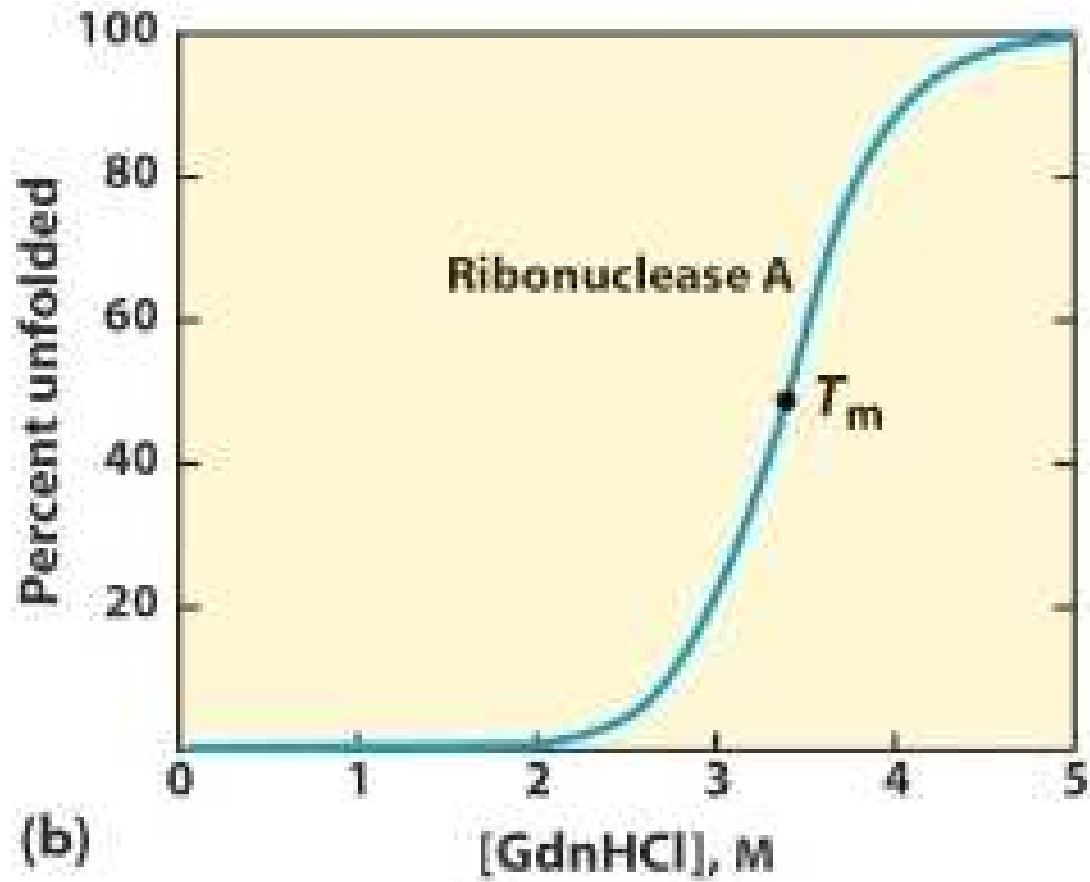
Others (fibre,...)

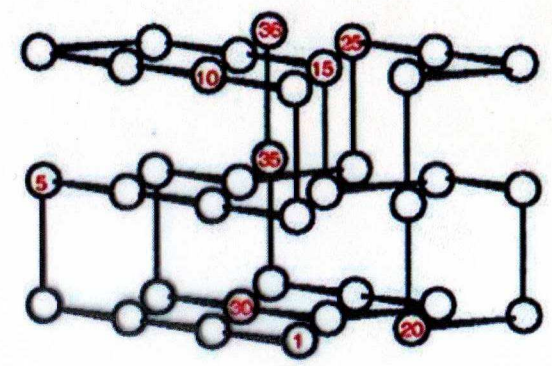
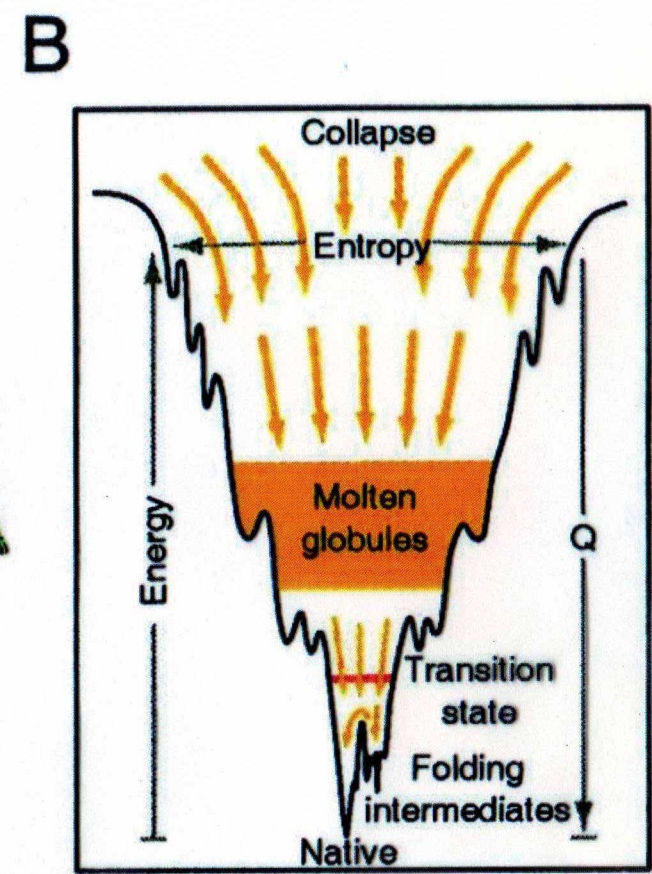
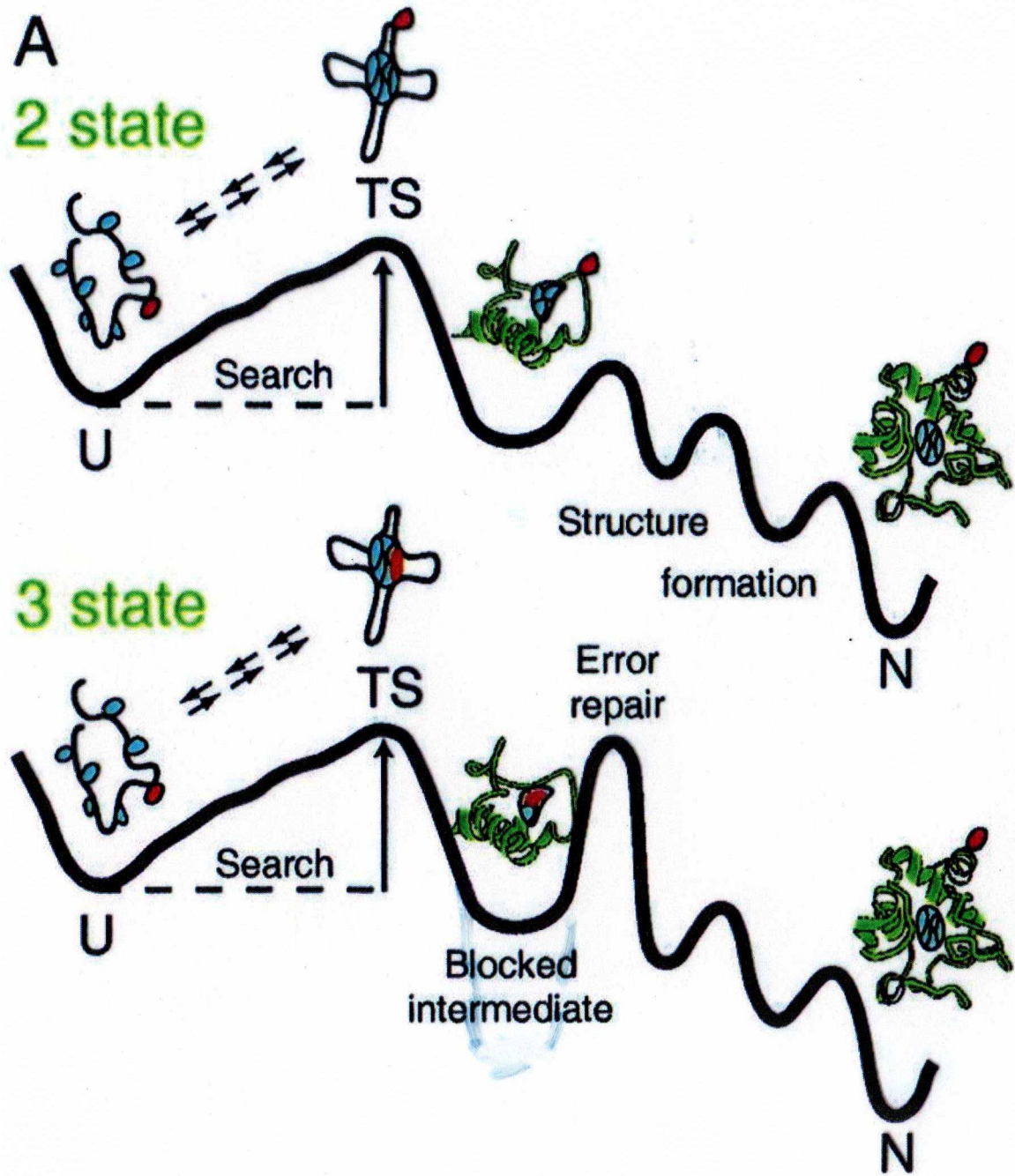
-Folding under chaperon control

