
1

GENERAL OVERVIEW OF BASIC CONCEPTS IN MOLECULAR BIOPHYSICS

This introductory chapter provides a brief overview of the basic concepts and current questions facing biophysicists in terms of the structural characterization of proteins, protein folding, and protein–ligand interactions. Although this chapter is not meant to provide an exhaustive coverage of the entire field of molecular biophysics, the fundamental concepts are explained in some detail to enable anyone not directly involved with the field to understand the important aspects and terminology.

1.1. COVALENT STRUCTURE OF BIOPOLYMERS

Biopolymers are a class of polymeric materials that are manufactured in nature. Depending on the building blocks (or *repeat units* using polymer terminology), biopolymers are usually divided into three large classes. These are (1) polynucleotides (built of nucleotides); (2) peptides and proteins (built of amino acids); and (3) polysaccharides (built of various saccharide units). This chapter only considers general properties of biopolymers using peptides and proteins as examples; questions related to polynucleotides and polysaccharides will be discussed in some detail in Chapter 8.

All polypeptides are linear chains built of small organic molecules called *amino acids*. There are 20 amino acids that are commonly considered *canonical* or *natural* (Table 1.1). This assignment is based upon the fact that these 20 amino acids correspond to 61 (out of total 64) codons within the triplet genetic code with three remaining codons functioning as terminators of protein synthesis (1,2), although there are at least as many other amino acids that occur less frequently in living organisms (Table 1.2). Noncanonical amino acids

are usually produced by chemical modification of a related canonical amino acid (e.g., oxidation of proline produces hydroxyproline), although at least two of them (selenocysteine and pyrrolysine) should be considered canonical based on the way they are utilized in protein synthesis *in vivo* by some organisms (3,4). Furthermore, new components can be added to the protein biosynthetic machinery of both prokaryotes and eukaryotes, which makes it possible to genetically encode unnatural amino acids *in vivo* (5,6). A peculiar structural feature of all canonical (with the exception of glycine) and most noncanonical amino acids is the presence of an asymmetric carbon atom (C_{α}), which should give rise to two different enantiomeric forms. Remarkably, all canonical amino acids are of the L-type. The D-forms of amino acids can also be synthesized *in vivo*, and are particularly abundant in fungi; however, these amino acids do not have access to the genetic code. The rise and persistence of homochirality in the living world throughout the entire evolution of life remains one of the greatest puzzles in biology; examples of homochirality at the molecular level also include almost exclusive occurrence of the D-forms of sugars in the nucleotides, while manifestations of homochirality at the macroscopic level range from specific helical patterns of snail shells to the chewing motions of cows (7,8).

Unlike most synthetic polymers and structural biopolymers (several examples of which will be presented in Chapter 8), peptides and proteins have a very specific sequence of monomer units. Therefore, even though polypeptides can be considered simply as highly functionalized linear polymers constituting a nylon-2 backbone, these functional groups, or *side chains*, are arranged in a highly specific order. All

TABLE 1.1. Chemical Structure and Masses of Natural (Canonical) Amino Acids

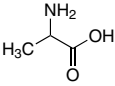
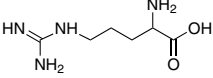
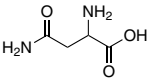
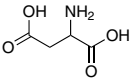
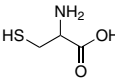
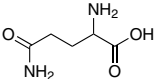
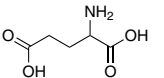
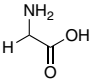
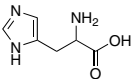
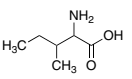
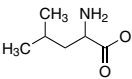
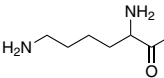
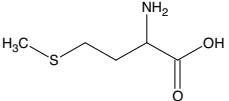
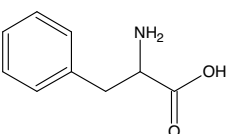
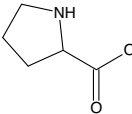
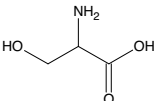
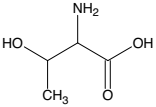
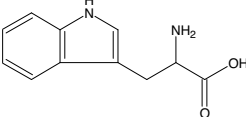
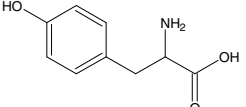
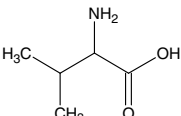
Symbol	Name	Molecular Formula (Residue)	Chemical Structure	Side-Chain Character	Monoisotopic Mass ^a (Residue)	Average Mass (Residue)
Ala (A)	Alanine	C ₃ H ₅ NO		Nonpolar	71.037	71.079
Arg (R)	Arginine	C ₆ H ₁₂ N ₄ O		Basic	156.101	156.188
Asn (N)	Asparagine	C ₄ H ₆ N ₂ O ₂		Polar	114.043	114.104
Asp (D)	Aspartic acid	C ₄ H ₅ NO ₃		Acidic	115.027	115.089
Cys (C)	Cysteine	C ₃ H ₅ NOS		Polar/acidic	103.009	103.145
Gln (Q)	Glutamine	C ₅ H ₈ N ₂ O ₂		Polar	128.059	128.131
Glu (E)	Glutamic acid	C ₅ H ₇ NO ₃		Acidic	129.043	129.116
Gly (G)	Glycine	C ₂ H ₃ NO		Nonpolar	57.021	57.052
His (H)	Histidine	C ₆ H ₇ N ₃ O		Basic	137.059	137.141
Ile (I)	Isoleucine	C ₆ H ₁₁ NO		Nonpolar	113.084	113.160
Leu (L)	Leucine	C ₆ H ₁₁ NO		Nonpolar	113.084	113.160
Lys (K)	Lysine	C ₆ H ₁₂ N ₂ O		Basic	128.095	128.174
Met (M)	Methionine	C ₅ H ₉ NOS		Nonpolar/ amphipathic	131.040	131.199
Phe (F)	Phenylalanine	C ₉ H ₉ NO		Nonpolar	147.068	147.177
Pro (P)	Proline	C ₅ H ₇ NO		Nonpolar	97.053	97.117
Ser (S)	Serine	C ₃ H ₅ NO ₂		Polar	87.032	87.078

TABLE 1.1. (Continued)

Symbol	Name	Molecular Formula (Residue)	Chemical Structure	Side-Chain Character	Monoisotopic Mass ^a (Residue)	Average Mass (Residue)
Thr (T)	Threonine	C ₄ H ₇ NO ₂		Polar/amphipathic	101.048	101.105
Trp (W)	Tryptophan	C ₁₁ H ₁₀ N ₂ O		Amphipathic	186.079	186.213
Tyr (Y)	Tyrosine	C ₉ H ₉ NO ₂		Amphipathic	163.063	163.176
Val (V)	Valine	C ₅ H ₉ NO		Nonpolar	99.068	99.133

^a See Chapter 3 for a definition of monoisotopic and average masses.

naturally occurring proteins consist of an exact sequence of amino acid residues linked by peptide bonds (Fig. 1.1a), which is usually referred to as the *primary structure*. Some amino acids can be modified after translation (termed *post-translational modification*), for instance, by phosphorylation, methylation, or glycosylation. Among these modifications, formation of the covalent bonds between two cysteine residues is particularly interesting, since such *disulfide bridges* can stabilize protein geometry, by bringing together residues that are distant in the primary structure into close proximity in three-dimensional (3D) space. The highly specific spatial organization of many (but not all) proteins under certain conditions is often referred to as *higher order structure* and is another point of distinction between them (as well as most biological macromolecules) and synthetic polymers. Although disulfide bridges are often important contributors to the stability of the higher order structure, correct protein folding does not necessarily require such covalent “stitches”. In fact, cysteine is one of the least abundant amino acids, and many proteins lack it altogether. As it turns out, relatively weak noncovalent interactions between functional groups of the amino acid side chains and the polypeptide backbone are much more important for the highly specific arrangement of the protein in 3D space. Section 1.2 provides a brief overview of such interactions.

1.2. NONCOVALENT INTERACTIONS AND HIGHER ORDER STRUCTURE

Just like all chemical forces, all inter- and intramolecular interactions involving biological macromolecules (both

covalent and noncovalent) are electrical in nature and can be described generally by the superposition of Coulombic potentials. In practice, however, the noncovalent interactions are subdivided into several categories, each being characterized by a set of unique features.

1.2.1. Electrostatic Interaction

The term *electrostatic interaction* broadly refers to a range of forces exerted among a set of stationary charges and/or dipoles. The interaction between two fixed charges q_1 and q_2 separated by a distance r is given by the Coulomb law:

$$E = \frac{q_1 q_2}{4\pi\epsilon_0\epsilon r} \quad (1-2-1)$$

where ϵ_0 is the absolute permittivity of vacuum [$8.85 \times 10^{-12} \text{ C}^2/\text{N}\cdot\text{m}$ in Système International (SI)] and ϵ is the dielectric constant of the medium. Although the numerical values of the dielectric constants of most homogeneous media are readily available, the use of this concept at the microscopic level is not very straightforward (9,10). The dielectric constant is a measure of the screening of the electrostatic interaction due to the polarization of the medium, hence the difficulty in defining a single constant for a protein, where such screening depends on the exact location of the charges, their environment, and so on. Although in some cases the values of the “effective” dielectric constants for specific protein systems can be estimated based on experimental measurements of the electrostatic interactions, such an approach has been disfavored by many for a long time (11). This book will follow the example set by Daune (12) and will write all expressions with $\epsilon = 1$.

TABLE 1.2. Chemical Structure and Masses of Some Less Frequently Occurring Natural (Noncanonical) Amino Acids

Symbol	Name	Molecular Formula (Residue)	Chemical Structure	Side-Chain Character	Monoisotopic Mass (Residue)	Average Mass (Residue)
Abu	2-Aminobutyric acid	C ₄ H ₇ NO		Nonpolar	85.053	85.106
Dha	Dehydroalanine	C ₃ H ₃ NO		Nonpolar	69.021	69.063
Hse	Homoserine	C ₄ H ₇ NO ₂		Polar	101.048	101.105
Hyp	Hydroxyproline	C ₆ H ₁₂ N ₂ O ₂		Polar	144.090	144.174
Nle	Norleucine	C ₆ H ₁₁ NO		Nonpolar	113.084	113.160
Orn	Ornithine	C ₅ H ₁₀ N ₂ O		Basic	114.079	114.147
Pyr	Pyroglutamic acid	C ₅ H ₅ NO ₂		Moderately polar	111.032	111.100
Pyl	Pyrrolysine	C ₁₁ H ₁₆ N ₃ O ₂ + R (NH ₂ , OH, or CH ₃)		Polar		
Sec	Selenocysteine	C ₃ H ₅ NOSe		Polar/acidic	144.960 (150.954 ^a)	150.039

^aMost abundant.

Interaction between a charge q and a permanent dipole p separated by a distance r is given by

$$E = -\frac{qp \cdot \cos\theta}{4\pi\epsilon_0 r^2} \quad (1-2-2)$$

where θ is the angle between the direction of the dipole and the vector connecting it with the charge q . If the dipole is not fixed directionally, it will align itself to minimize the energy Eq. (1-2-2), that is, $\theta = 0$. However, if such energy is small compared to thermal energy, Brownian motion will result in

MTTASTSQVR QNYHQDSEEA INAQINLELY ASYVYLSMSY YFDRDDVALK NFAKYFLHQS
 HEEREHAEKL MKLQNQRGGR IFLQDIKKPD CDDWESGLNA MECALHLEKN VNQSLELHK
 LATDKNDPHL CDFIETHYLN EQVKAIKELG DHVTNLRKMG APESGLAEYL FDKHTLGSD NES

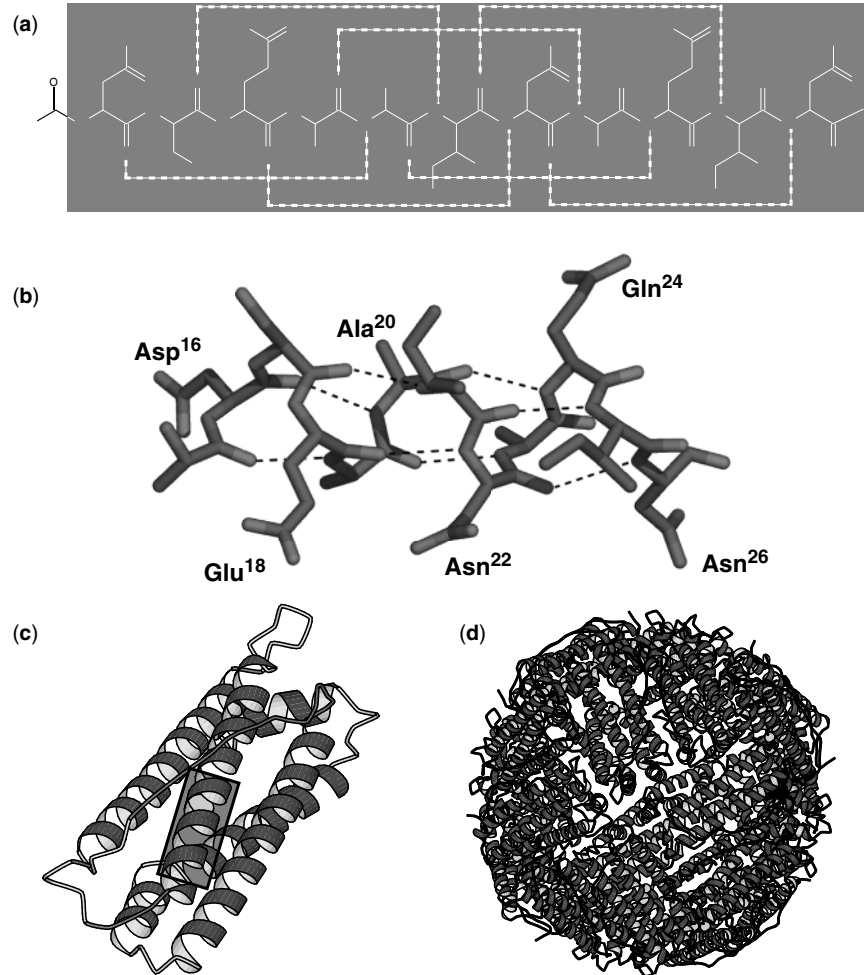


Figure 1.1. Hierarchy of structural organization of a protein (H-form of human ferritin). Amino acid sequence determines the primary structure (a). Covalent structure of the 11 amino acid residue long segment of the protein (Glu¹⁶ → Asn²⁶) is shown in the shaded box. A highly organized network of hydrogen bonds along the polypeptide backbone (shown with dotted lines) gives rise to secondary structure, α -helix (b). A unique spatial arrangement of the elements of the secondary structure gives rise to the tertiary structure, with the shaded box indicating the position of the (Glu¹⁶ → Asn²⁶) segment (c). Specific association of several folded polypeptide chains (24 in the case of ferritin) produces the quaternary structure (d).

the averaging of all values of θ with only a small preference for those that minimize the electrostatic energy, resulting in a much weaker overall interaction:

$$E = -\frac{q^2 p^2}{(4\pi\epsilon_0)^2 \cdot 3k_B T r^4} \quad (1-2-3)$$

where T = temperature and k_B is the Boltzmann constant.

