

# Introduction to protein folding for physicists

Pablo Echenique\*

Theoretical Physics Department, University of Zaragoza,  
Pedro Cerbuna 12, 50009, Zaragoza, Spain.

Institute for Biocomputation and Physics of Complex Systems (BIFI),  
Edificio Cervantes, Corona de Aragón 42, 50009, Zaragoza, Spain.

February 1, 2008

## Abstract

The prediction of the three-dimensional native structure of proteins from the knowledge of their amino acid sequence, known as the *protein folding problem*, is one of the most important yet unsolved issues of modern science. Since the conformational behaviour of flexible molecules is nothing more than a complex physical problem, increasingly more physicists are moving into the study of protein systems, bringing with them powerful mathematical and computational tools, as well as the sharp intuition and deep images inherent to the physics discipline. This work attempts to facilitate the first steps of such a transition. In order to achieve this goal, we provide an exhaustive account of the reasons underlying the protein folding problem enormous relevance and summarize the present-day status of the methods aimed to solving it. We also provide an introduction to the particular structure of these biological heteropolymers, and we physically define the problem stating the assumptions behind this (commonly implicit) definition. Finally, we review the ‘special flavor’ of statistical mechanics that is typically used to study the astronomically large phase spaces of macromolecules. Throughout the whole work, much material that is found scattered in the literature has been put together here to improve comprehension and to serve as a handy reference.

## 1 Why study proteins?

Virtually every scientific book or article starts with a paragraph in which the writer tries to persuade the readers that the topic discussed is very important for the future of humankind. We shall stick to that tradition in this work; but with the confidence that, in the case of proteins, the persuasion process will turn out to be rather easy and automatic.

---

\*E-mail address: [pnique@unizar.es](mailto:pnique@unizar.es) — Web page: <http://www.pabloechenique.com>

Proteins are a particular type of biological molecules that can be found in every single living being on Earth. The characteristic that renders them essential for understanding life is simply their versatility. In contrast with the relatively limited structural variations present in other types of important biological molecules, such as carbohydrates, lipids or nucleic acids, proteins display a seemingly infinite capability for assuming different shapes and for producing very specific catalytic regions on their surface. As a result, proteins constitute the working force of the chemistry of living beings, performing almost every task that is complicated. Quoting the first sentence of a section (which shares this section's title) in Lesk's book [1]:

*In the drama of life on a molecular scale, proteins are where the action is.*

Just to state a few examples of what is meant by 'action', in living beings, proteins

- are passive building blocks of many biological structures, such as the coats of viruses, the cellular cytoskeleton, the epidermal keratin or the collagen in bones and cartilages;
- transport and store other species, from electrons to macromolecules;
- as hormones, transmit information and signals between cells and organs;
- as antibodies, defend the organism against intruders;
- are the essential components of muscles, converting chemical energy into mechanical one, and allowing the animals to move and interact with the environment;
- control the passage of species through the membranes of cells and organelles;
- control gene expression;
- are the essential agents in the transcription of the genetic information into more proteins;
- together with some nucleic acids, form the ribosome, the large molecular organelle where proteins themselves are synthesized;
- as chaperones, protect other proteins to help them to acquire their functional three-dimensional structure.

Due to this participation in almost every task that is essential for life, protein science constitutes a support of increasing importance for the development of modern Medicine. On one side, the lack or malfunction of particular proteins is behind many pathologies; e.g., in most types of cancer, mutations are found in

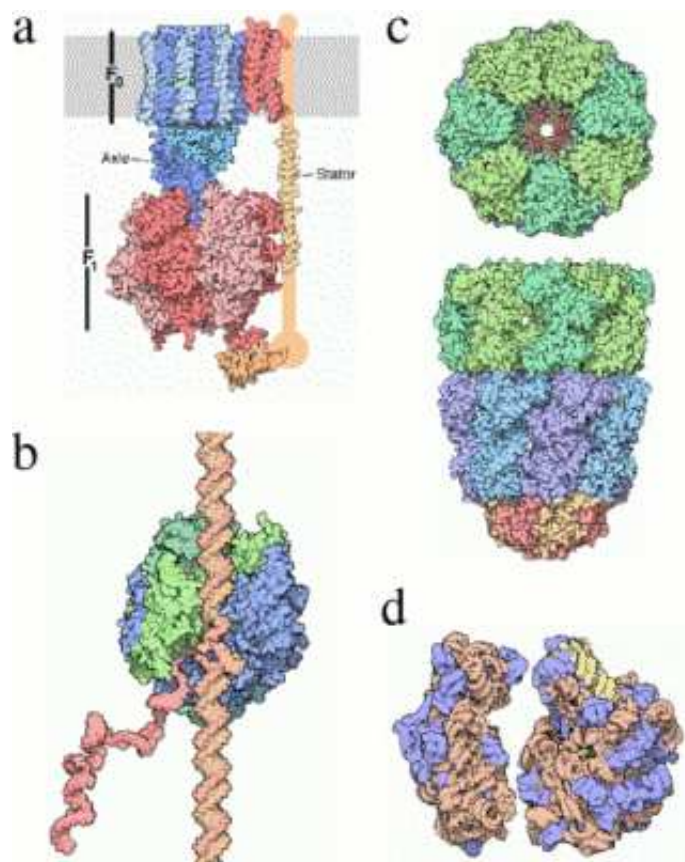


Figure 1.1: Four molecular machines formed principally by proteins. Figures taken from the *Molecule of the month* section of the RSCB Protein Data Bank (<http://www.pdb.org>), we thank the RSCB PDB and David S. Goodsell, from the Scripps Research Institute, for kind permission to use them. (a) *ATP synthase*: it acts as an energy generator when it is traversed by protons that make its two coupled engines rotate in reverse mode and the ATP molecule is produced. (b) *RNA polymerase*: it slides along a thread of DNA reading the base pairs and synthesizing a matching copy of RNA. (c) *GroEL-GroES complex*: it helps unfolded proteins to fold by sheltering them from the overcrowded cellular cytoplasm. (d) *Ribosome*: it polymerizes amino acids to form proteins following the instructions written in a thread of mRNA.

the tumor suppressor p53 protein [2]. Also, abnormal protein aggregation characterizes many neurodegenerative disorders, including Huntington, Alzheimer, Creutzfeld-Jakob (‘mad cow’), or motor neuron diseases [3–5]. Finally, to attack the vital proteins of pathogens (HIV [6, 7], SARS [8], hepatitis [9], etc.), or to block the synthesis of proteins at the bacterial ribosome [10], are common strategies to battle infections in the frenetic field of rational drug design [11].