-Thermodynamic hypothesis of folding (Anfinsen)

- Amyloid Fibril Formation

Monomeric Protein Folding: a simple experimental view





Monomeric Protein Folding: a simple view

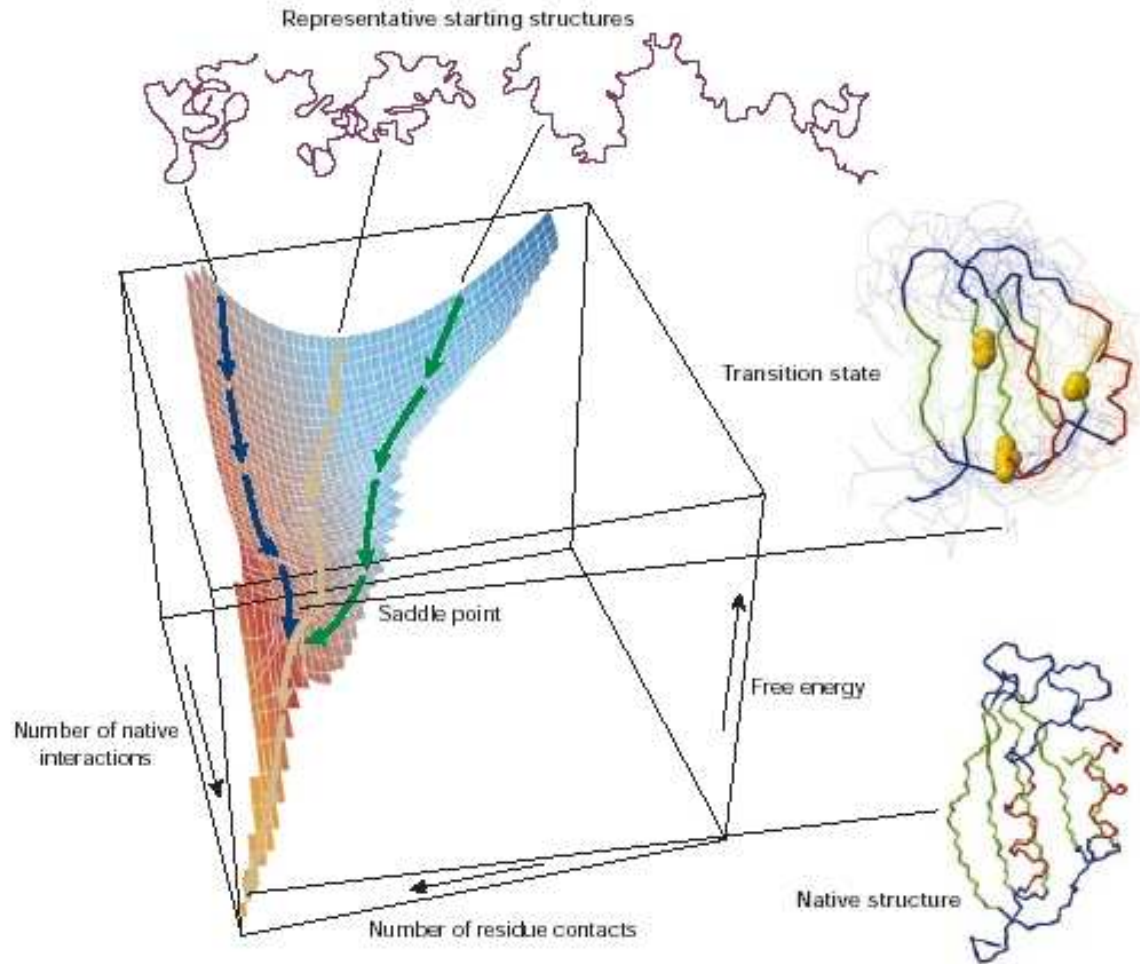
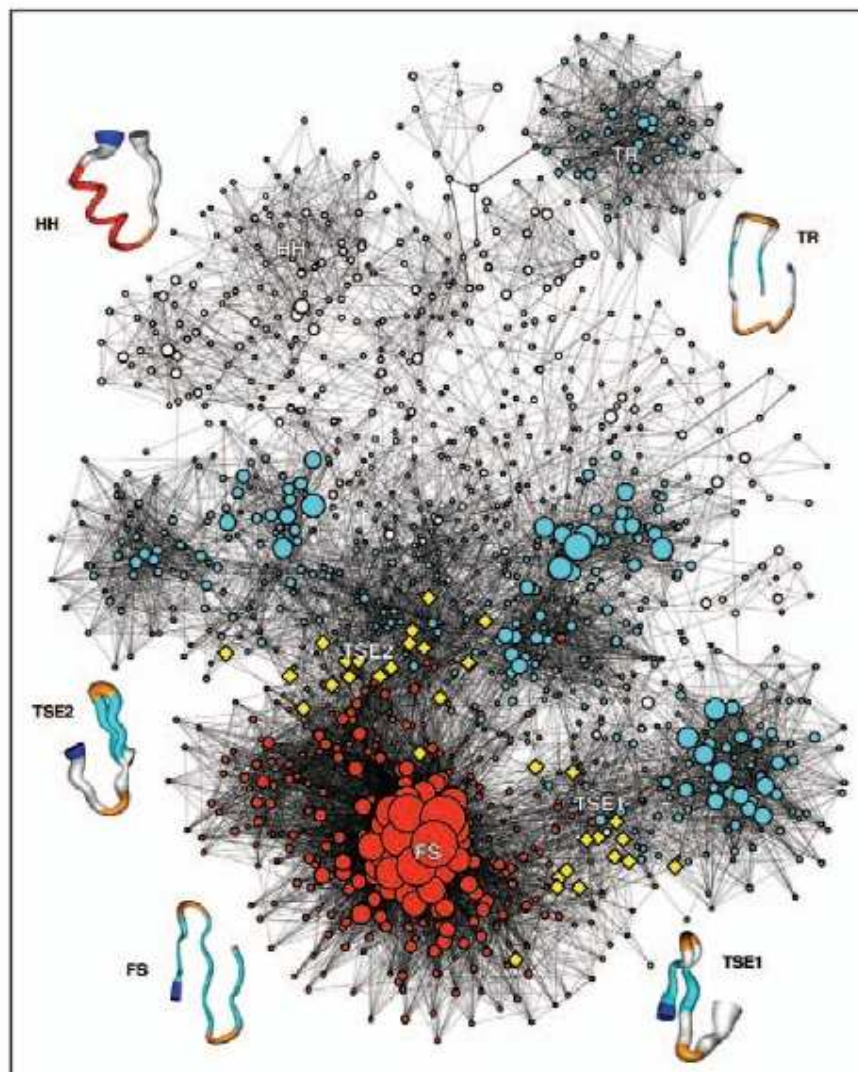