Interaction between two dipoles, p_1 and p_2 , separated by a distance r in this approximation will be given by

$$E = -\frac{2p_1^2 p_2^2}{(4\pi\epsilon_0)^2 \cdot 3k_B T r^6} \quad (1-2-4)$$

while the interaction between the two fixed dipoles will be significantly stronger ($\sim 1/r^3$).

Polarization of a molecule can also be viewed in terms of electrostatic interaction using a concept of *induced dipoles* (12). Such interaction is, of course, always an attractive force, which is inversely proportional to r^4 (for a charge-induced dipole interaction) or r^6 (for a permanent dipole-induced dipole interaction). Finally, interaction between two polarizable molecules can be described in terms of a weak induced dipole-induced dipole interaction.

1.2.2. Hydrogen Bonding

The electrostatic interactions considered in the preceding sections can be treated using classical physics. *Hydrogen bonding* is an example of a specific noncovalent interaction that cannot be treated within the framework of classical electrostatics. It refers to an interaction occurring between a proton donor group ($-\text{OH}$, $-\text{NH}_3^+$, etc.) and a proton acceptor atom that has an unshared pair of electrons. Although hydrogen-bond formation (e.g., $\text{R}=\ddot{\text{O}}:\cdots\text{H}-\text{NR}_2$) may look like a simple electrostatic attraction of the permanent dipole-induced dipole type, the actual interaction is more complex and involves charge transfer within the proton donor-acceptor complex. The accurate description of such exchange interaction requires the use of sophisticated apparatus of quantum mechanics.

The importance of hydrogen bonding as a major determinant and a stabilizing factor for the higher order structure of proteins was recognized nearly 70 years ago by Mirsky and Pauling, who wrote in 1936: “*the [native protein] molecule consists of one polypeptide chain which continues without interruption throughout the molecule . . . this chain is folded into a uniquely defined configuration, in which it is held by hydrogen bonds between the peptide nitrogen and oxygen atoms . . .*” (13). Considerations of the spatial arrangements that maximize the amount of hydrogen bonding within a polypeptide chain later led Pauling to predict the existence of the α -helix, one of the most commonly occurring local motifs of higher order structure in proteins (14). Hydrogen bonds can be formed not only within the macromolecule itself, but also between biopolymers and water molecules (the latter act as both proton donors and acceptors). Hydrogen bonding is also central for understanding the physical properties of water, as well as other protic solvents.

1.2.3. Steric Clashes and Allowed Conformations of the Peptide Backbone: Secondary Structure

Both electrostatic and hydrogen-bonding interactions within a flexible macromolecule would favor 3D arrangements of its atoms that minimize the overall potential energy. However, there are two fundamental restrictions that limit the

conformational freedom of the macromolecule. One is, of course, the limitation imposed by covalent bonding. The second is steric hindrance, which also restricts the volume of conformational space available to the biopolymer. This section considers the limits imposed by steric clashes on the conformational freedom of the polypeptide backbone.

The peptide amide bond is represented in Figure 1.1a as a single bond (i.e., C–N), however, it actually has a partial double-bond character in a polypeptide chain due to partial delocalization of electron density across the neighboring carbonyl group. The double-bond character of the C–N linkage, as well as the strong preference for the trans configuration of the amide hydrogen and carbonyl oxygen atoms,* result in four atoms lying coplanar. Figure 1.2 shows successive planes linked by the C_α atom of the i th amino acid residue. The two degrees of freedom at this junction are usually referred to as φ_i and ψ_i angles and the backbone conformation of the polypeptide composed of n amino acid residues can be described using $n - 1$ parameters (pairs of φ_i and ψ_i). Steric restrictions limit the conformational volume accessible to polypeptides, which is usually represented graphically on the (φ, ψ) plane using so-called *conformational maps* or *Ramachandran plots* (15). An example of such a diagram, shown in Figure 1.3, clearly indicates that only a very limited number of configurations of the polypeptide backbone are allowed sterically.

Several regions within the accessible conformational volume are of particular interest, since they represent the structures that are stabilized by highly organized networks of hydrogen bonds. The α -helix is one of such structures, where the carbonyl oxygen atom of the i th residue is hydrogen bonded to the amide of the $(i+4)$ th residue (Fig. 1.1b). This local motif, or spatial arrangement of a segment of the polypeptide backbone, is an example of a *secondary structure*, which is considered the first stage of macromolecular organization to form higher order structure. Another commonly occurring element of the secondary structure is located within a larger island of sterically allowed conformations on the Ramachandran plot. Such conformations [upper left corner on the (φ, ψ) plane in Fig. 1.3] are rather close to the fully extended configuration of the chain and, therefore, cannot be stabilized by local hydrogen bonds. Nevertheless, formation of strong stabilizing networks of hydrogen bonds becomes possible if two strands are placed parallel or antiparallel to each other, forming so-called *β -pleated sheets*.

*Proline is an exception to this rule. As an imino acid its side chain is also bonded to the nitrogen atom. Thus, the cis and trans forms are almost isoenergetic, leading to the possibility of *cis*-Xaa–Pro bonds in folded proteins, and statistically at the level of 5–30% in unstructured polypeptides.

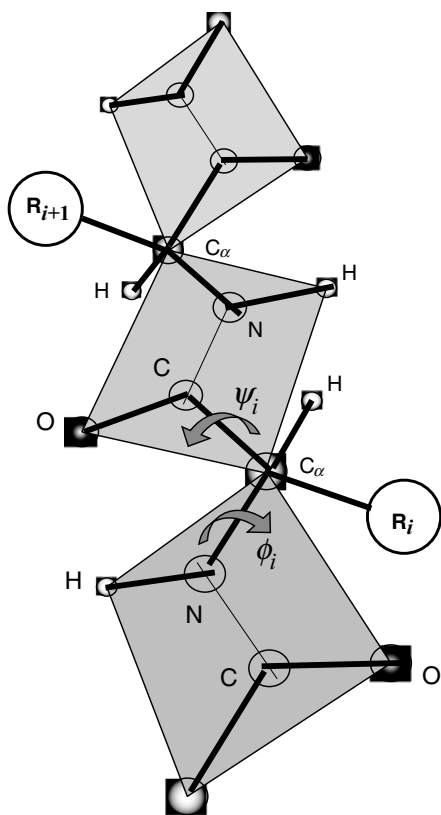


Figure 1.2. Peptide bond and the degrees of freedom determining the polypeptide backbone conformation.

The third important local structural motif is the *turn*, which causes a change in the chain direction within a folded protein. Whereas loops are generally flexible sections of chain, turn structures tend to be more rigid and are stabilized

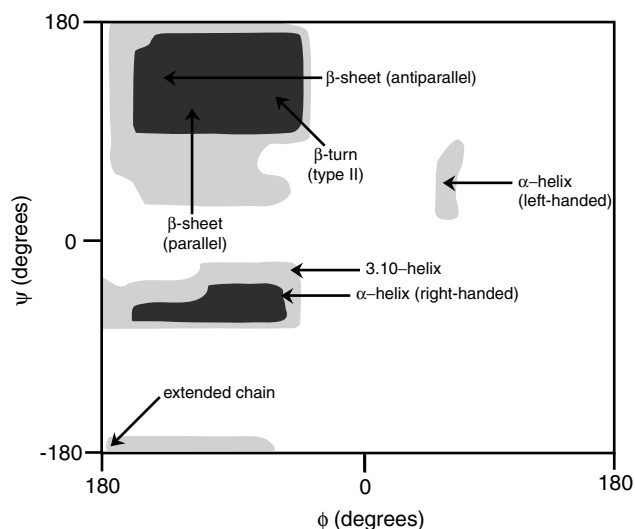


Figure 1.3. A schematic representation of the Ramachandran plot.

by hydrogen bonding or specific side-chain interactions. These turn structures can be highly important, particularly in antiparallel β -sheet structures, where a complete reversal of the chain is required to enable packing of adjacent strands. Other less frequently occurring elements of secondary structure (e.g., 3_{10} or π helices) can also be identified on the Ramachandran plot.

So far, we have largely ignored the contributions of the amino acid side chains to protein conformation. One obvious consequence of the existence of a variety of different side chains is the dependence of the Ramachandran plots for each particular (ϕ_i, ψ_i) pair on the identity of the i th amino acid residue. For example, a significantly larger conformational volume is available to glycine as compared to amino acid residues with bulky side chains. Furthermore, different side chains placed at “strategic” locations may exert a significant influence on the stability of the secondary structural elements. We will illustrate this point using the α -helix as an example. All hydrogen bonds in an α -helix are almost parallel to each other (and to the axis of the helix). This highly ordered pattern of hydrogen bonding results in a noticeable dipole moment, with the N-terminal end of the helix being a positive pole. Obviously, the presence of a positively charged residue at or near the N-terminal end of the helix will destabilize it due to the unfavorable charge–permanent dipole interaction (Eq. 1-2-2). On the other hand, the presence of a negatively charged residue will be energetically favorable and will increase the stability of the helix. Likewise, the presence of charged residues at or near the C-terminal end of the helix will also have a significant influence on the stability of this element of secondary structure. Note, however, that uncharged side chains may also be very important determinants of the higher order structure of proteins and polypeptides due to the so-called *hydrophobic interactions*. These will be considered in Section 1.2.4.

1.2.4. Solvent–Solute Interactions, Hydrophobic Effect, Side Chain Packing, and Tertiary Structure

The term *hydrophobic effect* (16–19) refers to a tendency of nonpolar compounds (e.g., nonpolar amino acid side chains, Table 1.1) to be sequestered from polar solutions (e.g., aqueous solution) into an organic phase. Such behavior is ubiquitous in nature and has been observed and described at least 2 millennia ago, although the term hydrophobic was coined only in 1915 (18). The initial view of the hydrophobic interaction was rather simplistic and implied attraction between like media (e.g., oil–oil attraction). A very different view, which is now commonly accepted, was proposed in the mid-1930s by Hartley, who suggested that nonpolar species are excluded from polar solvent because of their inability to compete with the strong interaction between the polar molecules themselves (20). In

Tanford's words, "*antipathy between hydrocarbon and water rests on the strong attraction of water for itself*" (21). An intriguing aspect of the hydrophobic interaction is that the placement of a hydrocarbon molecule in water may be enthalpically favorable. This fact was the basis for a widespread skepticism over the concept of hydrophobic interactions, although such views did not prevail (22). It is now understood that solvent–solute affinity is determined by the free energy (not the enthalpy alone), and it is the unfavorable free energy that leads to the observed disaffinity of water and nonpolar solutes.

Various microscopic explanations of the hydrophobic effect are usually based on the *frozen water patches* or *microscopic iceberg* model proposed originally by Frank and Evans (23). They suggested that placing a nonpolar solute in water creates a loose "cage" of first-shell water molecules around it. The creation of such a cage has a significant entropic price due to the forced ordering of water, hence the overall unfavorable free energy (despite a favorable enthalpic term). Readers interested in a more detailed account of the physics of hydrophobicity and related phenomena are referred to an excellent tutorial by Southall, Dill, and Haymet (18).

Although the initial work on the hydrophobic effect was focused on hydrocarbons, its main results and conclusions can be easily extended to nonpolar side chains of polypeptides and proteins, which are buried into a *hydrophobic core* of a folded or collapsed protein molecule in order to eliminate, or at least minimize, any contacts with the polar solvent. A very interesting historical account of the elucidation of the nature of the hydrophobic interaction and its role in protein folding can be found in an excellent review by Tanford (24). Hydrophobic side chains are generally more stable if sequestered away from the solvent in protein cores. Proteins tend to be very well-packed molecules so the side-chain atoms sequestered from the solvent must come into close contact with each other, hence the term *hydrophobic packing*. At the same time, hydrophilic residues usually decorate the solvent-exposed surface of the protein. This decoration is achieved by combining the elements of secondary structure (α -helices, β -sheets, and turns) in a unique 3D arrangement, or *tertiary structure*. It is the tertiary structure that affords proteins their unique biological function, whether it be purely structural, the precise spatial organization of side chains to effect catalysis of a reaction, presentation of a surface or loop for signaling or inhibition, creating a cavity or groove to bind ligand, or any of the other vast range of functions that proteins can perform.

Hydrophobic interaction is, of course, not the only driving force giving rise to a unique tertiary structure. Additional stabilization is afforded by the close proximity of acidic and basic residues, which is frequently observed in the folded structure, enabling the formation of salt

bridges. These can be viewed as charge–charge interactions (Eq. 1-2-1). We have already mentioned that certain elements of secondary structure have intrinsic (permanent) dipole moments. Favorable arrangement of such dipoles with respect to one another (e.g., in the so-called helical bundles) may also become a stabilizing factor (Eq. 1-2-2) in addition to the hydrophobic interaction. It is probably worth mentioning that in the vast majority of proteins, the interactions stabilizing the tertiary structure are cooperative. In other words, significant enthalpic gains are achieved only if several segments of the protein are in close proximity and interact with each other. All such factors have been evolutionarily optimized for each protein, but the important thing to realize is that any one natural protein sequence has only a single most stable conformation, and the genetically encoded primary sequence alone is necessary and sufficient to define the final folded structure of the protein (Fig. 1.4) (25).