Apart from Medicine, the rest of human technology may also benefit from the solutions that Nature, after thousands of millions<sup>1</sup> of years of ‘research’, has found to the typical practical problems. And that solutions are often proteins: New materials of extraordinary mechanical properties could be designed from the basis of the spider silk [12, 13], elastin [14] or collagen proteins [15]. Also, some attempts are being made to integrate these new biomaterials with living organic tissues and make them respond to stimuli [16]. Even further away on the road that goes from passive structural functions to active tasks, no engineer who has ever tried to solve a difficult chemical problem can avoid to experience a feeling of almost religious inferiority when faced to the speed, efficiency and specificity with which proteins cut, bend, repair, carry, link or modify other chemical species. Hence, it is normal that we play with the idea of learning to control that power and have, as a result, nanoengines, nanogenerators, nanoscissors, nanomachines in general [17]. The author of this work, in particular, felt a small sting of awe when he learnt about the pump and the two coupled engines of the principal energy generator in the cell, the *ATP synthase* (figure 1.1a); about the genetic Xerox machine, the *RNA polymerase* (figure 1.1b); about the hut where the proteins fold under shelter, the *GroEL-GroES* complex (figure 1.1c); or about the macromolecular factory where proteins are created, the *ribosome* (figure 1.1d), to mention four specially impressive examples. Agreeing again with Lesk [1]:

*Proteins are fascinating molecular devices.*

From a more academic standpoint, proteins are proving to be a powerful centre of interdisciplinary research, making many diverse fields and people with different formations come in contact<sup>2</sup>. Proteins force biologists, biochemists and chemists to learn more physics, mathematics and computation and force mathematicians, physicists and computer technicians to learn more biology, biochemistry and chemistry. This, indeed, cannot be negative.

In 2005, in a special section of Science magazine entitled ‘What don’t we know?’ [18], a selection of the hundred most interesting yet unanswered scientific questions was presented. What indicates the role of proteins, and particularly of the protein folding problem (treated in section 3), as focuses of interdisciplinary collaboration is not the inclusion of the question *Can we predict how proteins will fold?*, which was a must, but the large number of other questions which were

---

<sup>1</sup> Herein, we shall use the British convention for naming large numbers; in which  $10^9$ =‘a thousand million’,  $10^{12}$ =‘a billion’,  $10^{15}$ =‘a thousand billion’,  $10^{18}$ =‘a trillion’, and so on.

<sup>2</sup> The Institute for Biocomputation and Physics of Complex Systems, which the author is part of, constitutes an example of this rather new form of collaboration among scientists.

related to or even dependent on it, such as *Why do humans have so few genes?*, *How much can human life span be extended?*, *What is the structure of water?*, *How does a single somatic cell become a whole plant?*, *How many proteins are there in humans?*, *How do proteins find their partners?*, *How do prion diseases work?*, *How will big pictures emerge from a sea of biological data?*, *How far can we push chemical self-assembly?* or *Is an effective HIV vaccine feasible?*.

In this direction, probably the best example of the use that protein science makes of the existing human expertise, and of the positive feedback that this brings up in terms of new developments and resources, can be found in the machines that every one of us has on his/her desktops. In a first step, the enormous amount of biological data that emerges from the sequencing of the genomes of different living organisms requires computerized databases for its proper filtering. The NCBI GenBank database<sup>3</sup>, which is one of the most exhaustive repositories of sequenced genetic material, has doubled the number of deposited DNA bases approximately every 18 months since 1982 (see figure 1.2a) and has recently (in August 2005) exceeded the milestone of 100 Gigabases ( $10^{11}$ ) from over 165,000 species.

Among them, and according to the Entrez Genome Project database<sup>4</sup>, the sequencing of the complete genome of 366 organisms has been already achieved and there are 791 more to come in next few years. In the group of the completed ones, most are bacteria, and there are only two mammals: the poor laboratory mouse, *Mus Musculus*, and, notably [19], the *Homo Sapiens* (with  $\sim 3 \cdot 10^9$  bases and a mass-media-broadcast battle between the private firm Celera and the public consortium IHGSC).

However, not all the DNA encodes proteins (not all the DNA is genes). Typically, more than 95% of the genetic material in living beings is *junk DNA*, also called *non-coding DNA* (a more neutral term which seems recommendable in the light of some recent discoveries [20–22]). So, in a second step, the coding regions must be identified and each gene translated into the amino acid sequence of a particular protein<sup>5</sup>. The UniProt database<sup>6</sup> is, probably, the most comprehensive repository of these translated protein sequences and also of others coming from a variety of sources, including direct experimental determination [25, 26]. UniProt is comprised by two different sub-databases: the Swiss-Prot Protein Knowledgebase, which contains extensively human-annotated protein sequences with low redundancy; and TrEMBL, which contains computer-annotated sequences extracted directly from the underlying nucleotide entries at databases such as GenBank and where only the most basic redundancies have been removed.

The UniProt/Swiss-Prot database contains, at the moment (on 30 May

---

<sup>3</sup> <http://www.ncbi.nlm.nih.gov/Genbank/>

<sup>4</sup> <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=genomeprj>

<sup>5</sup> Note that many variations [23, 24] may occur before, during and after the process of gene expression, so that the relation gene-to-protein is not one-to-one. The size of the human proteome (the number of different proteins), for example, is estimated to be an order of magnitude or two larger than the size of the genome.

<sup>6</sup> <http://www.uniprot.org>

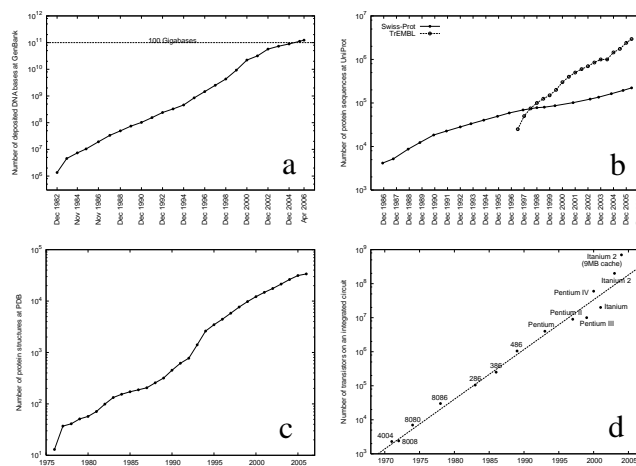


Figure 1.2: Recent exponential progress in genomics, proteomics and computer technology. **(a)** Evolution of the number of DNA bases deposited at the GenBank database. **(b)** Evolution of the number of protein sequences at the UniProt Swiss-Prot and TrEMBL databases. **(c)** Evolution of the number of protein three-dimensional structures at the Protein Data Bank. **(d)** Moore's Law: evolution of the number of transistors in the Intel CPUs.