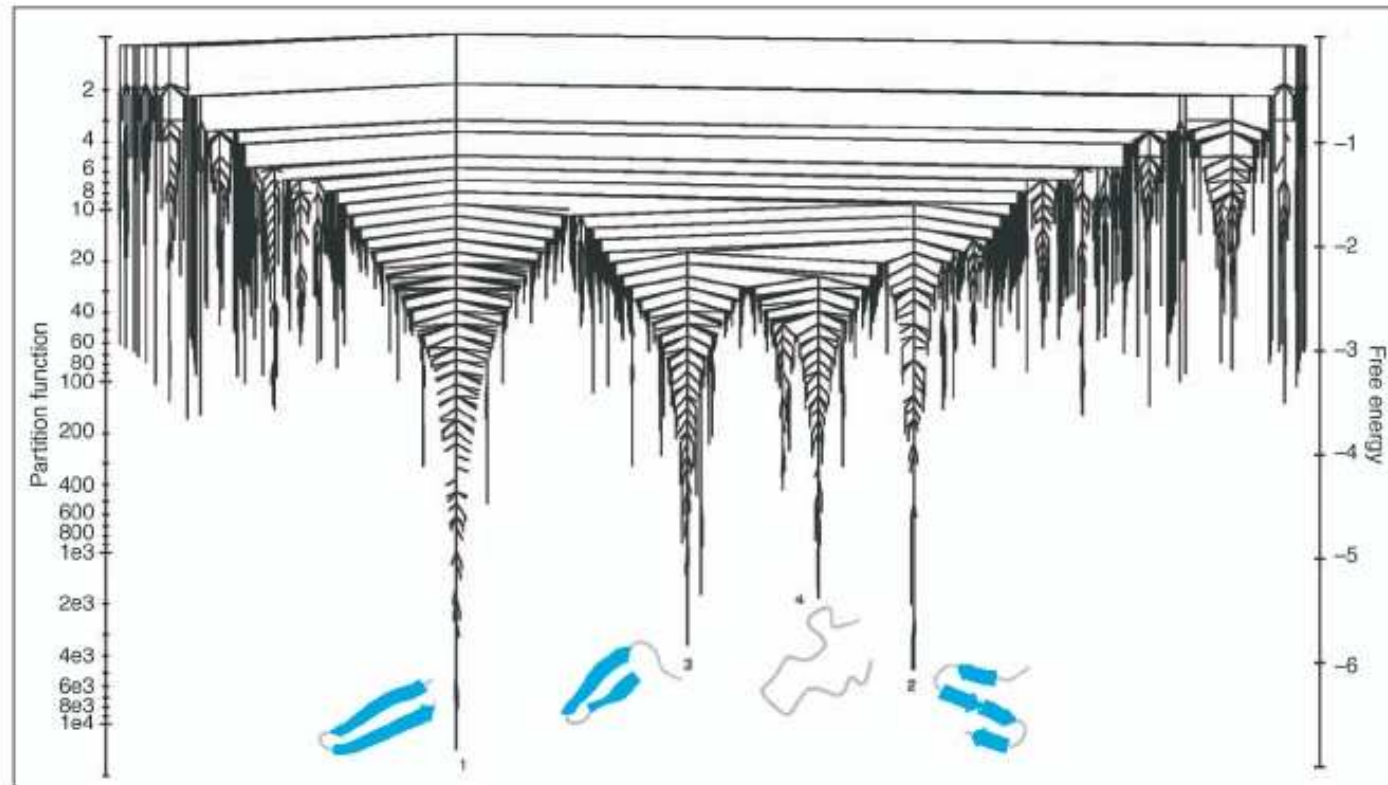
Figure 1 A schematic energy landscape for protein folding. The surface is derived from a computer simulation of the folding of a highly simplified model of a small protein. The surface 'funnels' the multitude of denatured conformations to the unique native structure. The critical region on a simple surface 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.

The funnel energy surface is oversimplified



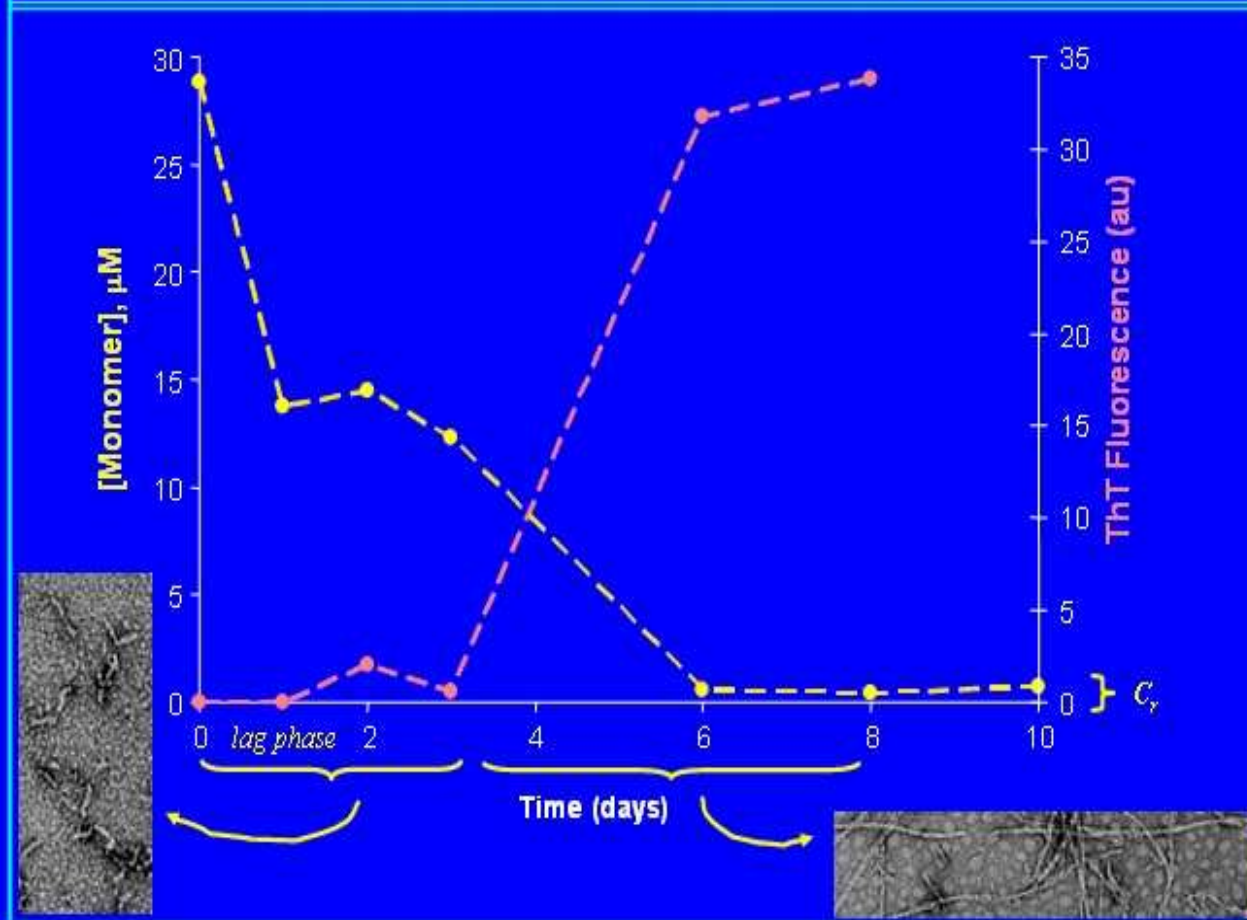
The funnel energy surface is oversimplified

Figure 3



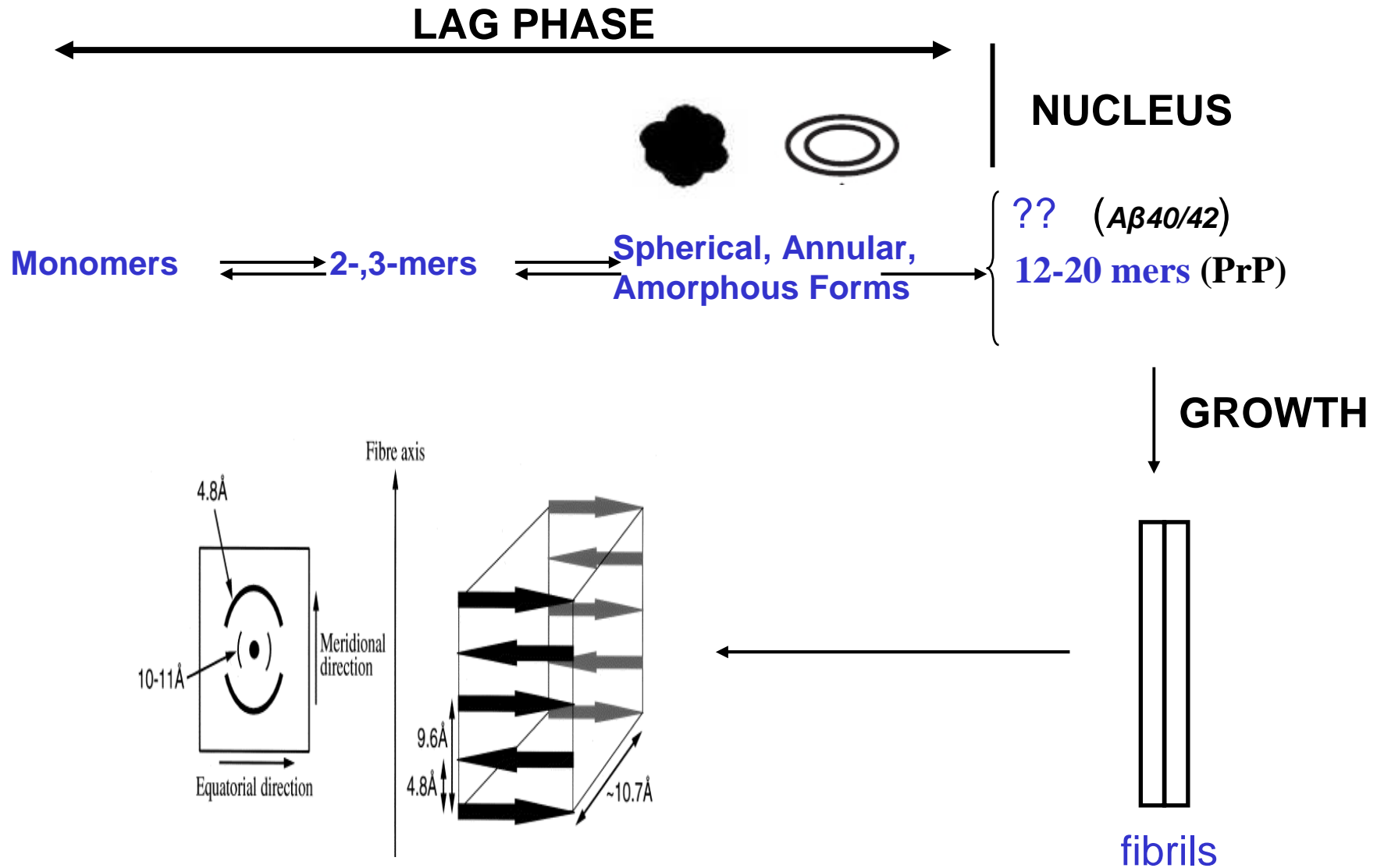
Transition disconnectivity graph of a β -hairpin (the C-terminal segment from the B1 domain of protein G). A total of 4 μ s implicit solvent molecular dynamics simulations at 360 K were sampled to obtain a sufficient number of folding-unfolding events [29]. Representative structures of the deepest free energy minima are shown and labeled 1–4. The left vertical axis shows the partition function of the minima and barriers. The right vertical axis shows the free energy of the minima and barriers. Reproduced with permission from [29].