Many proteins adopt similar common structural motifs resulting from combinations of secondary structure elements, such as the alternating $\beta\alpha\beta$ structure, 4-helix bundles, or β -barrels. As more and more protein structures are solved, the number of protein architectures increases, although it has been predicted that there are a limited number of fold motifs (26–30). This conclusion is based on the observations that (1) topological arrangements of the elements of secondary structure are highly skewed by favoring very few common connectivities and (2) folds can accommodate unrelated sequences [as a general rule, structure is more robust than sequence (31,32)]. Therefore, the fold universe appears to be dominated by a relatively small number of giant attractors, each accommodating a large number of unrelated sequences. In fact, the total number of folds is estimated to be <2000, of which 500 have been already characterized. Figure 1.5 represents the 15 most populated folds selected on the basis of a structural annotation of proteins from completely sequenced genomes of 20 bacteria, 5 Archaea, and 3 eukaryotes (33).

The existence of a "finite set of natural forms" in the protein world has inspired some to invoke the notion of Platonic forms that are "determined by natural law" (34), a suggestion that seems more poetic than explanatory. What has become clear though is that very similar tertiary structures can be adopted by quite dissimilar primary sequences (33). Protein primary sequences can be aligned and regions identified that are identical or homologous (meaning the chemical nature of the amino acid side chain is similar, e.g., polar, nonpolar, acidic, basic). However, even sequences with quite low homology can have a very similar overall fold, depending on the tertiary interactions that stabilize them. Although tertiary structure is sometimes viewed as the highest level of spatial organization of single-chain (i.e., monomeric) proteins, an even higher level of organization is often seen in larger proteins (generally,

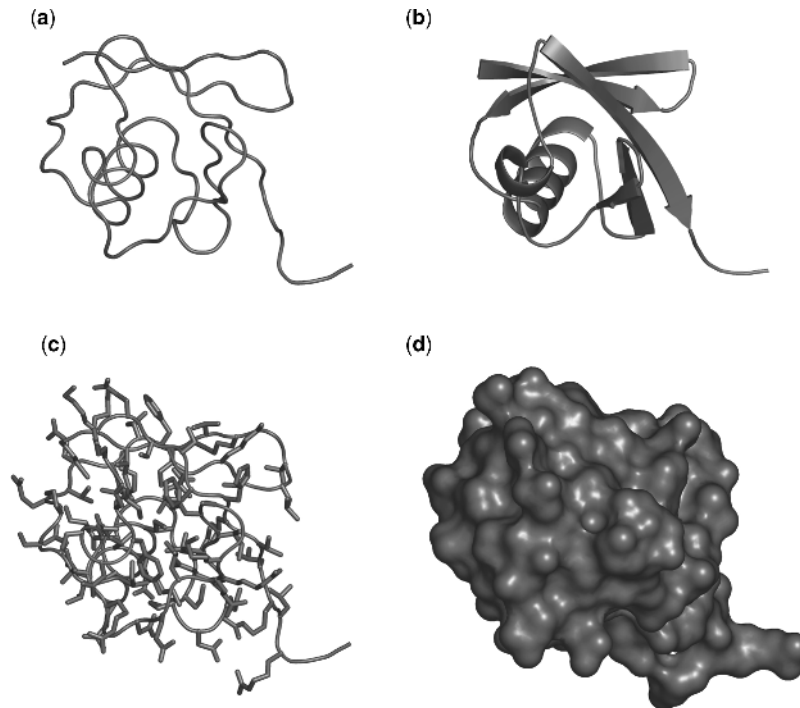


Figure 1.4. Different representations of the higher order structure of natively folded proteins.

>150 amino acid residues). Such proteins form clearly recognizable *domains*, which tend to be contiguous in primary structure and often enjoy a certain autonomy from one another.

1.2.5. Intermolecular Interactions and Association: Quaternary Structure

Above and beyond the folding of monomeric chains, many protein chains can also assemble to form multisubunit complexes, ranging from relatively simple homodimers (example are hemoglobin molecules of primitive vertebrates, e.g., lamprey and hagfish) to large homooligomers (e.g., the iron storage protein ferritin, comprised of 24 identical subunits) to assemblies of different proteins (e.g., ribosomes). Such assemblies are usually considered to be the highest level of molecular organization at the microscopic level, which is usually referred to as *quaternary structure*. Although covalent links are sometimes formed between the monomeric constituents of a multimeric protein assembly (e.g., in the form of disulfide bonds), the noncovalent interactions (discussed in the preceding sections) are usually much more important players.

The archetype of quaternary structure is mammalian hemoglobin, which is a noncovalent tetramer ($\alpha_2\beta_2$) consisting of two pairs of similar monomeric chains (α - and β -globins). The arrangement of monomers in the tetramer, which is in fact a dimer composed of two heterodimers, is crucial for the function of hemoglobin as an oxygen transporter. A tetramer composed of four iden-

tical globins (β_4) can also be formed and is indeed present in the blood of people suffering from some forms of thalassemia. However, this homotetramer (termed hemoglobin H or HbH), lacks the most important characteristic of the “normal” hemoglobin (HbA), namely, high cooperativity of oxygen binding.

1.3. THE PROTEIN FOLDING PROBLEM

1.3.1. What is Protein Folding?

Polymers can adopt different conformations in solution depending on functionality and the interaction with neighboring chains, other parts of the same chain, and the bulk solvent. However, almost all synthetic copolymers (i.e., polymers consisting of more than one type of repeat unit) consist of a range of different length chains and, in many cases, a nonspecific arrangement of monomer groups. On the other hand, the primary structure of a given protein is always the same, creating a homogeneous and highly monodisperse copolymer. Protein sequences are generally optimized to prevent nonspecific intermolecular interactions and individual molecules will fold to adopt a unique stable conformation governed solely by the primary sequence of amino acids. The ability of proteins to attain a unique higher order structure sets them apart from most random copolymers. Most proteins can fold reversibly *in vitro*, without being aided by any sophisticated cellular machinery (e.g., chaperones, which we will consider in

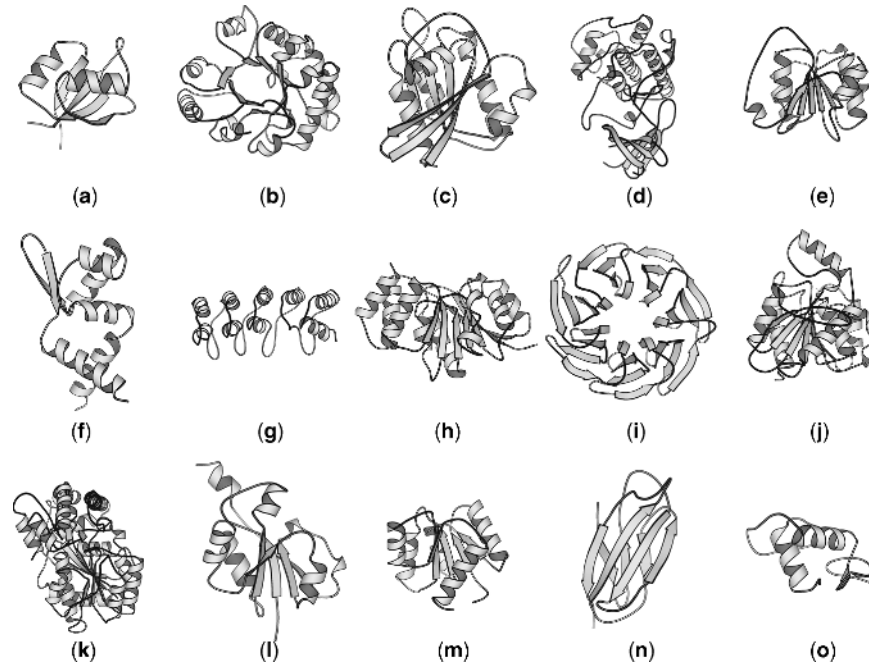


Figure 1.5. The 15 most populated folds selected on the basis of a structural annotation of proteins from the completely sequenced genomes of 20 bacteria, 5 Archaea, and 3 eukaryotes. From left to right and top to bottom, they are ferredoxin-like (4.45%) (a), TIM-barrel (3.94%) (b), P-loop containing nucleotide triphosphate hydrolase (3.71%) (c), protein kinases (PK) catalytic domain (3.14%) (d), NAD(P) (nicotinamide adenine dinucleotide phosphate)-binding Rossmann-fold domains (2.80%) (e), (deoxyribonucleic acid: ribonucleic acid) (DNA:RNA) binding 3-helical bundle (2.60%) (f), α - α superhelix (1.95%) (g), S-adenosyl-L-methionine-dependent methyltransferase (1.92%) (h), 7-bladed β -propeller (1.85%) (i), α/β -hydrolases (1.84%) (j), PLP-dependent transferase (1.61%) (k), adenine nucleotide α -hydrolase (1.59%) (l), flavodoxin-like (1.49%) (m), immunoglobulin-like β -sandwich (1.38%) (n), and glucocorticoid receptor-like (0.97%) (o). The values in parentheses are the percentages of annotated proteins adopting the respective folds. [Reprinted from (33). Copyright © 2001 with kind permission Springer Science+Business Media.]

Chapter 9), suggesting that the folding mechanism is solely determined by the primary structure of the protein, as well as the nature of the solvent. Folded proteins may remain stable indefinitely in most cases, suggesting that the native structures represent the global free energy minima among all kinetically accessible states (35).

Two classic puzzles are usually considered in connection with protein folding: (1) the *Blind Watchmaker's paradox* and (2) the *Levinthal paradox*. The former is named after a classic book by Dawkins (36), an outspoken critic of the *intelligent design* concept (37). It states that biological (function-competent) proteins could not have originated from random sequences. The Levinthal paradox states that the folded state of a protein cannot be found by a random search (38). Both paradoxes have been historically framed in terms of a random search through vast spaces (sequence space in the *Blind Watchmaker's paradox* and conformational space in the *Levinthal's paradox*), and the vastness of the searched space is equated with physical impossibility. Both paradoxes are elegantly solved within the framework

of the *energy landscape* description of the folding process by invoking the notion of a guided search (39). The concept of protein energy landscapes and its relevance to the protein folding problem will be considered in some detail in Section 1.4.

1.3.2. Why Is Protein Folding So Important?

First, one question is Why do we need to understand protein folding? In the post-genomic era, structure determination has become of paramount importance since it leads to a 3D picture of each gene product, and in many cases gives hints as to the function of the protein. However, the static structure only represents the end point of the chemical reaction of protein folding. Polypeptide chains are translated as extended structures from RNA on the ribosome of cells, but How does this unstructured sequence fold into its final biologically active structure? Are specific local structures present in the newly translated chain? Is there a specific pathway or reaction coordinate of protein folding?

The principles that govern the transitions of biopolymers from totally unstructured to highly ordered states, which often include several subunits assembled in a highly organized fashion, remain one of the greatest mysteries in structural biology (40,41). Deciphering this code is key to understanding a variety of biological processes at the molecular level (recognition, transport, signaling and biosynthesis, etc.), since the specificity of biological activity in proteins, as well as other biomolecules, is dictated by their higher order structure.

Aside from the obvious academic interest to biophysicists in discovering exactly how these biological machines work, there are many more practical implications. Only if we understand all of the processes that are involved in producing a biologically active protein can we hope to harness this power by designing proteins with specific functions. It may already be possible computationally to model an ideal binding site or even optimal arrangement of side chains to catalyze a chemical reaction, but without a thorough knowledge of how this site can be placed into an intact protein molecule, we cannot take advantage of the cellular machinery for the design of therapeutic protein drugs, or even molecules that can catalyze otherwise difficult chemical reactions. For instance, there are many enzymes in nature that catalyze reactions with extremely high specificity and efficiency, whereas chemists lag far behind. Hydrogenase enzymes, for example, catalyze the reduction of protons to produce diatomic hydrogen, a reaction that in a laboratory environment requires application of harsh reactants at elevated temperature or pressure, but that within the catalytic center of the protein occurs at physiological temperatures and with remarkably small energy requirements. Obviously, biological organisms have had a much longer time to optimize these processes relative to the chemical industry. If one can understand in detail the roles of each residue in a protein chain for both the folding and dynamics of the molecule, then the possibilities for protein engineering are boundless. Interestingly, manmade sequences quite often lead to proteins that either do not fold at all or are only marginally stable. This result clearly demonstrates the extremely fine balance of forces present, which can be destroyed by just a single amino acid residue substitution, deletion, or insertion.

Another important aspect of understanding protein folding is to find ways of preventing the process from going awry (42–45). An ever-increasing number of pathological conditions that result from misfolding of proteins in the cell are being identified (46–50). Amyloid plaques actually result from the undesirable formation of quaternary structure when a normally monomeric peptide folds incorrectly and self-assembles to form long proteinaceous fibers. Similarly, other proteins, that are not correctly folded may not present the correct binding surface for interaction with their physiological partners. Thus not only correct folding,

but also the correct assembly of proteins, is key to their correct biological function. Even relatively few mutations within a protein sequence may prevent folding to the native structure, and hence prove pathological. In other cases, mutation can reduce the efficiency of folding, or favor an alternative mode of folding that leads to aggregation and deposition of insoluble amyloid plaques within cells. We will consider the issues related to misfolding and aggregation later.