2006), around 200,000 protein sequences from about 10,000 species, and it has experienced an exponential growth (since 1986), doubling the number of records approximately every 41 months (see figure 1.2b). In turn, the UniProt/TrEMBL database contains almost 3 million protein sequences from more than 100,000 species, and its growth (from 1997) has also been exponential, doubling the number of records approximately every 16 months (see figure 1.2b).

After knowing the sequence of a protein, the next step towards the understanding of biological processes is the characterization of its three-dimensional structure. Most proteins perform their function under a very specific *native* shape which involves many twists, loops and bends of the linear chain of amino acids (see section 3). This spatial structure is much more important than the sequence for biochemists to predict and understand the mechanisms of life and it can be resolved, nowadays, by fundamentally two experimental techniques: for small proteins, nuclear magnetic resonance (NMR) [27,28] and, more commonly, for proteins of any size, x-ray crystallography [29–31]. The three-dimensional structures so obtained are deposited in a centralized public-access database called Protein Data Bank (PDB)<sup>7</sup> [32]. From the 13 structures deposited in 1976 to the 33,782 (from more than a thousand species) stored in June 2006, the growth of the PDB has been (guess?) exponential, doubling the number of records approximately every 3 years (see figure 1.2c).

<sup>7</sup> <http://www.rcsb.org/pdb/>

To summarize, in June 2006, we have sequenced partial segments of the genetic material of around 160,000 species, having completed the genomes of only 366; we know the sequences of some of the proteins of around 100,000 species and the three-dimensional structure of proteins in 1,103 species<sup>8</sup>. However, according to the UN Millennium Ecosystem Assessment<sup>9</sup>, the number of species formally identified is 1.7-2 million and the estimated total number of species on Earth ranges from 5 million to 30 million [33]. Therefore, we should expect that the exponential growth of genomic and proteomic data will continue to fill the hard-disks, collapse the broadband connexions and heat the CPUs of our computers at least for the next pair of decades.

Fortunately, the improvement of silicon technology behaves in the same way: In fact, in 1965, Gordon Moore, co-founder of Intel, made the observation that the number of transistors per square inch had doubled every year since the integrated circuit was invented, and predicted that this exponential trend would continue for the foreseeable future. This has certainly happened (although the doubling time seems to be closer to 18 months) and this empirical law, which is not expected to fail in the near future, has become to be known as *Moore's Law* (see figure 1.2d for an example involving Intel processors). So we do not have to worry about running short of computational resources!

Of course, information produces more information, and public databases do not end at the three-dimensional structures of proteins. In the last few years, a number of more specific web-based repositories have been created in the field of molecular biology. There is the Protein Model Database (PMDb)<sup>10</sup> [34], where theoretical three-dimensional protein models are stored (including all models submitted to last four editions of the CASP<sup>11</sup> experiment [35]); the ProTherm<sup>12</sup> and ProNIT<sup>13</sup> databases [36], where a wealth of thermodynamical data is stored about protein stability and protein-nucleic acid interactions, respectively; the dbPTM<sup>14</sup> database [37], that stores information on protein post-translational modifications; the PINT<sup>15</sup> database [38], with thermodynamical data on protein-protein interactions; and so on and so forth.

In addition to the use of computers for storage and retrieval of enormous quantities of data, the increasing numerical power of these machines is customarily used for a wide variety of applications that range from molecular visualization, to long simulations aimed to solve the equations governing biological systems (the central topic discussed more in detail in the rest of this work).

Indeed, as Richard Dawkins has stated [39]:

*What is truly revolutionary about molecular biology in the post-Watson-Crick era is that it has become digital.*

---

<sup>8</sup> <http://www.ebi.ac.uk/thornton-srv/databases/pdbsum/>

<sup>9</sup> <http://www.millenniumassessment.org>

<sup>10</sup> <http://www.caspar.it/PMDB/>

<sup>11</sup> <http://predictioncenter.gc.ucdavis.edu>

<sup>12</sup> <http://gibk26.bse.kyutech.ac.jp/jouhou/Protherm/protherm.html>

<sup>13</sup> <http://gibk26.bse.kyutech.ac.jp/jouhou/pronit/pronit.html>

<sup>14</sup> <http://dbPTM.mbc.nctu.edu.tw>

<sup>15</sup> <http://www.bioinfodatabase.com/pint/>



Finally, apart from all the convincing reasons and the appeals to authority given above, what is crystal-clear is that proteins are an unsolved and difficult enigma. And those are two irresistible qualities for any flesh and blood scientist.

## 2 Summary of protein structure

In spite of their diverse biological functions, summarized in the previous section, proteins are a rather homogeneous class of molecules from the chemical point of view. They are *linear heteropolymers*, i.e., unbranched chains of different identifiable monomeric units.

Before they are assembled into proteins, these building units are called *amino acids* and can exist as stand alone stable molecules. All amino acids are made up of a central  $\alpha$ -carbon with four groups attached to it: an amino group ( $-\text{NH}_2$ ), a carboxyl group ( $-\text{COOH}$ ), a hydrogen atom and a fourth arbitrary group ( $-\text{R}$ ) (see figure 2.2). In aqueous solvent and under physiological conditions, both the amino and carboxyl groups are charged, the first accepting one proton and getting a positive charge, and the second giving one proton away and getting a negative charge (compare figures 2.2a and 2.2c).

When the group  $-\text{R}$  is not equal to one of the other three groups attached to the  $\alpha$ -carbon, the amino acid is *chiral*, i.e., like our hands, it may exist in two different forms, which are mirror images of one another and cannot be superimposed by rotating one of them in space (you cannot wear the left-hand glove on your right hand). In chemical jargon, one says that the  $\alpha$ -carbon constitutes an *asymmetric centre* and that the amino acid may exist as two different *enantiomers* called *L-* (figure 2.2c) and *D-* (figure 2.2d) forms. It is common that, when used as prefixes, the L and D letters, which come from *levorotatory* and *dextrorotatory*, are written in small capitals, as in L- and D-. This nomenclature is based on the possibility of associating the amino acids to the optically active L- and D- enantiomers of glyceraldehyde, and could be related to the  $+/-$  or to the Cahn, Ingold and Prelog's R/S [41] notations. For us, it suffices to say that

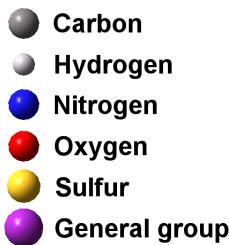


Figure 2.1: Color and size code for the atom types used in most of the figures in this section. All the figures have been made with the Gaussview graphical front-end of Gaussian03 [40] and then modified with standard graphical applications.