A β Amyloid Fibril Formation



Amyloid fibril formation is characterized by a polymerization-nucleation process

Details of the polymerization-nucleation process



Overall, aggregation very sensitive to sequence, pH, concentration, anions

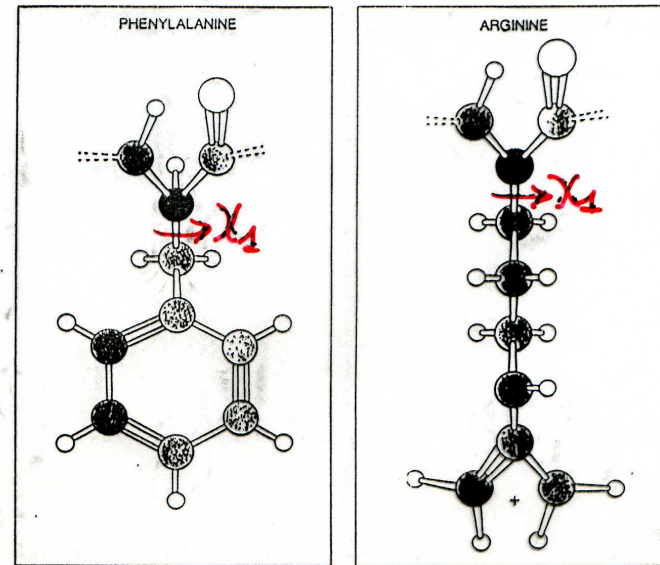
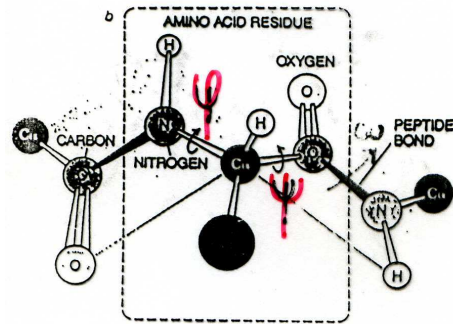
Prediction 3D structure from sequence

- Genomic Programs and the number of sequences**
- Structures are more conserved than sequences**
- Experimental costs for structure determination**

A small number of sequences associated with functions

Alphabet 20 AA.

$$N = 200 \text{ AA} \rightarrow 20^{200} \text{ seq} \gg \gg 10^{23}$$



DIFFERENCES in the shape, size and polarity of amino acids derive from differences in their side chains. In phenylalanine, for example, the side chain is nonpolar and cyclic, whereas the side chain of arginine is both strongly polar and linear.

A small number of 3D architectures
for water-soluble proteins, ≈ 450 known vs.
 $10^3 - 5 \cdot 10^3$ estimated

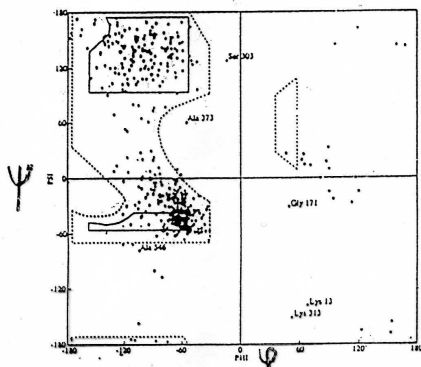
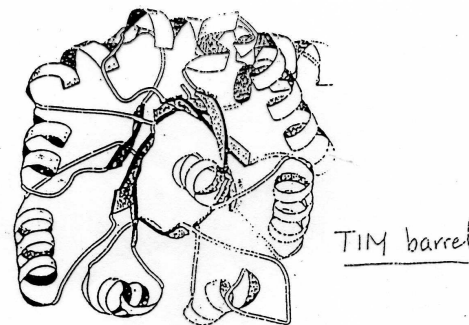
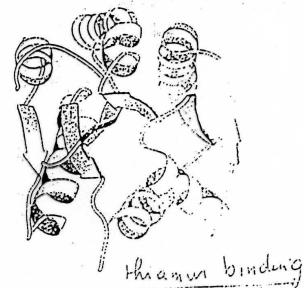
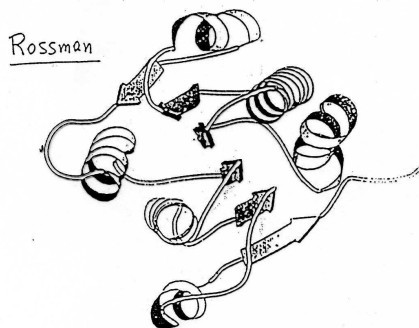
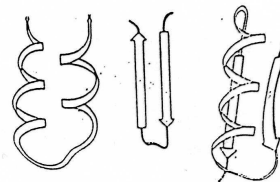


Fig. 1. Ramachandran plot of the refined hTIM structure. Open circles denote glycine residues, and the prolyl residues are indicated in triangles. Residues lying outside the allowed limits are marked (see text for details).



domain residues 213 to 326; Ammon *et al.*, 1900), (g) immunoglobulin, 3hr (chain B, domain residues 32 to 131; De Vos *et al.*, 1992); (h) UB, 1sha (chain A; Waksman *et al.*, 1992); (i) jelly roll, 2stv (Fridborg *et al.*, 1965; (j) plaitfold, 1aps (Saudek *et al.*, 1989).

Table 1. Summary of average % of secondary structure residues which are in supersecondary structural units and motifs in all the categories of domain structure

Fold	n	% Secondary structure residues in:		
		$\beta\beta$, $\alpha\alpha$, $\beta\alpha\beta$	$+(\beta\alpha\beta\beta)$	(β_4)
		88.0	88.0	88.0
Globin	4	90.3	90.3	90.3
UpDown	6	83.4	83.4	90.2
Trefoil	4	82.4	83.4	83.4
TIM barrel	13	76.7	76.7	81.4
OB folds	5	68.1	70.3	70.3
Doubly wound	38	66.7	66.7	85.6
Ig	17	55.3	55.3	61.1
UB roll	6	47.3	47.3	67.9
Jelly roll	6	47.3	47.3	64.0
Plaitfold	23	37.8	62.8	74.8
All superfolds	122	64.7	70.5	74.8
Non-superfolds	516	61.7	63.4	66.3

The *n* value is the number of domains per fold. The $\beta\beta$, $\alpha\alpha$, $\beta\alpha\beta$ column is the mean value for % of secondary structure residues in supersecondary structural units. The $+(\beta\alpha\beta\beta)$ column is the mean value for % of secondary structure residues in supersecondary structural units and in the $\beta\alpha$ -Greek key motif. The $+\beta_4$ column is the mean value for % of secondary structure residues in supersecondary structural units and in the $\beta\alpha$ -Greek key and in the β_4 Greek key motifs. This is the order in which the motif additions were made to the data.

PDB 1999 \approx 50% secondary structures

The Protein Folding problem and A limited number of folds. ($\sim 800 < N_{\text{fold}} < 5000$)

basic reasons:

obvious — structural stability

$$E_{\text{Nat}} \rightarrow P(E_{\text{Nat}}) = \frac{e^{-\frac{E_{\text{Nat}}}{kT}}}{\sum_i e^{-\frac{E_i}{kT}}} \gg P(E_i)$$

— functional property
(n.b. Identical function for different folds)

not obvious

— kinetic property
(‘fast folding seq’)
(‘folding nucleus’)

— AA mutation variability
and its correlation with the number
of folds

— AA mutation and its effect
on structural stability.