Finally, one more fundamental problem related to protein folding that has become a focal point of extensive research efforts is the prediction of the native structure and function of a protein based on its primary structure. Since the sequence of each natural protein effectively encodes a single tertiary structure, prediction of the latter is, in essence, a global optimization problem, which is similar to one encountered in crystallography and the physics of clusters (51). The complication that arises when such a global optimization methodology is applied to determine the position of the global energy minimum for a protein is the vastness of the system that precludes calculations based on first principles. So far, the most successful methods of structure prediction rely on the identification of a *template protein* of known structure, whose sequence is highly homologous to that of the protein in question. If no template structure can be identified, *de novo* prediction methods can be used, although it remains to be seen if such methods can predict structures to a resolution useful for biochemical applications (52). Prediction of protein function based on its sequence and structure is an even more challenging task, since homologous proteins often have different functions (53).

1.3.3. What Is the Natively Folded Protein and How Do We Define a Protein Conformation?

Before proceeding further with a description of protein folding it would be useful to define some terms commonly used in the field in order to avoid confusion. First, the *native state* of a protein is defined as the fully folded biologically active form of the molecule. This has generally been considered as a single state with a well-defined tertiary structure, as determined by crystallography or nuclear magnetic resonance (NMR) spectroscopy. More recently, researchers have come to appreciate the importance of dynamics within the protein structure. Even the native state is not a static single structure, but may in fact, depending on the protein, have small or even large degrees of flexibility that are important for its physiological function.

Unfortunately, the use of the term *protein conformation* in the literature has become rather inconsistent and often results in confusion. Historically, protein conformation referred to a specific “three-dimensional arrangement of its constituent atoms” (54). This definition, however, is rather

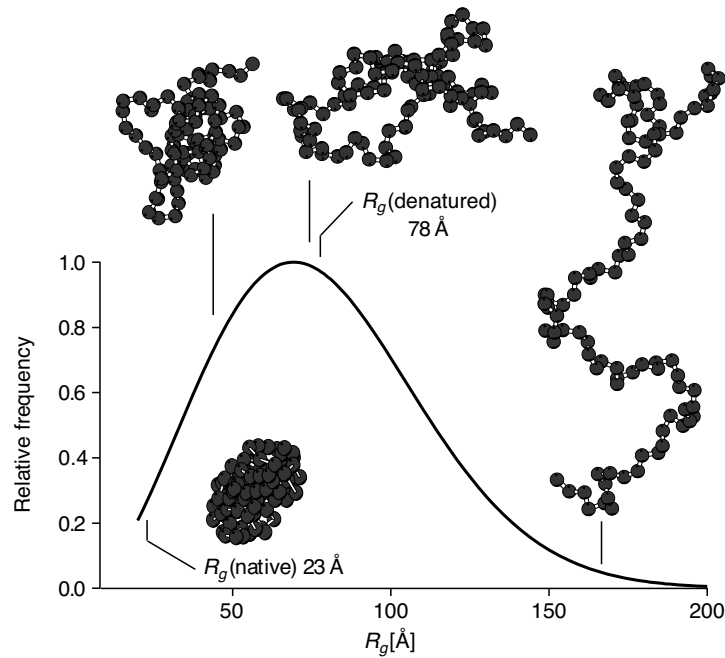


Figure 1.6. Representative configurations of a random coil (a freely joined chain of 100 hard spheres) and the distribution of its radius of gyration R_g . The R_g values of a model protein phosphoglycerate kinase are indicated for comparison. [Adapted from (55) Copyright © 1996 with kind permission from Elsevier.]

narrow, since it does not reflect adequately the dynamic nature of proteins. One particularly annoying complication that arises when conformation is defined using only microscopic terms (e.g., atomic coordinates) is due to the fact that a majority of proteins have segments lacking any stable structure even under native conditions. These could be either the terminal segments that are often invisible in the X-ray structures or flexible loops whose conformational freedom is often required for a variety of functions ranging from recognition to catalysis. In general, it is more than likely that any two randomly selected natively folded protein molecules will *not* have identical sets of atomic coordinates and, as a result, will not be assigned to one conformation if the geometry-based definition is strictly applied. Therefore, it seems that the thermodynamics-based definition of a protein conformation is a better choice. Throughout this book, we will refer to the protein conformation not as a specific microstate, but as a *macrostate*, which can be envisioned as a collection of microstates separated from each other by low energy barriers ($\leq k_B T$). In other words, if one microstate is accessible from another at room temperature, we will consider them as belonging to one conformation, even if there is a substantial difference in their configurations. According to this view, a protein conformation is a continuous subset of the conformational space (i.e., a continuum of well-defined configurations) that is accessible to a protein confined to a certain local minimum. The utility of this definition becomes obvious when

we consider non-native protein conformations, although unfortunately it is not without its own problems.*

1.3.4. What Are Non-Native Protein Conformations? Random Coils, Molten Globules, and Folding Intermediates

In the case of unfolded proteins, which are assumed to be completely nonrigid polypeptide chains, the *random coil* (55), we must consider the ensemble of molecules displaying an impressive variety of configurations (Fig. 1.6). In a truly random coil, as might be the case for a synthetic polymer with identical monomer units in a good solvent, there may well be no conformational preferences for the chain. However, proteins are decorated with side chains of a different chemical nature along their length, such that in water or even in a chemical denaturant one might expect there to be local preferences due to hydrophilic or hydrophobic interactions, and indeed steric effects. Thus for a number of proteins studied in solution, some persistent local and nonlocal conformational effects have been detected,

*For example, this definition is temperature dependent. Indeed, if any two local minima are separated by a high energy barrier ($> k_B T$), the interconversion between these two states does not occur readily at room temperature (T), and these two states should be viewed as two different conformations. However, raising the temperature significantly above room temperature will eventually make passage over this barrier possible, leading to a merging of the two microstates to a single conformation.

indicating that an unfolded protein generally is not in fact a truly random coil. On the other hand, the enthalpy of these interactions is very small in comparison to the entropy of the flexible chain so the overall free energy of each of these conformers will be very similar. On a free energy surface, these would be represented as shallow wells in the generally flat surface of unfolded state free energy.

The relative position of a local energy minimum with respect to the native state gives rise to a further set of descriptions of intermediate states. As a protein folds it may sample stabilizing conformations that contain persistent structure, constituting a local free energy minimum. At the earliest stages of folding there may be only a few interactions that may be very transient: These are termed *early intermediates*. By contrast, species may accumulate further along in the folding process that contain a large although incomplete number of native-like contacts. These are referred to as *late intermediates*, implying that they should form toward the end of the kinetic folding process. There is also the possibility that these local minima arise from stabilizing contacts that are not present in the native protein, and in fact need to be disrupted before the molecule can productively fold. These *off-pathway intermediates* may also arise from intermolecular interactions between folding chains and can lead to nonproductive aggregation that prevents further folding.

The above intermediate states form during folding in the “forward” direction from the unfolded to the native state and, since they are only partially stable, generally do not accumulate sufficiently to be detected other than transiently. It is also possible that such intermediates may form during the reverse process, that is protein unfolding, allowing them to be studied by other methods. Unfortunately, the conditions for unfolding (e.g., chemical denaturant, low pH, high temperature) are generally so harsh that once the stabilizing interactions in the native state have been removed, the unfolding process occurs with high cooperativity and without accumulation of intermediates. However, under mildly denaturing conditions, partially folded states have been detected at equilibrium for a number of proteins, and these have been termed *molten globules* (56). The original definition of the molten globule state was quite rigid: a structural state that has significant secondary structure, but with no fixed tertiary interactions. There are various biophysical tests for this, such as the ability of the protein to bind hydrophobic dyes, consistent with a significant amount of exposed hydrophobic surface area, as would be expected for a partially folded state. The definition has become somewhat relaxed to include many other partially folded ensembles observed, kinetically or at equilibrium, which almost fit the definition. What is clear is that the molten globule itself is a much more dynamic structure than previously thought. Several new concepts have been introduced to reflect the structural diversity and dynamic character of the molten globule state,

such as “a precursor of the molten globule” and “a highly structured molten globule” (57).

One common question that arises is whether the equilibrium molten globule intermediate is actually the same species as that detected in the folding pathway of proteins. Thermodynamically there is nothing to suggest they should be, since the equilibrium by definition is independent of the pathway (58,59). However, comparisons of the characteristics of transient intermediates with the corresponding equilibrium partially folded state have concluded that the similarities are very close, at least for the proteins studied (60–63). Also, a number of states transiently populated by the native state ensemble under mildly destabilizing conditions have been shown to have similarities to folding intermediates. Thus it seems likely that, at least in the later stages of folding, there is indeed some kind of folding–unfolding pathway with specific intermediate states visited in both the folding and unfolding directions.

1.3.5. Protein Folding Pathways

In Section 1.3.4, we began to use the term *folding pathway*, which is understood to be a series of structural changes leading from the fully denatured state of the protein to its native conformation. Introduction of the concept of a folding pathway resolves the Levinthal paradox mentioned earlier by suggesting that the folding process is a directed process involving conformational biases, rather than a merely random conformational search. Despite the vast number of degrees of freedom in macromolecules, the number of folding pathways was initially believed to be rather limited (64). A general scheme of protein folding within this paradigm is presented in terms of rapid equilibration of unfolded protein molecules between different conformations prior to complete refolding. Such equilibria favor certain compact conformations that have lower free energies than other unfolded conformations, and some of these favored conformations are important for efficient folding. The rate-limiting step is thought to occur late in the pathway and to involve a high-energy, distorted form of the native conformation. The latter is a single transition state through which essentially all molecules refold (65).

The classic folding pathway paradigm specifically states that “*proteins are not assembled via a large number of independent pathways, nor is folding initiated by a nucleation event in the unfolded protein followed by rapid growth of the folded structure*” (65). Nevertheless, a large body of experimental evidence now suggests the existence of a large number of folding routes. Furthermore, over the past several years it has become clear that the length of the polypeptide chain is an important factor in determining mechanistic details. Smaller proteins (< 100 residues) appear to prefer a nucleation-type mechanism that does not involve any specific metastable intermediate species (66). On the other

hand, for a number of larger proteins, intermediate states with specific regions of early formed stable structure have been established. If an intermediate is detected, then this argues strongly for a specific protein folding pathway, but in the protein lysozyme, for instance, parallel folding pathways were found, suggesting multiple possible trajectories. In smaller proteins, no observable intermediates accumulate, which might indicate a less directed folding mechanism, or indeed that partially formed structures are too labile to be detected. Based on these observations, a modified view of a folding pathway invokes several common stages of folding, consistent with the progressive development of structure and stability through an ever-slowing set of reactions (67).

Three major models have arisen to attempt to explain this narrowing of the conformational search process (68). The *simple framework* model suggests that secondary structure elements would first form based on their sequence-intrinsic propensities, followed by collision of these preformed structures to form tertiary interactions. Alternatively, *nucleation* of short regions of sequence to form transient secondary structure could act as a template upon which adjacent parts of the chain would condense and propagate structure. The third, *hydrophobic collapse* mechanism calls for the hydrophobic residues to conglomerate nonspecifically to minimize solvent exposure, followed by rearrangement to the final native structure. The actual mechanism of protein folding probably involves some or all of these processes and there is evidence for each mechanism from folding studies of different proteins.

The controversy in regard to whether protein folding follows a specific pathway or whether each molecule follows a completely different trajectory to achieve its final folded state has been elegantly resolved in the so-called *new view of protein folding* (69–71). The centerpiece of this theory is the concept of a protein energy landscape of conformational space, which is discussed in Section 1.4.

1.4. PROTEIN ENERGY LANDSCAPES AND THE FOLDING PROBLEM

1.4.1. Protein Conformational Ensembles and Energy Landscapes: Enthalpic and Entropic Considerations

Classically, a simple chemical reaction is considered to proceed along a reaction coordinate, in which chemical bonds are formed or broken in a well-defined manner, with a transition state at the highest point on the energy profile where bond breaking and formation is occurring, and possibly detectable metastable intermediate structures at local free energy minima. Transition state theory can be applied to relate reaction rates to the heights of the various free energy barriers along the reaction coordinate. Macromolecules are much more complex, however, and the folding of a

protein involves the formation of a large number of interactions that may be either local in nature or indeed involve regions that are quite distant in the polypeptide sequence. Nevertheless, for many years the protein folding reaction was assumed to occur via a similar sequential pathway, perhaps involving a number of intermediate species along the way, but for each unfolded molecule the mechanism of folding to the native state was identical.

Significant advances in theory during the past decade have changed our understanding of the basic principles that govern the protein folding process and have offered an elegant way to resolve the Levinthal paradox (70,72–76). In the case of small organic molecules in a reaction, there are only a small number of conformations available to the reactant species, but for an unfolded protein the conformational space sampled by the unstructured chain is vast by comparison, even for a relatively small protein. Thus it might seem difficult to imagine the chain becoming oriented in such a way as to proceed to fold via a single pathway, and the process would surely be extremely inefficient. Realization that folding may proceed through multiple parallel pathways, rather than a single route, has led to introduction of the concepts of protein energy landscapes (or folding funnels), a cornerstone of the “new view” of protein folding.