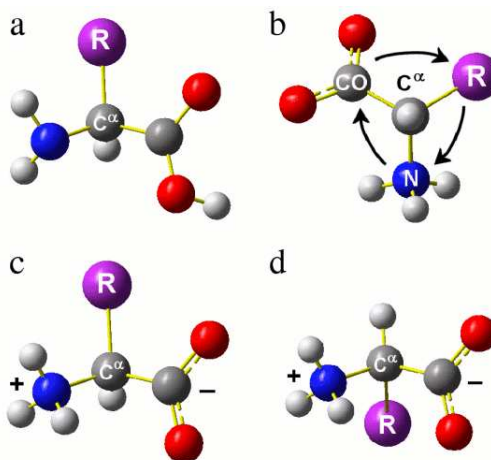


Figure 2.2: Amino acids. (a) Uncharged L-enantiomer. (b) CORN mnemotechnic rule to remember which one is the L-form. (c) Charged L-enantiomer (the predominant form found in living beings). (d) Charged D-enantiomer.

the D/L nomenclature is, by far, the most used one in protein science and the one that will be used in this work. For further details, take a look at the IUPAC recommendations at <http://www.chem.qmul.ac.uk/iupac/AminoAcid/>.

In principle, amino acids may be L- or D-, and the group  $\text{—R}$  may be anything provided that the resultant molecule is stable. However, for reasons that are still unclear [42], the vast majority of proteins in all living beings are made up of L-amino acids (as a rare exception, we may point out the fact that D-amino acids can be found in some proteins produced by exotic sea-dwelling organisms, such as *cone snails*) and the groups  $\text{—R}$  (called *side chains*) that are coded in the genetic material comprise a set of only twenty possibilities (depicted in figure 2.5).

A frequently quoted mnemotechnic rule for remembering which one is the L-form of amino acids is the so-called *CORN rule* in figure 2.2b. According to it, one must look from the hydrogen to the  $\alpha$ -carbon and, if the three remaining groups are labelled as in the figure, the word *CORN* must be read in the clockwise sense of rotation. The author of this work does not find this rule very useful, since normally he cannot recall if the sense is clockwise or counterclockwise. To know which form is the L- one, he draws the amino acid as in figure 2.2a or 2.2c, with the  $\alpha$ -carbon in the centre, the amino group on the left and the carboxyl group on the right, all of them in the plane of the paper (which is very natural and easy to remember because it matches the normal sense of writing with the fact that, conventionally, proteins start at  $\text{—NH}_3^+$  and end at  $\text{—COO}^-$ ). Finally, he must just remember that the side chain of the L-amino acid goes out of the paper approaching the reader (which is also natural

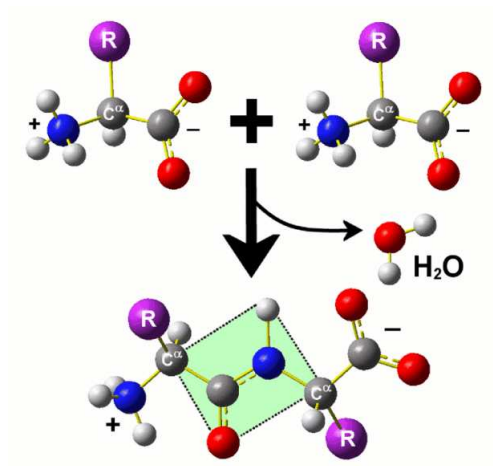


Figure 2.3: Peptide bond formation reaction. The peptide plane is indicated in green.

because the side chain is the relevant piece of information and we want to look at it closely).

The process through which amino acids are assembled into proteins (called *gene expression* or *protein biosynthesis*) is typically divided in two steps. In the first one, the *transcription*, the enzyme ARN polymerase (see figure 1.1b) binds to the DNA in the cellular nucleus and makes a copy of a section –the *gene*– of the base sequence into a messenger RNA (mRNA) molecule. In the second step, called *translation*, the mRNA enters the ribosome (see figure 1.1d) and is read stopping at each base triplet (called *codon*). Now, a specific molecule of transfer RNA (tRNA), which possesses the base triplet (called *anticodon*) that is complementary to the codon, links to the mRNA bringing with her the amino acid that is codified by the particular sequence of three bases. Each amino acid that arrives to the ribosome in this way is covalently attached to the previous one and so added to the nascent protein. In this reaction, the *peptide bond* is formed and a water molecule is released (see figure 2.3). This process continues until a stop codon is read and the transcription is complete.

The amino acid sequence of the resultant protein, read from the *amino terminus* to the *carboxyl terminus*, is called *primary structure*; and the amino acids included in such a polypeptide chain are normally termed *amino acid residues*, or simply *residues*, in order to distinguish them from their isolated form. The main chain formed by the repetition of  $\alpha$ -carbons and the C' and N atoms at the peptide bond is called *backbone* and the —R groups branching out from it are called *side chains*, as it has already been mentioned.

The specificity of each protein is provided by the different properties of the twenty side chains in figure 2.5 and their particular positions in the sequence. In textbooks, it is customary to group them in small sets according to different

criteria in order to facilitate their learning. Classifications devised on the basis of the physical properties of these side chains may be sometimes overlapping (e.g., tryptophan contains polar regions as well as an aromatic ring, which, in turn, could be considered hydrophobic but is also capable of participating in, say,  $\pi$ - $\pi$  interactions). Therefore, for a clearer presentation, we have chosen here to classify the residues according to the chemical groups contained in each side chain and discuss their physical properties individually.

Let us enumerate then the categories in figure 2.5 and point out any special remark regarding the residues in them:

#### Special residues:

- *Glycine* is the smallest of all the amino acids: its side chain contains only a hydrogen atom. So, since its  $\alpha$ -carbon has two hydrogens attached, glycine is the only achiral natural amino acid. Its affinity for water is mainly determined by the peptide groups in the backbone; therefore, glycine is hydrophilic.
- *Proline* is the only residue whose side chain is covalently linked to the backbone (the backbone is indicated in purple in figure 2.5), giving proline unique structural properties that will be discussed later. Since its side chain is entirely aliphatic, proline is hydrophobic.