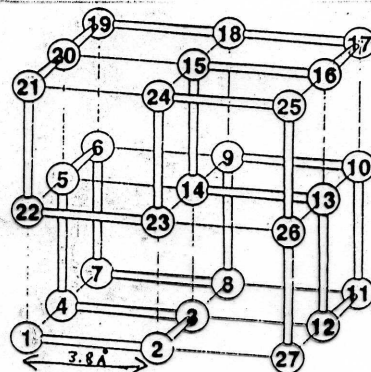
IDP: Intrinsically Disordered Proteins

- 30% of eukaryotic genome encoded proteins**
- disordered over the full or part of the sequence**
- fold upon binding on one or multiple partners**

Thermodynamic, kinetic and dynamic aspects

Mostly lattice Models.
(C_α)

Off-lattice C_α Models (Nymeyer, PNAS, 95, 1998)



Thermodynamic Properties

full enumeration N_c^* , GEM known

Canonical Ensemble

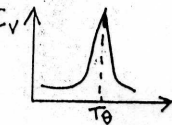
$$Q = \sum_c N_c e^{-\frac{E_c}{k_B T}}$$

$$\hookrightarrow A = -k_B T \ln Q = E - TS$$

$$\hookrightarrow P_{Nat} = \frac{e^{-\frac{E_{Nat}}{k_B T}}}{Q} \rightarrow T_f$$

$$\hookrightarrow C_v \rightarrow T_0$$

$$\hookrightarrow S \rightarrow T_g \quad \left(\frac{\partial S}{\partial T} = 0\right)$$



$$N_c = N \text{ (30-mers 2D, 15-mers 3D)} \\ = N_{compact} \text{ (27-mers, 3D)}$$

Kinetic Properties

run MC Metropolis

- $\min(1, e^{-\frac{\Delta E}{k_B T}})$
- moves

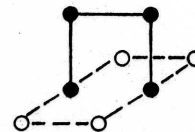


Figure 3. Cranksaft motion. The original is indicated by solid lines and filled circles. conformations are indicated by dashed line

- $\chi(t) = f(t)$ in function of sequence, T, PEF

What have we learned?

- Small Peptide model (3x3x3, 4x4x4)

$$Z = \frac{E_{NAT} - E_U}{\sigma_U}$$

\Rightarrow fast folding.
(Sali, Nature, 369, 1994)

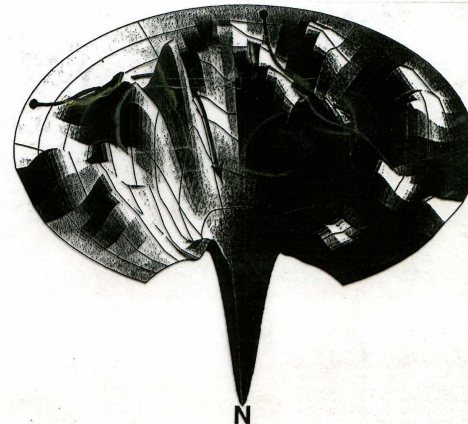
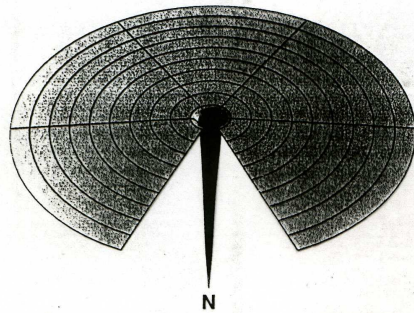
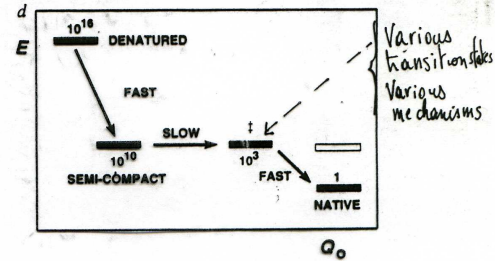
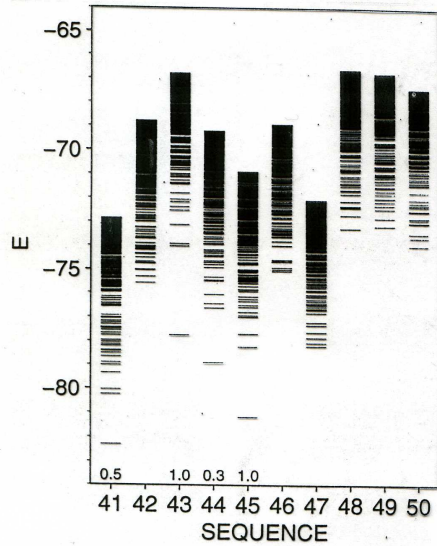


Fig. 1 The Levinthal 'golf-course' landscape. N is the native conformation. The chain searches for N randomly, that is, on a level playing field of energies.

(Dill, Nat. Struct. Biol., 4, 1997)

Coût d'une structure 3D

	Coût moyen	Change de succès
Protéines solubles bactériennes	\$140000	35%
Protéines humaines solubles (kinases, protéases, ...)	\$450000	35%
Protéines membranaires bactériennes	\$1,5 million	10%
Protéines membranaires humaines	\$2,5 million	10%

R.C. Stevens, Drug Discovery, 2003

Coût n'incluant pas les développements technologiques ni l'amortissement des équipements lourds

M1 Spécialité Bioinformatique Lecture 1B

P. Derreumaux

Predicting Protein Structures from AA:

Two perspectives

- 1. If you were a theoretical chemist or physicist**
- 2. If you were a pure bioinformatician**

From Quantum Mechanics.

$$H\psi = E\psi$$

to Molecular Mechanics
and Molecular Dynamics

↓
syst: ensemble of
springs

↓

$$\vec{F} = m \vec{a}$$

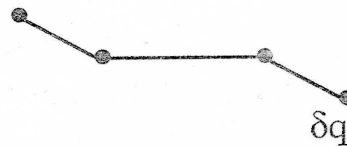
Niveau d'approximation en Mécanique classique

Born-Oppenheimer : les mouvements des électrons sont découplés des mouvements des noyaux.

Les électrons sont représentés par leurs effets : charges partielles, paramètres (distance de référence entre 2 atomes liés de façon covalente, ...).

Seuls, les mouvements des noyaux sont considérés et sont traités dans un modèle de *mécanique classique*.

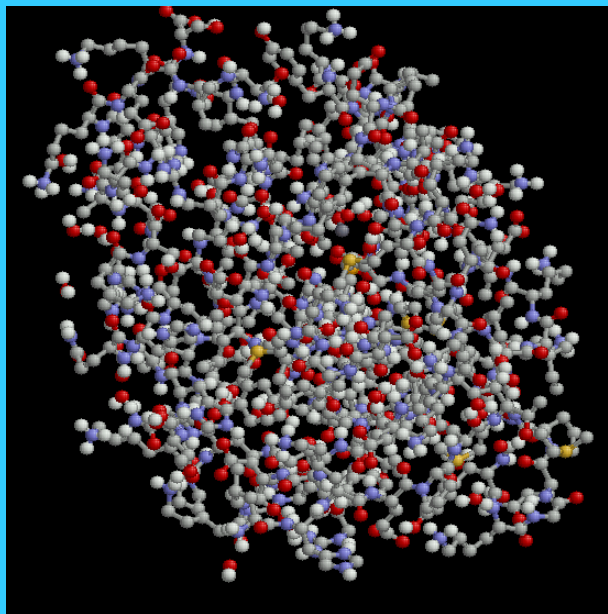
- champs de forces **semi-empiriques**
- grands systèmes, typiquement < 10000 atomes
- solutions numériques
- dynamique, limitée à ~ 1 ns (10^{-9} s)
- représentation du solvant et des contre-ions
- énergie conformationnelle
- pas de réactions chimiques, polarisation, transfert d'électrons



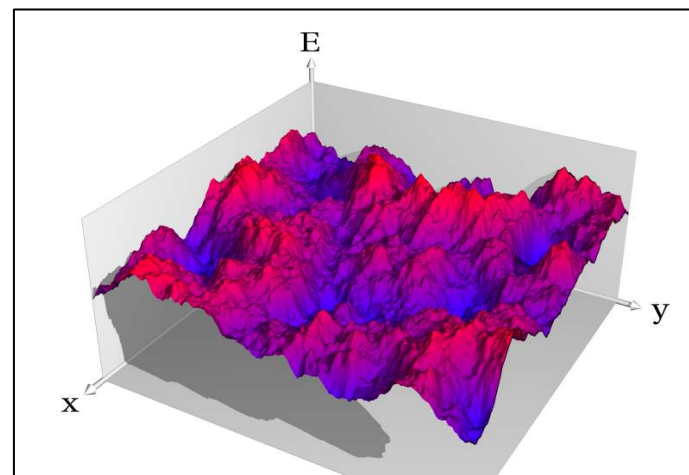
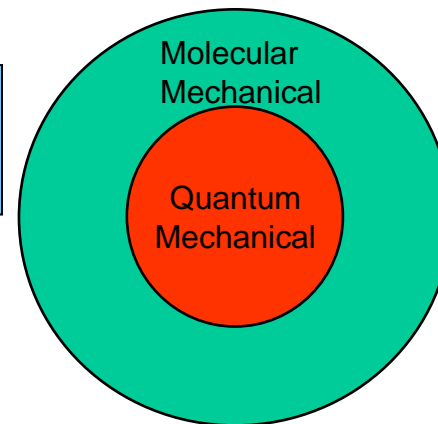
Capacité de transfert : des paramètres calculés et testés sur de petites molécules peuvent être utilisés pour des systèmes plus importants.