Protein folding can be viewed much like any other chemical reaction, which may be represented in 3D by a conformational energy surface: The trajectories on these surfaces lead from reactants (unstructured states) to products (the native state). Because entropy plays a much more significant role in protein folding reactions, it is necessary to consider the free energy, rather than simply potential energy (77). The enthalpic gain of forming hydrogen bonds and making favorable hydrophobic or hydrophilic contacts is compensated by a significant loss of entropy as the chain becomes more and more conformationally restricted. For simple molecules, the entropic term is generally far less significant, but in the case of a folding protein the overall free energy of stabilization in the folded protein may be only a few kilocalories per mole (kcal/mol), being the very small difference between large ΔH (formation of stabilizing interactions) and $T\Delta S$ (loss of entropy upon folding to the native state) terms. The conformational entropy loss for a protein, which continues to adopt a more well-defined 3D structure, is often defined on a “per residue” basis. It can be estimated by using only the backbone entropy [the entropy loss due to side chain packing is significantly less (75)]:

$$\Delta s = s_u - s_f \approx k_B \cdot \ln \left(\frac{\Omega_u}{\Omega_f} \right) \quad (1-4-1)$$

where $s_{u,f}$ and $\Omega_{u,f}$ represent the entropies and the numbers of microstates per residue in the unfolded and natively folded forms of the protein, respectively. Estimations of Δs for small proteins at room temperature give an entropy loss

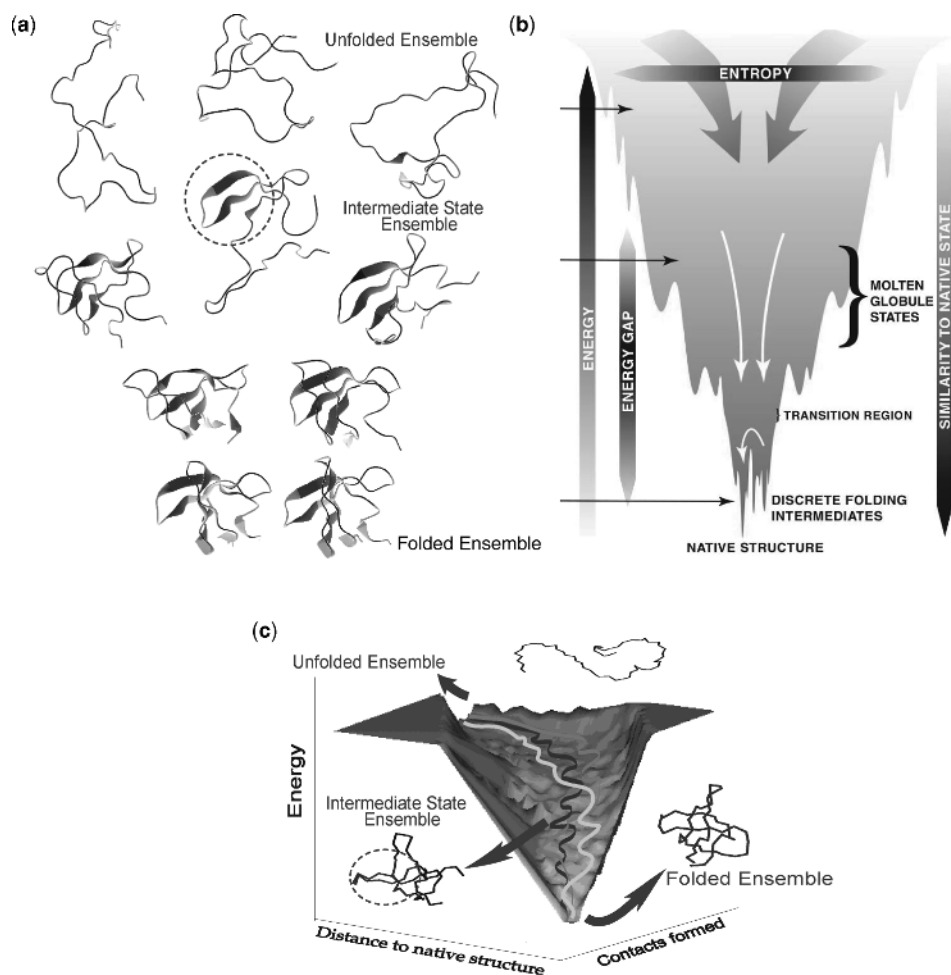


Figure 1.7. Schematic representation of a protein folding funnel. As the large ensemble of structures of an unfolded polypeptide compacts, forming native-like contacts through intermediate states and finally to the native state (a), the energy surface can be schematically represented as seen in panel (b). Multiplicity of folding routes is shown with different folded trajectories on the energy surface (c). However, asymmetry of the surface biases the trajectories toward the preferred route, which can be considered a folding pathway. [Reprinted with permission from (78). Copyright © 2001 C. Clementi and G. Bamberg.]

on the order of tens of joules per mole kelvin residue [$J/(\text{mol}\cdot\text{K}\cdot\text{residue})$] (75). To ensure fast folding despite the unfavorable entropy change, the corresponding free energy surface (or energy landscape) must have the form of a multidimensional “funnel” (Fig. 1.7), where the vertical axis (the depth of the funnel) represents the number of native contacts made (Q) or the relative free energy of the conformational space. The horizontal axis corresponds to the conformational entropy of the system. In this representation, it becomes clear that as the folding chain makes more native contacts, the chain entropy is reduced along with the overall free energy. The native state resides at the bottom of the potential well, characterized by low entropy and a global free energy minimum. A funneled energy landscape is robust to both environmental changes and sequence mutations, since most potentially competing low-energy states

are similar in structure (75), as represented by the multiple local minima at the bottom of the funnel in Figure 1.7.

An important feature of the folding funnel is that its slopes are not always monotonic, hence the competition between the “downhill slide” toward the native fold and the possibility of equally favorable excursions into local free energy minima, depending on the ruggedness of the energy surface. These local minima may represent transient formation of partially folded species, accumulation of intermediates or indeed misfolded forms, depending on the trajectory taken by the chain along the energy surface. This general model allows for multiple pathways with no specific order of structure formation, but subtle changes in the energy surface would lead to a far more directed approach by energetically favoring particular regions of conformational space. It also provides at least one possible solution

for Levinthal's problem of folding time scale. If a favored conformational space were visited on the energy landscape, then it would be sufficiently stabilized to reside there for a longer time, and hence restrict the conformational search for the most stable native structure.

As already stated, one of the most important features that distinguishes these folding funnel models from more traditional reaction coordinates is the key feature of conformational entropy. Thus we are forced to consider each stage along a folding reaction not as a single structure, but instead as a conformational ensemble. The unfolded state (the nearly flat plain at the top of the funnel) demonstrates the large number of microstates in the unfolded polypeptide chain, whereas the deep energy well of the native state may be a single structure or indeed a small ensemble of similar structures closely related in energy. Along the folding trajectory as the chain becomes more ordered, this ensemble become smaller and smaller, but throughout there is always some conformational flexibility that must be considered.

It must be stressed that these folding landscapes are purely theoretical. Based on experimental data it has proved possible to style different landscapes more as a representation than a true physical picture, showing local free energy wells for intermediate species, and alternate possible routes for parallel pathways. Even with the most powerful modern computing facilities, however, it is not yet possible to predict the folding pathway(s) of a protein, although some trajectories in a computational ensemble may appear to fit well with experimental data (79). Only with a complete understanding of the energetics involved in driving a protein to fold will we be able to computationally predict how a given protein will fold.

Semiquantitative models and experiments are revealing how the folding free energy surface is sculpted by the protein sequence and its environment. Although any downhill path from the unfolded state will eventually lead to a native conformation at the bottom of the funnel, the asymmetry (energetic heterogeneity) of the surface can bias the choice of folding routes (80). The existence of such "preferred" folding routes can be observed experimentally and interpreted in terms of folding pathways (see above). The sometimes conflicting demands of folding, structure, and function determine which folding pathways, if any, dominate (76).

1.4.2. Equilibrium and Kinetic Intermediates on the Energy Landscape

Under native conditions, folding of small proteins usually appears to be a highly cooperative process (81,82). By using standard biophysical techniques, only the native and unfolded states are generally sufficiently populated to be detectable. For instance, a titration of chemical denaturant studied by circular dichroism (CD) for most proteins will generate a series of spectra that can be deconvoluted to

contributions solely from the two end points of the unfolding reaction, namely, the native and denatured states. This finding led to the belief that the folding of proteins was not only cooperative, but also two-state, that is, without populated intermediate states.

In contrast, under certain conditions, significant accumulation of an intermediate conformational ensemble may occur if the free energy barrier is sufficiently high, referred to as an *equilibrium intermediate*. These species can be studied in great detail since the rate of conversion is low, allowing significant structural information to be obtained (83). Classic examples of these include the partially folded state of the apo-form of myoglobin (84), acid- and alcohol-induced A-state of ubiquitin (85), or the acid-induced molten globule of α -lactalbumin (86). As already mentioned, these equilibrium intermediate states in some cases may represent important conformations visited along the folding trajectory. Therefore, structural information about these states can give valuable clues as to the nature of the conformational search process. Transient formation or indeed accumulation of certain intermediates is usually induced *in vitro* by simply changing the protein's environment (by varying the solution pH, temperature, presence of chaotropes, etc.) (61,87). This change can result in significant alterations of the energy surface, decreasing the free energy of the intermediate states, and thus increasing their equilibrium population. One needs to be aware, however, that such equilibrium intermediates may differ significantly from the kinetically observed species. As pointed out by Fersht and co-workers (59), an equilibrium intermediate need not by definition be on the preferred kinetic folding pathway since thermodynamic (equilibrium) parameters are independent of mechanism. The goal here though is to make these elusive states more amenable to study using a variety of biophysical techniques, in order to determine not only the conformational preferences of a partially folded protein, but also the possible conformations transiently visited by the native state that may be vital to its *in vivo* function.

Refolding of large proteins often does not conform to the simple two-state model discussed in the beginning of this section. Nevertheless, such proteins may be spontaneously refolded by rapid dilution from a chemically denatured state into native conditions. Kinetic studies of these processes have enabled detection of transient *kinetic intermediates*, which serve as "resting points" in the protein folding process. Since the energy difference between the global minimum and any of the surrounding local minima is usually quite high, the Boltzmann weight of the states that correspond to the local minima is very low. Under native conditions, kinetic intermediates become populated only transiently during refolding experiments. These species have been the focus of close experimental scrutiny, since their structure and behavior may reveal many intimate details of the protein folding process. In order to be detected,

these species must have characteristics different from either the native or unfolded states. For instance, the aromatic residues may experience an environment that makes them hyperfluorescent, or a secondary structure may form in advance of tertiary interactions making the intermediate detectable using stopped-flow optical techniques (e.g., CD). Alternatively, the secondary structure may protect certain amide protons against exchange with bulk solvent, which can be detected by NMR or mass spectrometry (MS), as we will see in the following chapters. The experimental identification and characterization of kinetic intermediates has been the focus of a great deal of research over the past couple of decades, as researchers attempt to glean the determinants of protein folding. Due to their transient nature, however, kinetic intermediates often cannot be observed directly and their properties can only be inferred from indirect measurements.

If the energy barrier separating a kinetic intermediate from the native conformation is high, a significant accumulation of such intermediates may occur. These are referred to as *kinetically trapped* or *metastable* intermediate states. Under certain conditions, excessive accumulation of metastable kinetic intermediates in the course of the folding process may trigger nonspecific interactions among them and, in extreme cases, aggregation. Likewise, aggregation can also be caused by incorrect folding (e.g., due to a sequence mutation). Aggregation processes *in vivo* are prevented by chaperones, a special class of proteins that bind and sequester the misfolded and partially folded polypeptides (88–91). It is important to note, however, that the chaperones only assist, but do not in themselves direct protein folding.

1.5. PROTEIN DYNAMICS AND FUNCTION

1.5.1. Limitations of the Structure–Function Paradigm

Proteins carry out their functions by interacting with other proteins, as well as other molecules ranging from giant biopolymers (e.g., DNA) to small organic molecules and monatomic ligands. In all cases, protein–ligand binding is the first stage of the interaction, which can be followed by a variety of processes ranging from sophisticated chemical transformation of the ligand (e.g., enzyme catalysis) to simple release of the ligand in the presence of other cofactors (e.g., transport proteins). In this section, we will only consider the main characteristics of the protein–ligand binding process.

Binding has been traditionally considered within a framework of the structure–function paradigm, a cornerstone of molecular biology for many years. It was > 100 years ago that Fischer coined the term *lock-and-key* to emphasize the requirement for a stereochemical fit between

an enzyme and its substrate in order for binding to occur (92). The limitations of this view of the binding process became obvious in the middle of the twentieth century, when a large body of newly acquired information on enzyme kinetics appeared to be in conflict with the notion that “*the enzyme was a rather rigid negative of the substrate and that the substrate had to fit into this negative to react*” (93). The revision of the lock-and-key theory by Koshland (92) led to a rise of the so-called *induced fit* theory, whose major premise was that the “*reaction between the enzyme and substrate can occur only after a change in protein structure induced by the substrate itself*”. Conformational changes occurring during enzyme catalysis are relatively small scale and affect mostly the catalytic site (94). Similarly, the conformational changes occurring in other proteins as a result of induced fit-type binding usually affect a limited fraction of the protein structure. An example of this is ferric ion binding by the iron-transport protein transferrin, an event that results in the repositioning of two protein domains within each lobe of the protein (Fig. 1.8). Although the overall effect of such repositioning is quite significant (and results in closing the cleft between the two domains), the number of affected amino acid residues does not exceed a dozen (95). More recently, numerous examples of large-scale conformational changes induced by ligand binding have been reported. The most extreme case is represented by the so-called *intrinsically disordered* proteins, which actually lack stable structure under native conditions in the absence of the ligand (96).

The above considerations strongly suggest that structure is not the sole determinant of protein function. As elegantly put by Onuchic and Wolynes (80), “*the twentieth century’s fixation on structure catapulted folding to center stage in molecular biology. The lessons learned about folding may, in the future, increase our understanding of many functional motions and large-scale assembly processes*”. In Section 1.5.2, we will consider various aspects of protein dynamics under native conditions that may be important modulators or even determinants of function.