#### Sulfur-containing residues:

- *Cysteine* is a very important structural residue because, in a reaction catalyzed by *protein disulfide isomerases* (PDIs), it may form, with another cysteine, a very stable covalent bond called *disulfide bond* (see figure 2.4). Curiously, all the L-amino acids are S-enantiomers according to the Cahn, Ingold and Prelog rules [41] except for cysteine, which is R-. This is probably the reason that makes the D/L nomenclature favourite among protein scientists [24]. Cysteine is a polar residue.

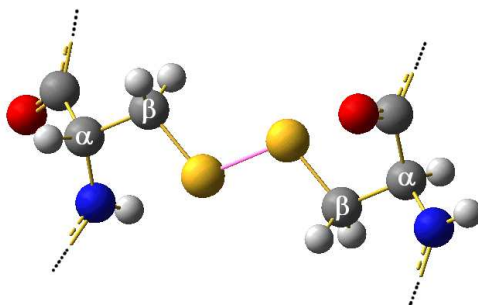


Figure 2.4: Disulfide bond between two cysteine residues.

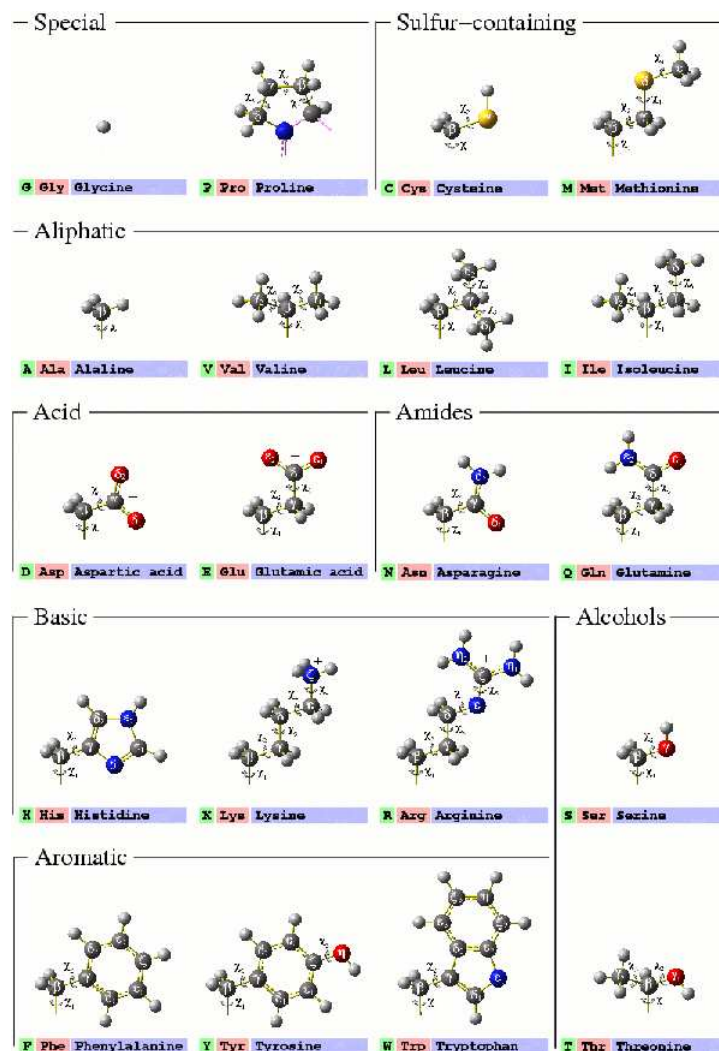


Figure 2.5: Side chains of the twenty amino acid residues encoded in the genetic material of living beings. They have been classified according to the chemical groups they contain. The rotameric degrees of freedom  $\chi_i$  are indicated with small arrows over the bonds. The name of the heavy atoms and the numbering of the branches comply with the IUPAC rules <http://www.chem.qmul.ac.uk/iupac/AminoAcid/>). Below the molecular structure, the one letter code (green), the three letter code (red) and the complete name (blue) of each amino acid may be found. In the case of proline, the N and the  $\alpha$ -carbon have been included in the scheme, and the backbone bonds have been coloured in purple. The titratable residues Asp, Glu, Lys and Arg have been represented in their charged forms, which is the most common one in aqueous solvent under physiological conditions. Histidine is shown in its neutral  $\varepsilon_2$ -tautomeric form.

- *Methionine* is mostly aliphatic and, henceforth, apolar.

#### Aliphatic residues:

- *Alanine* is the smallest chiral residue. This is the fundamental reason for using alanine models, more than any other ones, in the computationally demanding ab initio studies of peptides that are customarily performed in quantum chemistry [43–53]. It is hydrophobic, like all the residues in this group.
- *Valine* is one of the three  $\beta$ -branched residues (i.e., those that have more than one heavy atom attached to the  $\beta$ -carbon, apart from the  $\alpha$ -carbon), together with isoleucine and threonine. It is hydrophobic.
- *Leucine* is hydrophobic.
- *Isoleucine*'s  $\beta$ -carbon constitutes an asymmetric centre and the only enantiomer that occurs naturally is the one depicted in the figure. Only isoleucine and threonine contain an asymmetric centre in their side chain. Isoleucine is  $\beta$ -branched and hydrophobic.

#### Acid residues:

- *Aspartic acid* is normally charged under physiological conditions. Hence, it is very hydrophilic.
- *Glutamic acid* is just one  $\text{CH}_2$  larger than aspartic acid. Their properties are very similar.

#### Amides:

- *Asparagine* contains a chemical group similar to the peptide bond. It is polar and can act as a hydrogen bond donor or acceptor.
- *Glutamine* is just one  $\text{CH}_2$  larger than asparagine. Their physical properties are very similar.

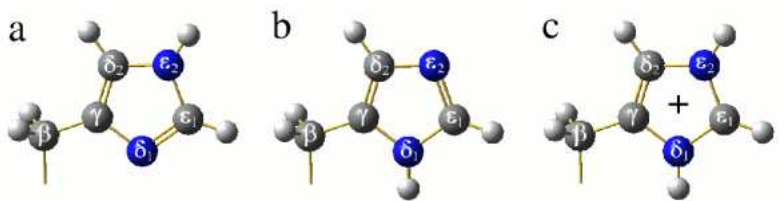


Figure 2.6: Three forms of histidine found in proteins. (a) Neutral  $\epsilon_2$ -tautomer. (b) neutral  $\delta_1$ -tautomer. (c) Charged form.