Computer Simulation - Basic Principles

Model System

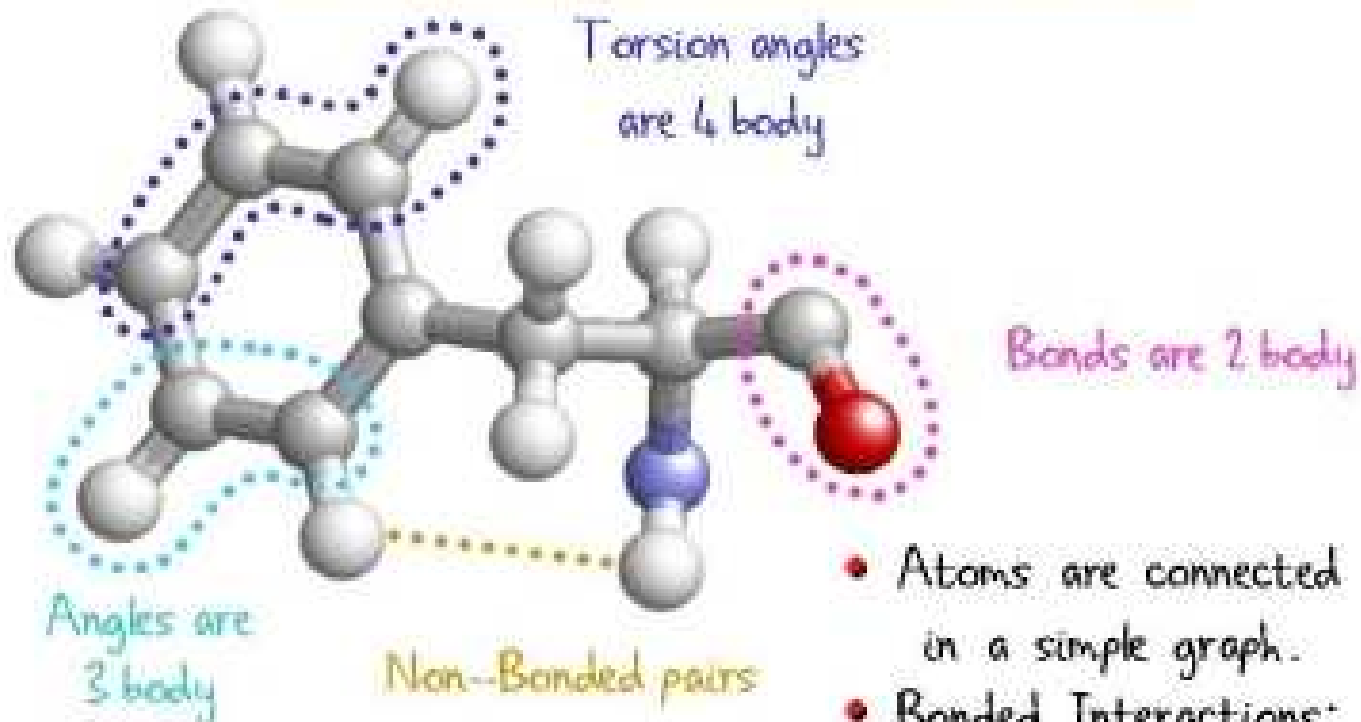


Coarse-grained
Force field



Techniques to explore the
energy landscape
(structures, thermodynamics)

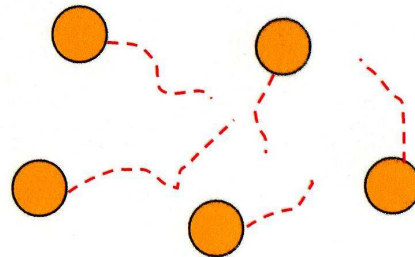
WHAT IS A MOLECULE?



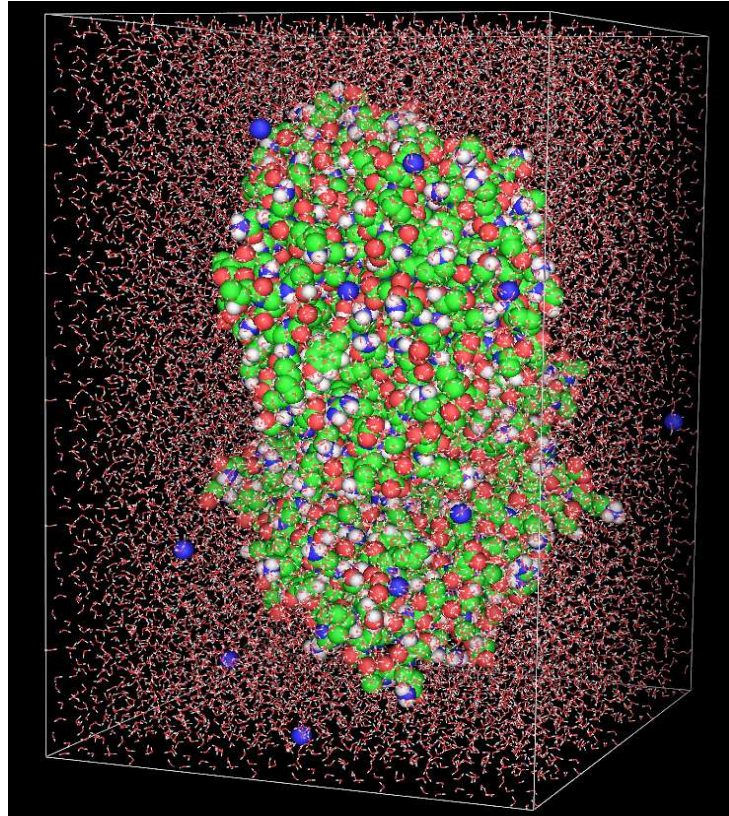
- Atoms are connected in a simple graph.
- Bonded Interactions: 2, 3, 4 body.
- Non-Bonded Interactions.

What is a molecular dynamics simulation?

- Simulation that shows how the atoms in the system move with time
- Typically on the nanosecond timescale
- Atoms are treated like hard balls, and their motions are described by Newton's laws.