1.5.2. Protein Dynamics Under Native Conditions

With very few exceptions, protein structure under native conditions is not a rigid crystalline state, but undergoes local breathing motions, involving anything from side-chain rotation to rearrangement of secondary structure elements relative to each other. Although the existence of such motions within the native state of the protein can be detected with a variety of experimental techniques, their exact nature remains the subject of discussion in the literature. The commonly accepted models of local dynamics within natively folded proteins invoke the notions of *structural fluctuation* (localized transient unfolding affecting only few atoms within the protein) (97) or a *mobile defect*, which

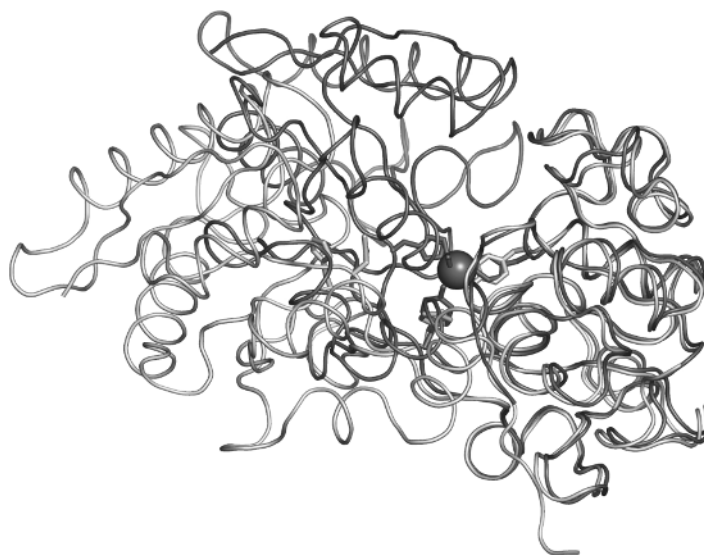


Figure 1.8. Superimposed crystal structures of the apo and holo forms of the *N*-lobe of human serum transferrin.

considers not only the emergence and dissipation of local disorder, but also the possibility of its propagation through the protein structure. Although the latter model has not enjoyed as much attention as the former, thorough theoretical considerations suggest that local perturbations of the secondary structure may in fact propagate through certain elements of the secondary structure (e.g., α -helices) in the form of a soliton (98). Alternatively, local dynamics can be described with a *solvent penetration* model (slow diffusion of the solvent molecules into and out of the protein interior) (99). Such description is actually very similar to the mobile defect model, as applied to the integral solute–solvent system, instead of the protein molecule alone.

Above and beyond local structural fluctuations, the dynamics of proteins under native conditions is exemplified by transiently sampling alternative (higher energy or “activated”) conformations. Such activated (non-native) states are often functionally important despite their low Boltzmann weight (100,101). An example of such behavior can be seen in cellular retinoic acid binding protein I (CRABP I), which sequesters and transports insoluble all-*trans* retinoic acid (RA) in the cytosol. The structures of the apo and the holo forms of this protein are very similar, consistent with the lock-and-key type of binding. However, the native structure of CRABP I provides no clue as to how the ligand gets access to the internal protein cavity, which is its binding site (Fig. 1.9). Obviously, in order to provide entrance into the cavity, a fraction of the native structure has to be lost transiently, an event consistent with the notion of sampling an activated protein state. Realization of the importance of transient non-native protein structures for their function has not only greatly advanced our understanding of processes as diverse as

recognition, signaling, and transport, but also has had profound practical implications, particularly for the design of drugs targeting specific proteins (102).

Many proteins use dynamics as a means of communication between different domains. This process, by which a signal, such as a binding event in one domain triggers a conformational change in another domain, is known as *allostery*. The paradigm for this effect is hemoglobin, a tetrameric protein mentioned earlier in this chapter. Binding a molecule of oxygen at the heme site of the α -chain induces a change in the oxygen affinity of the β -chain binding site by rearrangement of interdomain interactions (103,104). Another example is the chaperone protein DnaK, which assists in preventing the misfolding of nascent chains as they emerge from the ribosome. This 70-kDa protein consists of an ATPase domain joined via a short linker region to a peptide-binding domain. Binding of adenosine triphosphate (ATP) causes a conformational change in the peptide-binding domain that increases its affinity for substrate. Subsequent hydrolysis of nucleotide in the ATPase domain signals a conformational change in the adjacent domain that releases the unfolded polypeptide and allows it to begin to refold. The exact mechanism by which this allosteric communication occurs is still poorly understood. It is clear, however, that it must involve dynamic events at the interdomain interface that transmits the signal between the two binding sites (105).

1.5.3. Is Well-Defined Structure Required for Functional Competence?

A very interesting class of proteins that came to prominence in the early 2000s relies on dynamics even more heavily

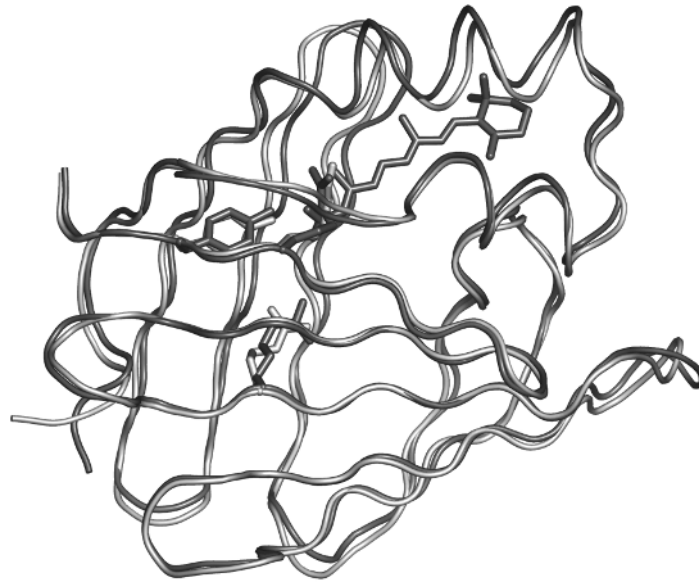


Figure 1.9. Overlaid crystal structures of the apo- and holo- forms of cellular retinoic acid binding protein I.

compared to the examples considered in Section 1.5.2. *Intrinsically disordered* proteins are remarkable in that they appear to have very little stable folded structure in isolation, contrary to our classical view of proteins as folded species. Several hundred proteins have now been identified that contain large segments of disorder even under native conditions (106), and many of these proteins seem to serve a wide variety of functions *in vivo* (107–109). A number of these are involved in activation or inhibition of transcription or translation, and while these proteins appear unstructured under native conditions they undergo a structural transition in the presence of their cognate substrate, whether this be another protein or a recognition site on a molecule of DNA or RNA (96,110). One important aspect of intrinsic disorder may be the necessity for this class of proteins to recognize and bind to multiple sites. Whereas a highly structured protein may only have a limited and very specific binding site available to it, one that has a highly dynamic structure should be able to adapt to a variety of different structural motifs. This intrinsic disorder phenomenon also seems contrary to the paradigm that an unstructured protein should be targeted for proteolysis, or else degradation by the proteasome. It seems this class of proteins manages to avoid such scenarios either by having regions that are sterically inaccessible or else do not contain residues that are sensitive to proteases. Indeed one other observation is that many intrinsically disordered proteins have a relatively short lifetime in the cell: They are expressed as needed in response to a signal and then rapidly removed by degradation. This would provide an efficient mechanism to switch on or off a cellular process for only a short period of time. A larger number of proteins in the eukaryotic genome have

been predicted to have disordered regions compared to prokaryotes, perhaps indicating the need for higher organisms to adjust more rapidly to environmental changes.

1.5.4. Biomolecular Dynamics and Binding from the Energy Landscape Perspective

The development of the folding funnel concept also has far-reaching consequences for our understanding of how proteins interact with each other and with other ligands. One theory offered as a general scheme of protein folding and binding implies that the only difference between the two processes is chain connectivity (111,112), namely, that monomeric protein folding represents an energy funnel for a single chain, whereas protein–protein association and peptide binding is a similar landscape, but with discontinuous backbone connections. In the more general case, however, the concept can be extended to encompass the chemical nature of the ligand and the energetics of the binding process, whether it be a noncovalent interaction or a chemical process as in the case of enzymatic catalysis; the energy funnel concept can be applied theoretically to describe the process by which a protein recognizes and binds to another molecule. Rigid proteins that bind ligands via a lock-and-key type mechanism presumably do not require significant dynamic events, so they will have few local minima similar in energy to the native state. In contrast, those proteins that utilize an induced fit binding mechanism may have a rugged energy surface characterized by a number of local minimum conformational states at the bottom of the folding–binding funnel (111).

The idea of a binding funnel has also been demonstrated computationally by Zhang et al. (113) to explain the fast

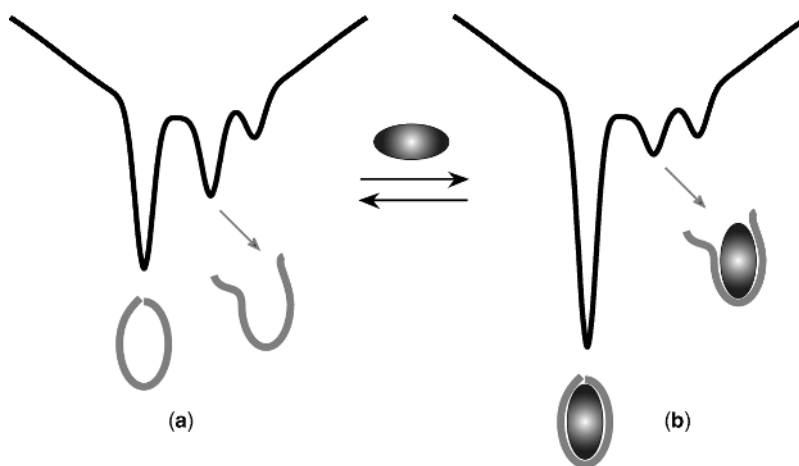


Figure 1.10. Schematic representations of the energy landscapes for the apo (*a*) and holo (*b*) forms of a protein whose ligand-binding behavior conforms to the lock-and-key type interaction.

protein–ligand association rates exhibited by many proteins. In this model, the initial collision event is accompanied by favorable interactions to form a long-lived encounter complex that significantly limits the search process to the ligand-bound conformation. The funnel energy landscape of protein binding may be a common feature in protein–protein associations (114,115). Knowledge of the relative energies and structural features of the local conformational minima available to proteins is clearly key to the understanding of what makes proteins efficient in binding their physiological ligands. Likewise, an extension of the protein folding problem is to understand how protein monomers assemble as functional multimers or other macromolecular assemblies. In the cellular environment, this must be an efficient process aided by dynamics and specific recognition events, all of which may be described in terms of the energetics of accessible conformational space.

All visibly different modes of binding (lock-and-key, induced fit, and binding of intrinsically disordered proteins) appear to have only quantitative differences within the framework provided by the binding funnel concept. This point is illustrated in Figures 1.10–1.12, which represent hypothetical folding funnels for proteins of each class in the absence and presence of their respective ligands. A protein whose ligand-binding behavior conforms to the lock-and-key type interaction (e.g., CRABP I considered earlier) is suggested to have an activated state whose structural features increase the rate of ligand entry into the binding site (Fig. 1.10). The protein molecules sample this activated state relatively frequently due to its relatively low energy. Once the ligand enters the binding site, its interaction with the protein increases the stability of the native conformation. Although “visitations” to the activated state are still possible, they do not occur as frequently due to

the increased energy difference between the two states. As a result, the protein can acquire the ligand relatively easily, but does not release it unless the energy landscape is altered again, (e.g., by a competing receptor of the ligand).

A similar analysis can be carried out for proteins conforming to the induced-fit type behavior (we will use the N-lobe of human serum transferrin as an example). Although the X-ray data suggest the existence of two distinct conformations of the protein depending on the presence of the ligand (open conformation for the apo form and closed for the holo form of the protein), there is experimental evidence suggesting that both conformations coexist in solution in equilibrium (116). The open conformation is, of course, favored in the absence of the ligand, while iron binding shifts the equilibrium toward the closed state (Fig. 1.11). Such a shift is qualitatively similar to the one considered for the lock-and-key interaction. The only difference is that the protein state corresponding to the global energy minimum in the absence of the ligand becomes “downgraded” to the status of an activated state (local energy minimum) as a result of the ligand binding.

Finally, folding of an intrinsically unstructured protein in the process of ligand binding can be viewed as a preferential stabilization of one particular conformation among many available to the protein in the ligand-free form (Fig. 1.12). An example of such behavior is presented by the β -chain of mammalian hemoglobins, which populate at least four different states (only one of them appears to be very close to the compact natively folded conformation) in solution in the absence of its binding partner, α -globin (117). Again, the general features of binding in this scheme appear to be very similar to those seen in the previous two examples (Figs. 1.10 and 1.11), the major distinct feature being the absence of the preferred conformation in the absence of the ligand.

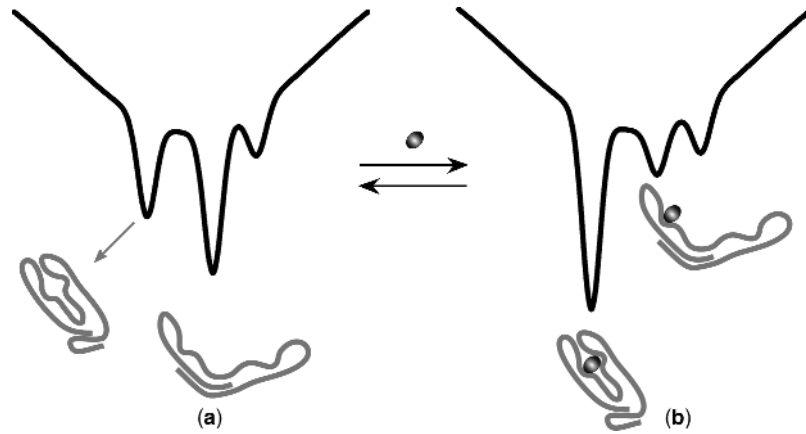


Figure 1.11. Schematic representations of the energy landscapes for the apo (a) and holo (b) forms of a protein whose ligand-binding behavior conforms to the induced fit type interaction.