### Basic residues:

- *Histidine* is a special amino acid: in its neutral form, it may exist as two different tautomers, called  $\delta_1$  and  $\varepsilon_2$ , depending on which nitrogen has an hydrogen atom attached to it. The  $\varepsilon_2$ -tautomer has been found to be slightly more stable in model dipeptides [54], although both forms are found in proteins. Histidine can readily accept a proton and get a positive charge, in fact, it is the only side chain with a pKa in the physiological range, so non-negligible proportions of both the charged and uncharged forms are typically present. Of course, histidine is hydrophilic.
- *Lysine*'s side chain is formed by a rather long chain of  $\text{CH}_2$  with an amino group at its end, which is nearly always positively charged. Therefore, lysine is very polar and hydrophilic.
- *Arginine*'s properties are similar to those of lysine, although its terminal guanidinium group is a stronger basis than the amino group and it may also participate in hydrogen bonds as a donor.

### Alcohols:

- *Serine* is one of the smallest residues. It is polar due to the hydroxyl group.
- *Threonine*'s  $\beta$ -carbon constitutes an asymmetric centre; the enantiomer that occurs in living beings is the one shown in the figure. The physical properties of threonine are very similar to those of serine.

### Aromatic residues:

- *Phenylalanine* is the smallest aromatic residue. Its benzyl side chain is largely apolar and interacts unfavourably with water. It may also participate in specific  $\pi$ -stacking interactions with other aromatic groups.
- *Tyrosine*'s properties are similar to those of phenylalanine, being only slightly more polar due to the presence of a hydroxyl group.
- *Tryptophan*, with 17 atoms in her side chain, is the largest residue. It is mainly hydrophobic, although it contains a small polar region and it can also participate in  $\pi$ - $\pi$  interactions, like all the residues in this category.

After having introduced the building blocks of proteins, some qualifying remarks about them are worth to be done: On one side, why amino acids encoded in DNA codons are the ones in the list or why there are exactly twenty of them are questions that are still subjects of controversy [55,56]. In fact, although the side chains in figure 2.5 seem to confer enough versatility to proteins in most cases, there are also rare exceptions in which other groups are needed to perform a particular function. For example, the amino acid *selenocysteine* may be incorporated into some proteins at an UGA codon (which normally indicates

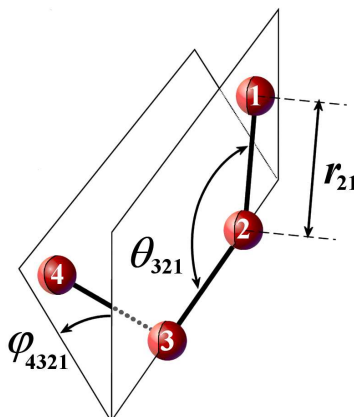


Figure 2.7: Typical definition of internal coordinates.  $r_{21}$  is the *bond length* between atoms 2 and 1.  $\theta_{321}$  is the *bond angle* formed by the bonds (2,1) and (3,2), it ranges from 0 to  $180^\circ$ . Finally  $\varphi_{4321}$  is the *dihedral angle* describing the rotation around the bond (3,2); it is defined as the angle formed by the plane containing atoms 1, 2 and 3 and the plane containing atoms 2, 3 and 4; it ranges either from  $-180^\circ$  to  $180^\circ$  or from  $0^\circ$  to  $360^\circ$ , depending on the convention; the positive sense of rotation for  $\varphi_{4321}$  is the one indicated in the figure. Also note that the definition is symmetric under a complete change in the order of the atoms, in such a way that, quite trivially,  $r_{21} = r_{12}$  and  $\theta_{123} = \theta_{321}$ , but also, not so trivially,  $\varphi_{4321} = \varphi_{1234}$ . (See reference [51] for further information.)

a stop in the transcription), or the amino acid *pyrrolysine* at an UAG codon (which is also a stop indication in typical cases). In addition, the arginine side chain may be post-translationally converted into *citrulline* by the action of a family of enzymes called *peptidylarginine deiminases* (PADs).

On the other hand, the chemical (covalent) structure of the protein chain may suffer from more complex modifications than just the inclusion of non-standard amino acid residues: A myriad of organic molecules may be covalently linked to specific points, the chain may be cleaved (cut), chemical groups may be added or removed from the N- or C-termini, disulfide bonds may be formed between cysteines, and the side chains of the residues may undergo chemical modifications just like any other molecule [54]. The vast majority of these changes either depend on the existence of some chemical agent external to the protein, or are catalyzed by an enzyme.

In this work, our interest is in the folding of proteins. This problem, which will be discussed in detail in the next section, is so huge and so difficult that, in the opinion of the author, there is no point in worrying about details, such as the ones mentioned in the two preceding paragraphs, before the big picture is at least preliminarily understood. Therefore, when we talk about the folding



of proteins in what follows, we will be thinking about single polypeptide chains, made up of L-amino acids, in water and without any other reagent present, with the side chains chosen from the set in figure 2.5, and having undergone no post-translational modifications nor any chemical change on their groups. Finally, although some simple modifications, such as the formation of disulfide bonds or the  $\text{trans} \rightarrow \text{cis}$  isomerization of Xaa-Pro peptide bonds (see what follows), could be more easily included in the first approach to the problem, we shall also leave them for a later stage.

Now, with this considerations, we have fixed the covalent structure of our molecule as well as the enantiomerism of the asymmetric centres it may contain. This information is enough to specify the three-dimensional arrangement of the atoms of small rigid molecules. However, long polymers and, particularly, proteins, possess degrees of freedom (termed *soft*) that require small amounts of energy to be changed while drastically altering the relative positions of groups and atoms. In a first approximation, all bond lengths, bond angles and dihedral angles describing rotations around triple, double and partial double bonds (see figure 2.7) may be considered to be determined by the covalent structure. Whereas dihedral angles describing rotations around single bonds may be considered to be variable and soft. The non-superimposable three-dimensional arrangements of the molecule that correspond to different values of the soft degrees of freedom are called *conformations*.

In proteins, some of these soft dihedrals are located at the side chains; they are the  $\chi_i$  in figure 2.5 and, although they are important in the later stages of

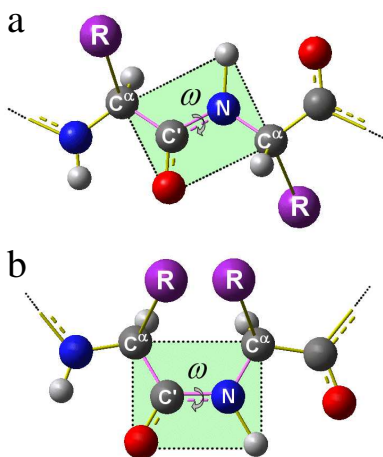


Figure 2.8: Trans and cis conformations of the peptide plane. The bonds defining the peptide bond dihedral angle  $\omega$  are indicated in purple. (a) *Trans* conformation ( $\omega \simeq \pm 180^\circ$ ). The most common one in proteins. (b) *Cis* conformation ( $\omega \simeq 0^\circ$ ). Significantly found only in Xaa-Pro bonds.