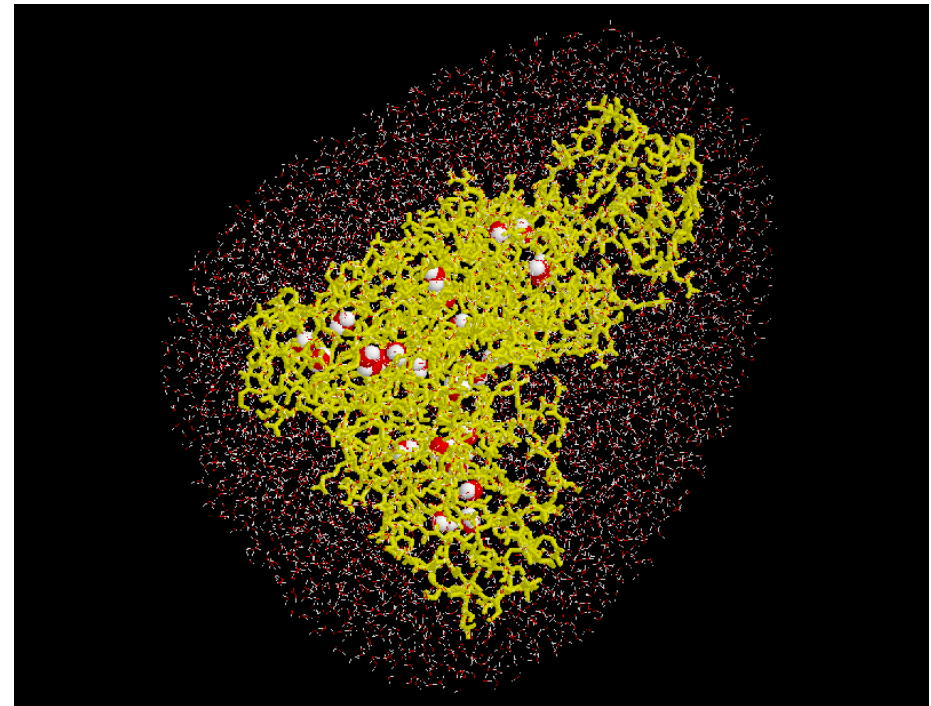


Treatment solvent

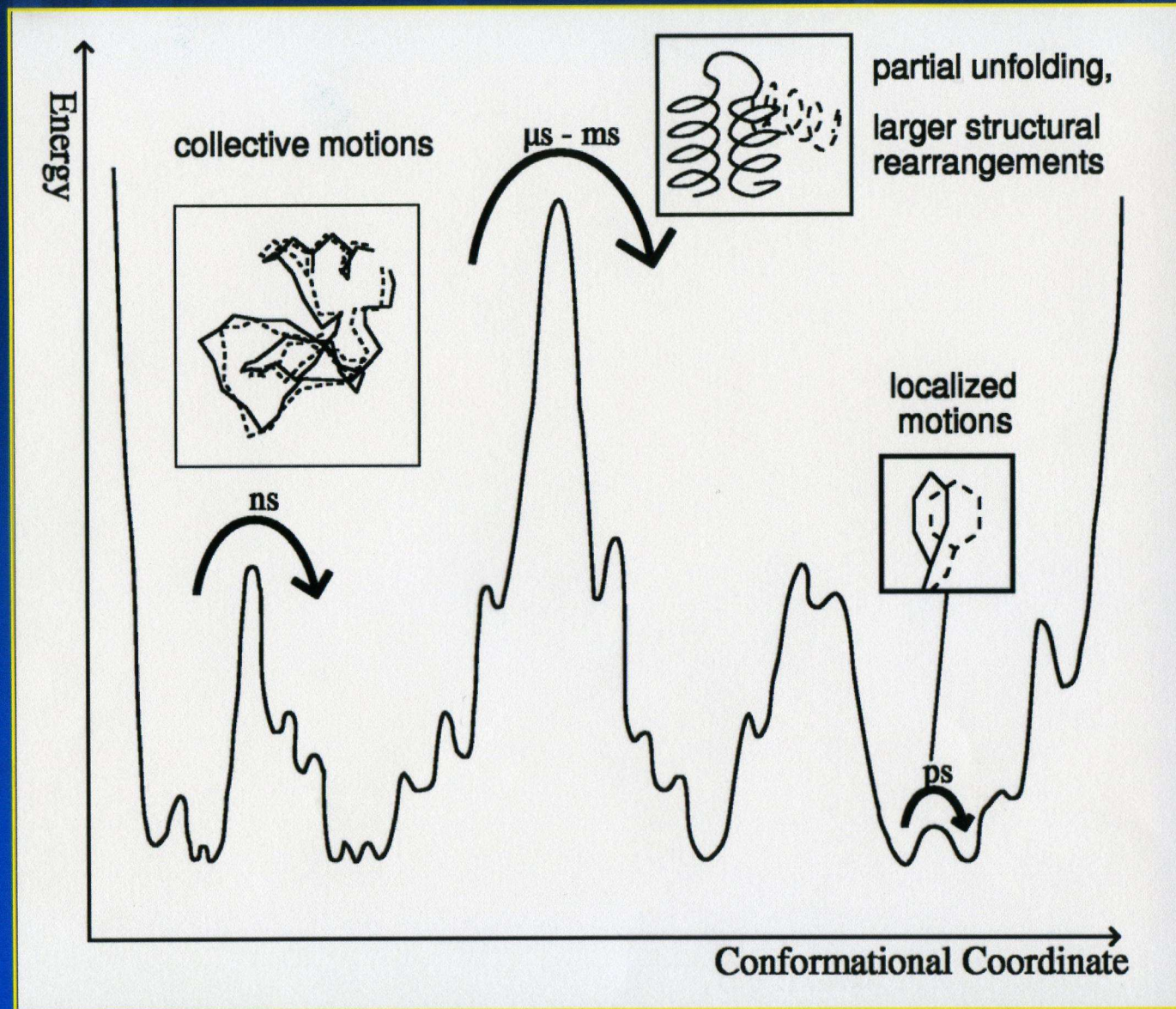


*Box or
Droplet ?
SPC, TIP3P..?*

*explicit
or
implicit?*



Proteins jump between many, hierarchically ordered conformational substates



Protein Folding: Fast Folders

Time Scale:



- Trp-cage, designed mini-protein (20 aa): $4\mu\text{s}$
- β -hairpin of C-terminus of protein G (16 aa) : $6\mu\text{s}$
- Engrailed homeodomain (En-HD) (61 aa): $\sim 27\mu\text{s}$
- WW domains (38-44 aa): $>24\mu\text{s}$
- Fe(II) cytochrome b_{562} (106 aa): extrapolated $\sim 5\mu\text{s}$
- B domain of protein A (58 aa): extrapolated $\sim 8\mu\text{s}$

Ideally: MD simulations in solvent.

MD simulation: a simple application.

$$F = m a = m \frac{dv}{dt} = m \frac{dx}{dt^2}$$

supp $a = dv/dt = \frac{dv}{dt}$

$$v = at + v_0$$

$$x = vt + x_0 = at^2 + v_0 t + x_0$$

with $a = -\frac{1}{m} \frac{dV}{dx}$.

To calculate a trajectory, one needs.

- initial positions of the particles
- initial distribution of velocities
- the gradient of the potential energy function

Molecular Dynamics Simulation

Molecule: (classical) N-particle system

Newtonian equations of motion:

$$m_i \frac{d^2}{dt^2} \vec{r}_i = \vec{F}_i(\vec{r})$$

$$\vec{F}_i(\vec{r}) = -\nabla_i V(\vec{r})$$

with

$$\vec{r} = (\vec{r}_1, \dots, \vec{r}_N)$$

Integrate numerically via the „leapfrog“ scheme

$$\begin{aligned} \mathbf{v}(t + \frac{\Delta t}{2}) &= \mathbf{v}(t - \frac{\Delta t}{2}) + \frac{\mathbf{F}(t)}{m} \Delta t \\ \mathbf{r}(t + \Delta t) &= \mathbf{r}(t) + \mathbf{v}(t + \frac{\Delta t}{2}) \Delta t \end{aligned}$$

with
 $\Delta t \approx 1\text{fs!}$

(equivalent to the Verlet algorithm)

$\Delta t = 1-2 \text{ fs}$

Molecular dynamics

36 amino acids
+3,000 water
molecules

1- μs simulation
(exptl. 5-10 μs)

256 parallel processors
on a CRAY T3E

2 months of computer
time

Y.Duan and
P.A.Kollman, *Science*,
282,740 (1998)

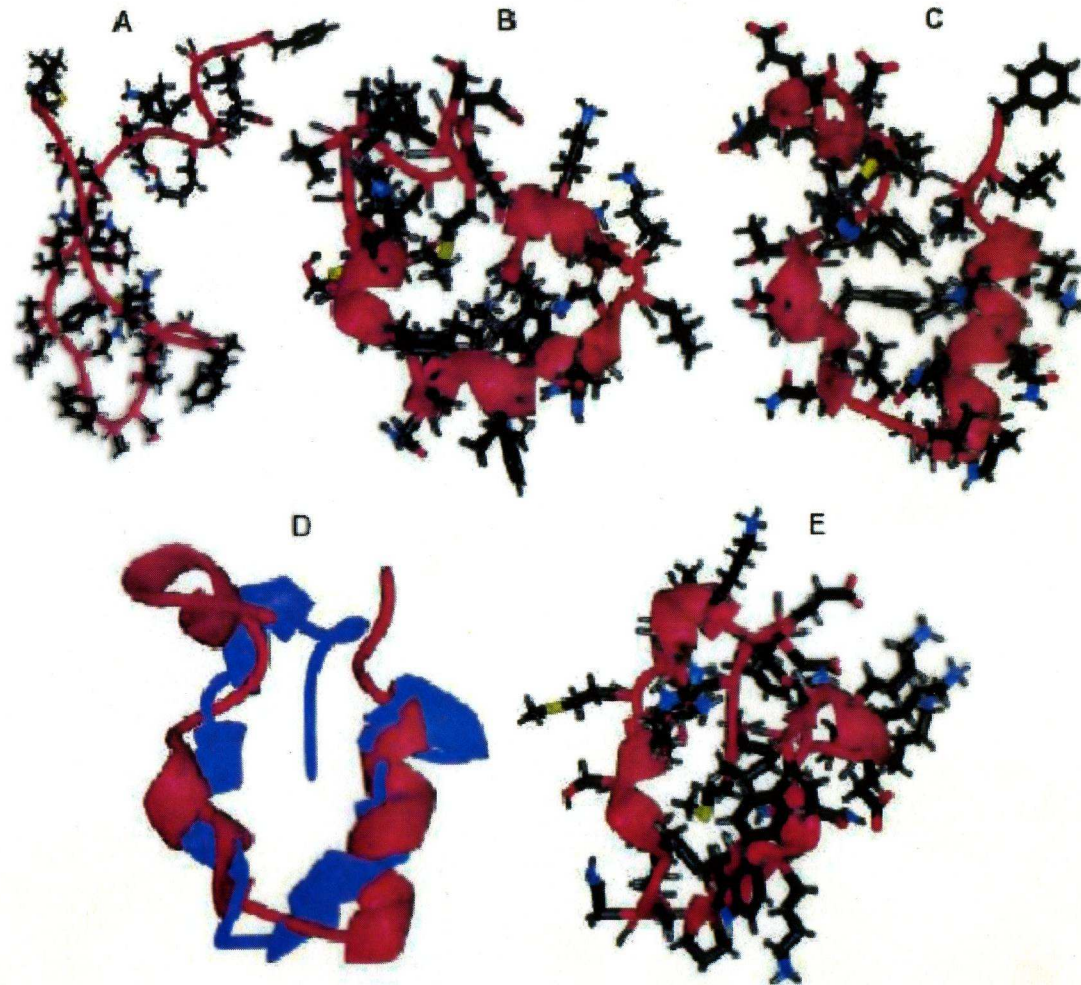


Fig. 1. Ribbon representations of (A) the unfolded, (B) partially folded (at 980 ns), and (C) native structures, and (E) a representative structure of the most stable cluster and (D) the overlap of the native (red) and the most stable cluster (blue) structures, generated with UCSF MidasPlus. Color code [except (D)]: red, main chain atoms and oxygen; black, non-main chain carbon; blue, non-main chain nitrogen; gray, hydrogen; yellow, sulfur.