1.5.5. Energy Landscapes Within a Broader Context of Nonlinear Dynamics: Information Flow and Fitness Landscapes

It becomes increasingly clear that the significance of the concept of energy landscapes extends well beyond the fields of protein folding and even molecular biophysics. Huang and Ingber (118) questioned the validity of the commonly accepted paradigm of cell regulation as a collection of pathways that link receptors with genes by asking “Can identification of all these signaling proteins and their assignment into distinct functional pathways lead to the full understanding of developmental control and cell fate regulation?” The existence of distributed information within cellular signaling, as well as the fact that a single signaling molecule may activate several genes and even produce opposite effects depending on its microenvironment, led to a suggestion that the concept of *linear signaling pathways* is inappropriate. The suggestion that the signaling and regulatory pathways are not simple linear connections

between receptors and genes is similar to the earlier realization of the limits of the classic concept of protein folding pathways. The paradox of signaling nonlinearity is resolved by the introduction of a concept of *cellular states*, and the switches between these states are viewed as *biological phase transitions* (118). In this new view, cell fates are viewed as common end programs or attractors within the entire regulatory network, which can be visualized as a *potential landscape* with multiple minima. Each minimum corresponds to a certain fate of the cell (e.g., differentiation, proliferation, apoptosis). A similar idea was recently proposed as the basis of a quantitative model of carcinogenesis, in which normal cells *in vivo* occupy a ridge-shaped maximum in a well-defined *tissue fitness landscape*, a configuration that allows cooperative coexistence of multiple cellular populations (119).

Finally, *dynamic fitness landscapes* have proven to be a valuable concept in evolutionary biology, molecular evolution (120,121), combinatorial optimization, and the physics

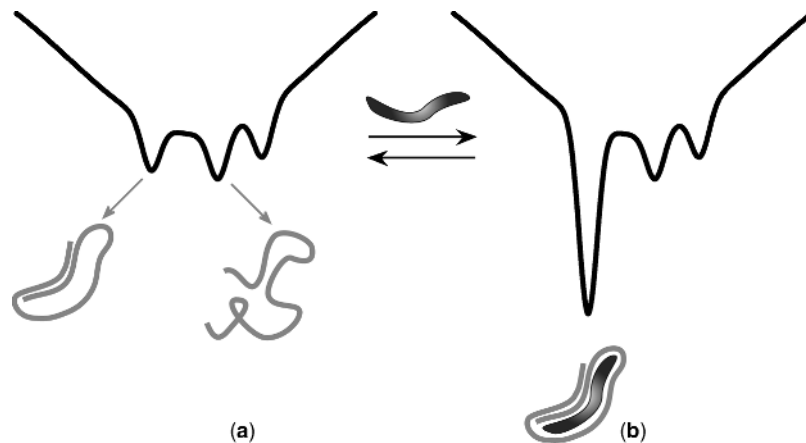


Figure 1.12. Schematic representations of the energy landscapes for the apo (a) and holo (b) forms of an intrinsically unstructured protein whose folding is induced by the ligand.

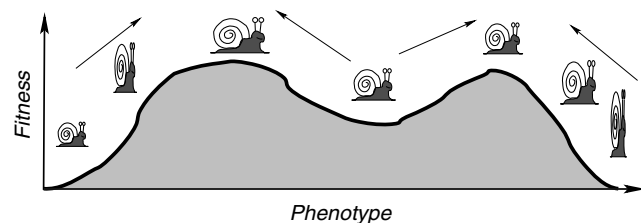


Figure 1.13. Schematic representation of a fitness landscape, where the fitness of a snail species depends on its shape. Mutations continuously produce variants that are selected if their fitness is larger than the fitness of the current “wild-type” snail. As a consequence, the shape of the snails changes over time until it reaches a maximum of the fitness landscape. [Reprinted with permission from (122). Copyright © 2002 Society for Industrial and Applied Mathematics. All rights reserved.]

of disordered systems (122). In evolutionary biology, this concept is often used to visualize the relationship between genotypes (or phenotypes) and replication success, an idea initially put forward by Wright (123). An evolving population typically climbs uphill in the fitness landscape, until it reaches a local optimum (Fig. 1.13), where it then remains, unless a rare mutation opens a path to a new fitness peak.

1.6. PROTEIN HIGHER ORDER STRUCTURE AND DYNAMICS FROM A BIOTECHNOLOGY PERSPECTIVE

In the first edition of this book, the discussion of protein folding and conformation was focused primarily on fundamental aspects, and the goal was to understand some general principles of protein (and, more broadly, biopolymer) behavior *in vitro* with the hope of being able to also apply this knowledge to *in vivo* situations. Spectacular progress has been made in the field of biotechnology in the past several years, which has resulted in a dramatic explosion of both the number of protein-based drugs and the range of the diseases they can treat (124,125), bringing to the forefront another very important aspect of protein folding and dynamics. Indeed, the large size of protein therapeutics (from several kDa to nearly 1 MDa, well beyond the molecular weight range of classical small molecule medicinal drugs) leads to an important distinction between small molecule drugs (where covalent structure is the sole determinant of the 3D structure and, ultimately, the therapeutic properties of the drug) and protein pharmaceuticals (where the large physical size makes a multitude of noncovalent contacts not only inevitable but, in fact, the defining element of their 3D structure). Correct folding is vital not only for the ability of a protein to execute its biological function, but also for many other aspects of its behavior (126). Failure to fold or maintain the native conformation obviously has a negative impact on the efficacy of the protein drug, since the recognition of a range of physiological targets requires that

the native conformation be maintained throughout the life-cycle of a protein molecule. Proteins that are not folded properly are prone to aggregation both *in vitro* and *in vivo*, and are targeted by proteases both inside and outside the cell, which obviously impacts bioavailability of the protein drug. Furthermore, misfolding and aggregation may trigger an immune response, thereby adversely affecting the safety profile of the protein drug. Critical dependence of the protein drug’s potency, stability, and safety on conformation makes its characterization an essential element throughout the drug development process from design to manufacturing to postapproval monitoring. Throughout this book, we will provide several examples of how MS can be used to probe various aspects of conformation, dynamics, and stability of protein pharmaceuticals, as well as their interactions with physiological partners and therapeutic targets.

REFERENCES

1. Thieffry, D. and Sarkar, S. (1998) Forty years under the central dogma, *Trends Biochem. Sci.* 23, 312–316.
2. Nirenberg, M. (2004) Historical review: Deciphering the genetic code—a personal account, *Trends Biochem. Sci.* 29, 46–54.
3. Hatfield, D. L. and Gladyshev, V. N. (2002) How selenium has altered our understanding of the genetic code, *Mol. Cell. Biol.* 22, 3565–3576.
4. Ibba, M. and Soll, D. (2002) Genetic code: Introducing pyrrolysine, *Curr. Biol.* 12, R464–R466.
5. Budisa, N., Minks, C., Alefelder, S., Wenger, W., Dong, F., Moroder, L., and Huber, R. (1999) Toward the experimental codon reassignment *in vivo*: protein building with an expanded amino acid repertoire, *FASEB J.* 13, 41–51.
6. Chin, J. W., Cropp, T. A., Anderson, J. C., Mukherji, M., Zhang, Z., and Schultz, P. G. (2003) An expanded eukaryotic genetic code, *Science* 301, 964–967.
7. Avetisov, V., and Goldanskii, V. (1996) Mirror symmetry breaking at the molecular level, *Proc. Natl. Acad. Sci. USA* 93, 11435–11442.
8. Podlech, J. (2001) Origin of organic molecules and biomolecular homochirality, *Cell. Mol. Life Sci.* 58, 44–60.
9. Warshel, A. and Papazyan, A. (1998) Electrostatic effects in macromolecules: fundamental concepts and practical modeling, *Curr. Opin. Struct. Biol.* 8, 211–217.
10. Schutz, C. N. and Warshel, A. (2001) What are the dielectric “constants” of proteins and how to validate electrostatic models?, *Proteins* 44, 400–417.
11. Harvey, S. C. (1989) Treatment of electrostatic effects in macromolecular modeling, *Proteins* 5, 78–92.
12. Daune, M. (1999) *Molecular biophysics: structures in motion*, Oxford University Press, Oxford, New York.
13. Mirsky, A. E. and Pauling, L. (1936) On the structure of native, denatured, and coagulated proteins, *Proc. Natl. Acad. Sci. USA* 22, 439–447.

14. Pauling, L., Corey, R. B., and Branson, H. R. (1951) The structure of proteins—2 hydrogen-bonded helical configurations of the polypeptide chain, *Proc. Natl. Acad. Sci. USA* 37, 205–211.
15. Ramachandran, G. N., Ramakrishnan, C., and Sasisekharan, V. (1963) Stereochemistry of polypeptide chain configurations, *J. Mol. Biol.* 7, 95–99.
16. Tanford, C. (1978) Hydrophobic effect and organization of living matter, *Science* 200, 1012–1018.
17. Scheraga, H. A. (1998) Theory of hydrophobic interactions, *J. Biomol. Struct. Dyn.* 16, 447–460.
18. Southall, N. T., Dill, K. A., and Haymet, A. D. J. (2002) A view of the hydrophobic effect, *J. Phys. Chem. B* 106, 521–533.
19. Widom, B., Bhimalapuram, P., and Koga, K. (2003) The hydrophobic effect, *Phys. Chem. Chem. Phys.* 5, 3085–3093.
20. Hartley, G. S. (1936) *Aqueous solutions of paraffin-chain salts; a study in micelle formation*, Hermann & Cie, Paris.
21. Tanford, C. (1979) Interfacial free-energy and the hydrophobic effect, *Proc. Natl. Acad. Sci. USA* 76, 4175–4176.
22. Hildebrandt, J. H., Nemethy, G., Scheraga, H. A., and Kauzmann, W. (1968) A criticism of term hydrophobic bond, *J. Phys. Chem.* 72, 1841–1842.
23. Frank, H. S. and Evans, M. W. (1945) Free volume and entropy in condensed systems. III. Entropy in binary liquid mixtures; partial molar entropy in dilute solutions; structure and thermodynamics in aqueous electrolytes, *J. Chem. Phys.* 13, 507–532.
24. Tanford, C. (1997) How protein chemists learned about the hydrophobic factor, *Protein Sci.* 6, 1358–1366.
25. Anfinsen, C. B. (1973) Principles that govern the folding of protein chains, *Science* 181, 223–230.
26. Chothia, C., Hubbard, T., Brenner, S., Barns, H., and Murzin, A. (1997) Protein folds in all- α and all- β classes, *Annu. Rev. Biophys. Biomol. Struct.* 26, 597–627.
27. Zhang, C. and DeLisi, C. (1998) Estimating the number of protein folds, *J. Mol. Biol.* 284, 1301–1305.
28. Govindarajan, S., Recabarren, R., and Goldstein, R. A. (1999) Estimating the total number of protein folds, *Proteins* 35, 408–414.
29. Hou, J., Sims, G. E., Zhang, C., and Kim, S.-H. (2003) A global representation of the protein fold space, *Proc. Natl. Acad. Sci. USA* 100, 2386–2390.
30. Caetano-Anolles, G. and Caetano-Anolles, D. (2003) An evolutionarily structured universe of protein architecture, *Genome Res.* 13, 1563–1571.
31. Rost, B. (1997) Protein structures sustain evolutionary drift, *Fold. Des.* 2, S19–S24.
32. Wood, T. C. and Pearson, W. R. (1999) Evolution of protein sequences and structures, *J. Mol. Biol.* 291, 977–995.
33. Zhang, C. and DeLisi, C. (2001) Protein folds: molecular systematics in three dimensions, *Cell. Mol. Life Sci.* 58, 72–79.
34. Denton, M. J., Marshall, C. J., and Legge, M. (2002) The protein folds as platonic forms: New support for the pre-Darwinian conception of evolution by natural law, *J. Theor. Biol.* 219, 325–342.
35. Epstein, C. J., Goldberger, R. F., and Anfinsen, C. B. (1963) Genetic control of tertiary protein structure—studies with model systems, *Cold Spring Harbor Symp. Quant. Biol.* 28, 439–439.
36. Dawkins, R. (1996) *The blind watchmaker: why the evidence of evolution reveals a universe without design*, Norton, New York.
37. Pennock, R. T. (2001) *Intelligent design creationism and its critics: philosophical, theological, and scientific perspectives*, MIT Press Cambridge, MA.
38. Levinthal, C. (1968) Are there pathways for protein folding, *J. Chim. Phys. Phys. Chim. Biol.* 65, 44–45.
39. Dill, K. (1999) Polymer principles and protein folding, *Protein Sci.* 8, 1166–1180.
40. Radford, S. E. and Dobson, C. M. (1999) From computer simulations to human disease: emerging themes in protein folding., *Cell* 97, 291–298.
41. Jaenicke, R. (1998) Protein self-organization in vitro and in vivo: partitioning between physical biochemistry and cell biology., *Biol. Chem.* 379, 237–243.
42. Jaenicke, R. (1995) Folding and association versus misfolding and aggregation of proteins., *Philos. Trans. R. Soc. London Ser. B* 348, 97–105.
43. Rochet, J.-C. and Lansbury, P. T. (2000) Amyloid fibrillogenesis: themes and variations, *Curr. Opin. Struct. Biol.* 10, 60–68.
44. Uversky, V. N. (2003) Protein folding revisited. A polypeptide chain at the folding-misfolding-nonfolding cross-roads: which way to go?, *Cell. Mol. Life Sci.* 60, 1852–1871.
45. Dobson, C. M. (2004) Principles of protein folding, misfolding and aggregation, *Semin. Cell Dev. Biol.* 15, 3–16.
46. Carrell, R. W. and Lomas, D. A. (1997) Conformational disease, *Lancet* 350, 134–138.
47. Bellotti, V., Mangione, P., and Stoppini, M. (1999) Biological activity and pathological implications of misfolded proteins, *Cell. Mol. Life Sci.* 55, 977–991.
48. Thompson, A. J. and Barrow, C. J. (2002) Protein conformational misfolding and amyloid formation: characteristics of a new class of disorders that include Alzheimer's and Prion diseases, *Curr. Med. Chem.* 9, 1751–1762.
49. Bucciantini, M., Giannoni, E., Chiti, F., Baroni, F., Formigli, L., Zurdo, J., Taddei, N., Ramponi, G., Dobson, C. M., and Stefani, M. (2002) Inherent toxicity of aggregates implies a common mechanism for protein misfolding diseases, *Nature (London)*, 416, 507–511.
50. Caughey, B. and Lansbury, P. T. (2003) Protofibrils, pores, fibrils, and neurodegeneration: Separating the responsible protein aggregates from the innocent bystanders, *Annu. Rev. Neurosci.* 26, 267–298.
51. Wales, D. J. and Scheraga, H. A. (1999) Global optimization of clusters, crystals, and biomolecules, *Science* 285, 1368–1372.
52. Schonbrun, J., Wedemeyer, W. J., and Baker, D. (2002) Protein structure prediction in 2002, *Curr. Opin. Struct. Biol.* 12, 348–354.