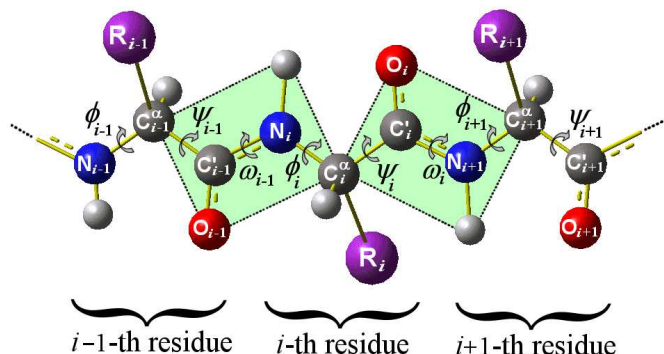


Figure 2.9: Numeration of the heavy atoms and the dihedrals angles describing rotations around backbone bonds. In agreement with IUPAC recommendations (see <http://www.chem.qmul.ac.uk/iupac/AminoAcid/>). The peptide planes are indicated as green rectangles.

the folding process and must be taken into account in any ambitious model of the system, their variation only alters the conformation locally. On the contrary, a small change in the dihedral angles located at the backbone of the polypeptide chain may drastically modify the relative position of many pairs of atoms and they must be given special attention.

That is why, the special properties of the peptide bond, which is the basic building block of the backbone, are very important to understand the conformational behaviour of proteins. These properties arise from the fact that there is an electron pair delocalized between the C—N and C—O bonds (using the common chemical image of *resonance*), which provokes that neither bond is single nor double, but *partial double bonds* that have a mixed character. In particular, the partial double bond character of the peptide bond is the cause that the six atoms in the green plane depicted in figures 2.3, 2.8 and 2.9 have a strong tendency to be coplanar, forming the so-called *peptide plane*. This coplanarity allows for only two different conformations: the one called *trans* (corresponding to  $\omega \simeq \pm 180^\circ$ ), in which the  $\alpha$ -carbons lie at different sides of the line containing the C—N bond; and the one called *cis* (corresponding to  $\omega \simeq 0^\circ$ ), in which they lie at the same side of that line (see figure 2.8).

Although the quantitative details are not completely elucidated yet and the very protocol of protein structure determination by x-ray crystallography could introduce spurious effects in the structures deposited in the PDB [57], it seems clear that a great majority of the peptide bonds in proteins are in the *trans* conformation. Indeed, a superficial look at the two forms in figure 2.8 suggests that the steric clashes between substituents of consecutive  $\alpha$ -carbons will be more severe in the *cis* case. When the second residue is a Proline, however, the special structure of its side chain makes the probability of finding the *cis* con-

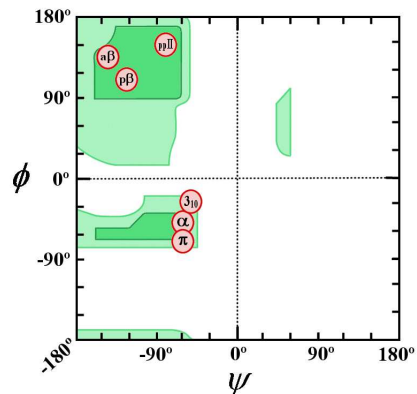


Figure 2.10: Original Ramachandran plot drawn by Ramachandran and Ramakrishnan in 1963 [60]. In dark-green, the fully allowed regions, calculated by letting the atoms approach to the average clashing distance; in light-green, the partially allowed regions, calculated by letting the atoms approach to the minimum clashing distance; in white, the disallowed regions. Some points representing secondary structure elements are shown as red circles at the ideal  $(\phi, \psi)$ -positions in table 1: ( $\alpha$ )  $\alpha$ -helix. ( $\pi$ )  $\pi$ -helix. ( $3_{10}$ )  $3_{10}$ -helix. ( $a\beta$ ) Antiparallel  $\beta$ -sheet. ( $p\beta$ ) Parallel  $\beta$ -sheet. ( $ppII$ ) Polyproline II.

former significantly higher: For Xaa-nonPro peptide bonds in native structures, the trans form is more common than the cis one with approximately a 3000:1 proportion; while this ratio decreases to just 15:1 if the bond is Xaa-Pro [57].

In any case, due to the aforementioned partial double bond character of the C—N bond, the rotation barrier connecting the two states is estimated to be of the order of  $\sim 20$  kcal/mol [58], which is about 40 times larger than the thermal energy at physiological conditions, thus rendering the spontaneous trans  $\rightarrow$  cis isomerization painfully slow. However, mother Nature makes use of every possibility that she has at hand and, sometimes, there are a few peptide bonds that must be cis in order for the protein to fold correctly or to function properly. Since all peptide bonds are synthesized trans at the ribosome [59], the trans  $\rightarrow$  cis isomerization must be catalyzed by enzymes (called *peptidylprolyl isomerases* (PPIs)) and, in the same spirit of the post-translational modifications discussed before, this step may be taken into account in a later refinement of the theoretical models.

Therefore, we shall assume in what follows that all peptide bonds (even the Xaa-Pro ones) are in the trans state and, henceforth, the conformation of the protein will be essentially determined by the values of the  $\phi$  and  $\psi$  angles, which describe the rotation around the two single bonds next to each  $\alpha$ -carbon (see figure 2.9 for a definition of the dihedral angles associated to the backbone).

This assumption was introduced, as early as 1963, by Ramachandran and

	$\phi$	$\psi$
$\alpha$ -helix	-57	-47
$3_{10}$ -helix	-49	-26
$\pi$ -helix	-57	-70
polyproline II	-79	149
parallel $\beta$ -sheet	-119	113
antiparallel $\beta$ -sheet	-139	135

Table 1: Ramachandran angles (in degrees) of some important secondary structure elements in polypeptides. Data taken from reference [1].