Atomic-Level Characterization of the Structural Dynamics of Proteins

David E. Shaw,^{1,2*} Paul Maragakis,^{1†} Kresten Lindorff-Larsen,^{1†} Stefano Piana,^{1†} Ron O. Dror,¹ Michael P. Eastwood,¹ Joseph A. Bank,¹ John M. Jumper,¹ John K. Salmon,¹ Yibing Shan,¹ Willy Wriggers¹

Molecular dynamics (MD) simulations are widely used to study protein motions at an atomic level of detail, but they have been limited to time scales shorter than those of many biologically critical conformational changes. We examined two fundamental processes in protein dynamics—protein folding and conformational change within the folded state—by means of extremely long all-atom MD simulations conducted on a special-purpose machine. Equilibrium simulations of a WW protein domain captured multiple folding and unfolding events that consistently follow a well-defined folding pathway; separate simulations of the protein's constituent substructures shed light on possible determinants of this pathway. A 1-millisecond simulation of the folded protein BPTI reveals a small number of structurally distinct conformational states whose reversible interconversion is slower than local relaxations within those states by a factor of more than 1000.

Many biological processes involve functionally important changes in the three-dimensional structures of proteins. Conformational changes associated with protein

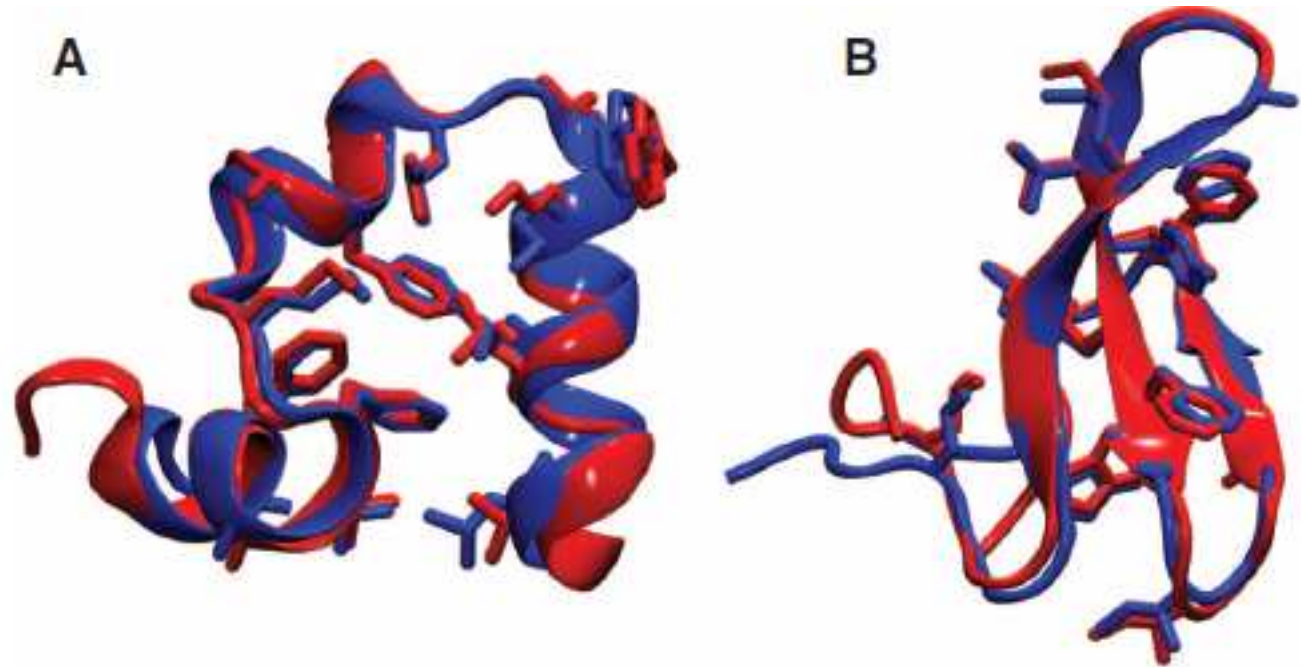
folding (1), signal transduction (2), the catalytic cycles of enzymes (3), and the operation of molecular machines and motor proteins (4) often involve transitions among two or more structur-

ally distinct states characterized as “basins” in the energy landscape. Substantial progress has come from both experimental and computational techniques, in characterizing the states and the ways they interconvert. It has proved difficult to experimentally characterize these states and to elucidate the mechanisms and timescales of the transitions between states.

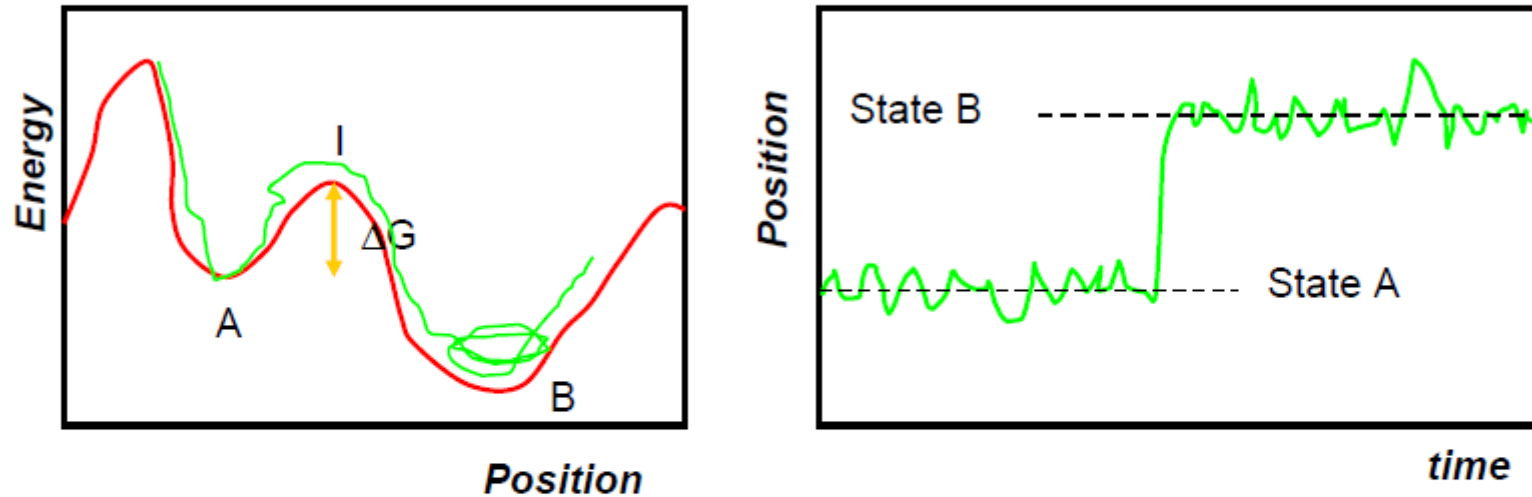
All-atom molecular dynamics simulations are designed to provide a detailed view of the protein's conformational space (5), providing a direct view of the protein's potential energy landscape. Long simulation shots generate a large number of conformational snapshots.

¹D. E. Shaw Research, IBM Research, Armonk, New York 10516, USA. ²Center for Computational Molecular Science, Columbia University, New York, New York 10027, USA. *To whom correspondence should be addressed. Email: david.shaw@des.berkeley.edu. †These authors contributed equally to this work.

Fig. 1. Folding proteins at x-ray resolution, showing comparison of x-ray structures (blue) (15, 24) and last frame of MD simulation (red): (A) simulation of villin at 300 K, (B) simulation of FiP35 at 337 K. Simulations were initiated from completely extended structures. Villin and FiP35 folded to their native states after 68 μ s and 38 μ s, respectively, and simulations were continued for an additional 20 μ s after the folding event to verify the stability of the native fold.



Crossing energy barriers



The actual transition time from A to B is very quick (a few pico seconds).

What takes time is waiting. The average waiting time for going from A to B can be expressed as:

$$\tau_{A \rightarrow B} = C e^{\frac{\Delta G}{kT}}$$