53. Whisstock, J. C. and Lesk, A. M. (2003) Prediction of protein function from protein sequence and structure, *Q. Rev. Biophys.* 36, 307–340.
54. Price, N. C. (2000) Conformational issues in the characterization of proteins, *Biotechnol. Appl. Biochem.* 31, 29–40.
55. Smith, L. J., Fiebig, K. M., Schwalbe, H., and Dobson, C. M. (1996) The concept of a random coil. Residual structure in peptides and denatured proteins, *Fold. Des.* 1, R95–R106.
56. Ptitsyn, O. B. (1995) Molten globule and protein folding, *Adv. Protein Chem.* 47, 83–229.
57. Uverskii, V. N. (1998) [How many molten globules states exist?], *Biofizika* 43, 416–421.
58. Clarke, J. and Fersht, A. R. (1996) An evaluation of the use of hydrogen exchange at equilibrium to probe intermediates on the protein folding pathway, *Fold. Des.* 1, 243–254.
59. Clarke, J., Itzhaki, L. S., and Fersht, A. R. (1997) Hydrogen exchange at equilibrium: a short cut for analysing protein-folding pathways?, *Trends Biochem. Sci.* 22, 284–287.
60. Jennings, P. A. and Wright, P. E. (1993) Formation of a molten globule intermediate early in the kinetic folding pathway of apomyoglobin, *Science* 262, 892–896.
61. Bai, Y., Sosnick, T. R., Mayne, L., and Englander, S. W. (1995) Protein folding intermediates: native state hydrogen exchange., *Science* 269, 192–196.
62. Bai, Y. (1999) Equilibrium amide hydrogen exchange and protein folding kinetics, *J. Biomol. NMR* 15, 65–70.
63. Parker, M. J. and Marqusee, S. (2001) A kinetic folding intermediate probed by native state hydrogen exchange, *J. Mol. Biol.* 305, 593–602.
64. Creighton, T. E. (1984) Pathways and mechanisms of protein folding, *Adv. Biophys.* 18, 1–20.
65. Creighton, T. E. (1988) Toward a better understanding of protein folding pathways, *Proc. Natl. Acad. Sci. USA* 85, 5082–5086.
66. Gunasekaran, K., Eyles, S. J., Hagler, A. T., and Gierasch, L. M. (2001) Keeping it in the family: folding studies of related proteins, *Curr. Opin. Struct. Biol.* 11, 83–93.
67. Matthews, C. R. (1993) Pathways of protein folding, *Annu. Rev. Biochem.* 62, 653–683.
68. Daggett, V. and Fersht, A. R. (2003) Is there a unifying mechanism for protein folding?, *Trends Biochem. Sci.* 28, 18–25.
69. Baldwin, R. L. (1995) The nature of protein folding pathways: the classical versus the new view., *J. Biomolec. NMR* 5, 103–109.
70. Dill, K. A. and Chan, H. S. (1997) From Levinthal to pathways to funnels., *Nat. Struct. Biol.* 4, 10–19.
71. Pande, V. S., Grosberg, A., Tanaka, T., and Rokhsar, D. S. (1998) Pathways for protein folding: is a new view needed?, *Curr. Opin. Struct. Biol.* 8, 68–79.
72. Bryngelson, J. D., Onuchic, J. N., Socci, N. D., and Wolynes, P. G. (1995) Funnels, pathways and the energy landscape of protein folding: a synthesis., *Proteins* 21, 167–195.
73. Onuchic, J. N., Luthey-Schulten, Z., and Wolynes, P. G. (1997) Theory of protein folding: the energy landscape perspective, *Annu. Rev. Phys. Chem.* 48, 545–600.
74. Brooks, C. L., 3rd, Gruebele, M., Onuchic, J. N., and Wolynes, P. G. (1998) Chemical physics of protein folding, *Proc. Natl. Acad. Sci. USA* 95, 11037–11038.
75. Plotkin, S. S. and Onuchic, J. N. (2002) Understanding protein folding with energy landscape theory. Part I: Basic concepts, *Q. Rev. Biophys.* 35, 111–167.
76. Gruebele, M. (2002) Protein folding: the free energy surface, *Curr. Opin. Struct. Biol.* 12, 161–168.
77. Dinner, A. R., Sali, A., Smith, L. J., Dobson, C. M. and Karplus, M. (2000) Understanding protein folding via free-energy surfaces from theory and experiment., *Trends Biochem. Sci.* 25, 331–339.
78. Brooks, C. L., 3rd, Onuchic, J. N., and Wales, D. J. (2001) Statistical thermodynamics. Taking a walk on a landscape, *Science* 293, 612–613.
79. Fersht, A. R. and Daggett, V. (2002) Protein folding and unfolding at atomic resolution, *Cell* 108, 573–582.
80. Onuchic, J. N. and Wolynes, P. G. (2004) Theory of protein folding, *Curr. Opin. Struct. Biol.* 14, 70–75.
81. Jackson, S. E. (1998) How do small single-domain proteins fold?, *Fold. Des.* 3, R81–91.
82. Myers, J. K. and Oas, T. G. (2002) Mechanisms of fast protein folding, *Annu. Rev. Biochem.* 71, 783–815.
83. Ptitsyn, O. B. (1994) Kinetic and equilibrium intermediates in protein folding., *Protein Eng.* 7, 593–596.
84. Eliezer, D., Jennings, P. A., Dyson, H. J., and Wright, P. E. (1997) Populating the equilibrium molten globule state of apomyoglobin under conditions suitable for structural characterization by NMR, *FEBS Lett.* 417, 92–96.
85. Brutscher, B., Bruschweiler, R., and Ernst, R. R. (1997) Backbone dynamics and structural characterization of the partially folded A state of ubiquitin by ¹H, ¹³C, and ¹⁵N nuclear magnetic resonance spectroscopy, *Biochemistry* 36, 13043–13053.
86. Baum, J., Dobson, C. M., Evans, P. A., and Hanley, C. (1989) Characterization of a partly folded protein by NMR methods: studies on the molten globule state of guinea pig alpha-lactalbumin, *Biochemistry* 28, 7–13.
87. Bai, Y., Englander, J. J., Mayne, L., Milne, J. S., and Englander, S. W. (1995) Thermodynamic parameters from hydrogen exchange measurements., *Methods Enzymol.* 259, 344–356.
88. Fink, A. L. (1999) Chaperone-mediated protein folding, *Physiol. Rev.* 79, 425–449.
89. Hartl, F. U. and Hayer-Hartl, M. (2002) Molecular chaperones in the cytosol: from nascent chain to folded protein, *Science* 295, 1852–1858.
90. Trombetta, E. S. and Parodi, A. J. (2003) Quality control and protein folding in the secretory pathway, *Annu. Rev. Cell Dev. Biol.* 19, 649–676.
91. Mogk, A. and Bukau, B. (2004) Molecular chaperones: structure of a protein disaggregase, *Curr. Biol.* 14, R78–R80.
92. Fischer, E. (1894) Einfluss der configuration auf die wirkung derenzyme, *Ber. Dt. Chem. Ges.* 27, 2985–2993.

93. Koshland, D. E. (1958) Application of a theory of enzyme specificity to protein synthesis, *Proc. Natl. Acad. Sci. USA* 44, 98–104.
94. Koshland, D. E., Jr. (1998) Conformational changes: how small is big enough?, *Nat. Med.* 4, 1112–1114.
95. Jeffrey, P. D., Bewley, M. C., MacGillivray, R. T., Mason, A. B., Woodworth, R. C., and Baker, E. N. (1998) Ligand-induced conformational change in transferrins: crystal structure of the open form of the N-terminal half-molecule of human transferrin, *Biochemistry* 37, 13978–13986.
96. Dunker, A. K., Brown, C. J., Lawson, J. D., Iakoucheva, L. M., and Obradovic, Z. (2002) Intrinsic disorder and protein function, *Biochemistry* 41, 6573–6582.
97. Maity, H., Lim, W. K., Rumbley, J. N., and Englander, S. W. (2003) Protein hydrogen exchange mechanism: local fluctuations, *Protein Sci.* 12, 153–160.
98. d’Ovidio, F., Bohr, H. G., and Lindgard, P. A. (2003) Solitons on H bonds in proteins, *J. Phys. Condens. Mat.* 15, S1699–S1707.
99. Miller, D. W. and Dill, K. A. (1995) A statistical mechanical model for hydrogen exchange in globular proteins, *Protein Sci.* 4, 1860–1873.
100. Tsai, C. D., Ma, B., Kumar, S., Wolfson, H., and Nussinov, R. (2001) Protein folding: binding of conformationally fluctuating building blocks via population selection, *Crit. Rev. Biochem. Mol. Biol.* 36, 399–433.
101. Papoian, G. A. and Wolynes, P. G. (2003) The physics and bioinformatics of binding and folding—an energy landscape perspective, *Biopolymers* 68, 333–349.
102. Carlson, H. A. (2002) Protein flexibility and drug design: how to hit a moving target, *Curr. Opin. Chem. Biol.* 6, 447–452.
103. Perutz, M. F. (1972) Stereochemical mechanism of cooperative effects in haemoglobin, *Biochimie* 54, 587–588.
104. Perutz, M. F., Wilkinson, A. J., Paoli, M., and Dodson, G. G. (1998) The stereochemical mechanism of the cooperative effects in hemoglobin revisited, *Annu. Rev. Biophys. Biomol. Struct.* 27, 1–34.
105. Pellicchia, M., Montgomery, D. L., Stevens, S. Y., Vander Kooi, C. W., Feng, H. P., Gierasch, L. M., and Zwietering, E. R. (2000) Structural insights into substrate binding by the molecular chaperone DnaK, *Nat. Struct. Biol.* 7, 298–303.
106. Le Gall, T., Romero, P. R., Cortese, M. S., Uversky, V. N., and Dunker, A. K. (2007) Intrinsic disorder in the Protein Data Bank, *J. Biomol. Struct. Dyn.* 24, 325–342.
107. Dunker, A. K., Cortese, M. S., Romero, P., Iakoucheva, L. M., and Uversky, V. N. (2005) Flexible nets. The roles of intrinsic disorder in protein interaction networks, *FEBS J.* 272, 5129–5148.
108. Radivojac, P., Iakoucheva, L. M., Oldfield, C. J., Obradovic, Z., Uversky, V. N., and Dunker, A. K. (2007) Intrinsic disorder and functional proteomics, *Biophys. J.* 92, 1439–1456.
109. Dunker, A. K., Oldfield, C., Meng, J., Romero, P., Yang, J., Chen, J., Vacic, V., Obradovic, Z., and Uversky, V. (2008) The unfoldomics decade: an update on intrinsically disordered proteins, *BMC Genomics* 9, S1.
110. Iakoucheva, L. M., Brown, C. J., Lawson, J. D., Obradovic, Z., and Dunker, A. K. (2002) Intrinsic disorder in cell-signaling and cancer-associated proteins, *J. Mol. Biol.* 323, 573–584.
111. Tsai, C.-J., Xu, D., and Nussinov, R. (1998) Protein folding via binding and vice versa., *Fold. Des.* 3, R71–R80.
112. Tsai, C.-J., Kumar, S., Ma, B., and Nussinov, R. (1999) Folding funnels, binding funnels, and protein function, *Protein Sci.* 8, 1181–1190.
113. Zhang, C., Chen, J., and DeLisi, C. (1999) Protein-protein recognition: exploring the energy funnels near the binding sites, *Proteins* 34, 255–267.
114. Tovchigrechko, A. and Vakser, I. A. (2001) How common is the funnel-like energy landscape in protein–protein interactions?, *Protein Sci.* 10, 1572–1583.
115. Boehr, D. D., Nussinov, R., and Wright, P. E. (2009) The role of dynamic conformational ensembles in biomolecular recognition, *Nat. Chem. Biol.* 5, 789–796.
116. Baker, H. M., Anderson, B. F., and Baker, E. N. (2003) Dealing with iron: Common structural principles in proteins that transport iron and heme, *Proc. Natl. Acad. Sci. USA* 100, 3579–3583.
117. Griffith, W. P. and Kaltashov, I. A. (2003) Highly asymmetric interactions between globin chains during hemoglobin assembly revealed by electrospray ionization mass spectrometry, *Biochemistry* 42, 10024–10033.
118. Huang, S. and Ingber, D. E. (2000) Shape-dependent control of cell growth, differentiation and apoptosis: switching between attractors in cell regulatory networks, *Exp. Cell Res.* 261, 91–103.
119. Gatenby, R. A. and Vincent, T. L. (2003) An evolutionary model of carcinogenesis, *Cancer Res.* 63, 6212–6220.
120. Drossel, B. (2001) Biological evolution and statistical physics, *Adv. Phys.* 50, 209–295.
121. Wilke, C. O., Ronnewinkel, C., and Martinetz, T. (2001) Dynamic fitness landscapes in molecular evolution, *Phys. Rep.* 349, 395–446.
122. Reidys, C. M. and Stadler, P. F. (2002) Combinatorial landscapes, *SIAM Rev.* 44, 3–54.
123. Wright, S. (1967) “Surfaces” of selective value, *Proc. Natl. Acad. Sci. USA* 58, 165–172.
124. Walsh, G. (2007) *Pharmaceutical biotechnology: concepts and applications*, John Wiley & Sons, Hoboken, NJ.
125. Walsh, G. (2010) Biopharmaceutical benchmarks 2010, *Nat. Biotech.* 28, 917–924.
126. Korotchkina, L. G., Ramani, K., and Balu-Iyer, S. V. (2008) Folding considerations for therapeutic protein formulations, In *Molecular Biology of Protein Folding, Part A*, pp. 255–270, Elsevier Academic Press Inc., San Diego.