Ramakrishnan [60] and the  $\phi$  and  $\psi$  coordinates are commonly named *Ramachandran angles* after the first one of them. In their famous paper [60], they additionally suppose that the bond lengths, bond angles and dihedral angles on double and partial double bonds are fixed and independent of  $\phi$  and  $\psi$ , they define a typical distance up to which a specific pair of atoms may approach and also a minimum one (taken from statistical studies of structures) and they draw the first *Ramachandran plot* (see figure 2.10): A depiction of the regions in the  $(\phi, \psi)$ -space that are energetically allowed or disallowed on the basis of the local sterical clashes between atoms that are close to the  $\alpha$ -carbon.

One of the main advantages of this type of diagrams as ‘thinking tools’ lies in the fact that (always in the approximation that the non-Ramachandran variables are fixed) some very common repetitive structures found in proteins may be ideally depicted as a single point in the plot. In fact, these special conformations, which are regarded as the next level of protein organization after the primary structure and are said to be elements of *secondary structure*, may be characterized exactly like that, i.e., by asking that a certain number of consecutive residues present the same values of the  $\phi$  and  $\psi$  angles. In the book by Lesk [1], for example, one may find a table with the most common of these repetitive patterns, together with the corresponding  $(\phi, \psi)$ -values taken from statistical investigations of experimentally resolved protein structures (see table 1).

However, the non-Ramachandran variables are not really constant, and the elements of secondary structure do possess a certain degree of flexibility. Moreover, the side chains may interact and exert different strains at different points of the chain, which provokes that, in the end, the secondary structure elements gain some stability by slightly altering their ideal Ramachandran angles. Therefore, it is more appropriate to characterize them according to their hydrogen-bonding pattern, which, in fact, is the feature that makes these structures prevalent, providing them with more energetic stability than other repetitive conformations which are close in the Ramachandran plot.

The first element of secondary structure that was found is the  $\alpha$ -*helix*. It is a coil-like<sup>16</sup> structure, with  $\sim 3.6$  residues per turn, in which the carbonyl

<sup>16</sup> Here, we use the word ‘coil’ to refer to the twisted shape of a telephone wire, a corkscrew

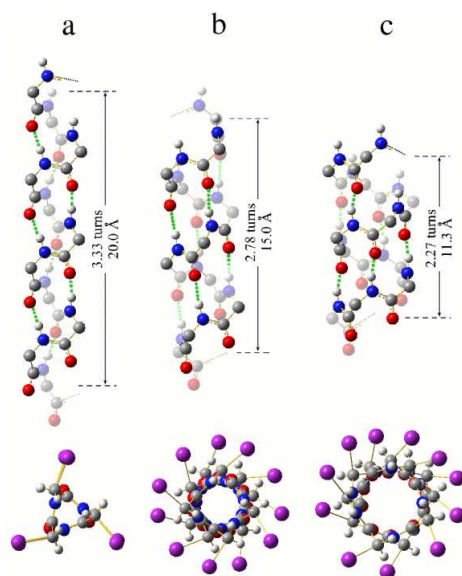


Figure 2.11: The three helices found in protein native structures. **(a)**  $3_{10}$ -helix, **(b)**  $\alpha$ -helix, and **(c)**  $\pi$ -helix. In the three cases, the helices shown are 11-residues long. In the standing views (above), the hydrogen bonds are depicted as green dotted lines and the distance and number of turns spanned by 10 residues are indicated at the right of the structures. Whereas in the standing views, the side chains and  $\alpha$ -hydrogens have been removed for visual convenience, in the zenithal views (below), they are included.

group (C=O) of each  $i$ -th residue forms a hydrogen bond with the amino group (N-H) of the residue  $i + 4$  (see figure 2.11b). According to a common notation, in which  $x_y$  designates a helix with  $x$  residues per turn and  $y$  atoms in the ring closed by the hydrogen bond [61], the  $\alpha$ -helix is also called  $3.6_{13}$ -helix.

She was theoretically proposed in 1951 by Pauling, Corey and Branson [62], who used precise information about the geometry of the non-Ramachandran variables, taken from crystallographic studies of small molecules, to find the structures compatible with the additional constraints that: (i) the peptide bond is planar, and (ii) every carbonyl and amino group participates in a hydrogen bond.

The experimental confirmation came from Max Perutz, who, together with Kendrew and Bragg, had proposed in 1950 (one year before Pauling's paper) a series of helices with an integer number of residues per turn [61] that are not so commonly found in native structures of proteins (see however, the discussion

---

or the solenoid of an electromagnet. Although this is common English usage, the same word occurs frequently in protein science to designate different (and sometimes opposed) concepts. For example, a much used ideal model of the denatured state of proteins is termed *random coil*, and a popular statistical description of helix formation is called *helix-coil theory*.

about the  $3_{10}$ -helix below). Perutz read Pauling, Corey and Branson’s paper one Saturday morning [63] in spring 1951 and realized immediately that their helix looked very well: free of strain and with all donor and acceptor groups participating in hydrogen bonds. So he rushed to the laboratory and put a sample of horse hair (rich in keratin, a protein that contains  $\alpha$ -helices) in the x-ray beam, knowing that, according to diffraction theory, the regular repeat of the ‘spiral staircase steps’ in Pauling’s structure should give rise to a strong x-ray reflection of 1.5 Å spacing from planes perpendicular to the fiber axis. The result of the experiment was positive<sup>17</sup> and, in the last years of the 50s, Perutz and Kendrew saw again the same signal in myoglobin and hemoglobin, when they resolved, for the first time in history, the structure of these proteins [64,65].

However, despite its being, by far, the most common, the  $\alpha$ -helix is not the only coil-like structure that can be found in native proteins [66–68]. If the hydrogen bonds are formed between the carbonyl group (C=O) of each  $i$ -th residue with the amino group (N–H) of the residue  $i + 3$ , one obtains a  $3_{10}$ -helix, which is more tightly wound and, therefore, longer than an  $\alpha$ -helix of the same chain length (see figure 2.11a). The  $3_{10}$ -helix is the fourth most common conformation for a single residue after the  $\alpha$ -helix,  $\beta$ -sheet and reverse turn<sup>18</sup> [67] but, remarkably, due to its having an integer number of residues per turn, it seemed more natural to scientists with crystallographic background and was theoretically proposed before the  $\alpha$ -helix [61,69]. On the other hand, if the hydrogen bonds are formed between the carbonyl group (C=O) of each  $i$ -th residue with the amino group (N–H) of the residue  $i + 5$ , one obtains a  $\pi$ -helix (or  $4.4_{16}$ -helix), which is wider and shorter than an  $\alpha$ -helix of the same length (see figure 2.11c). It was originally proposed by Low and Baybutt in 1952 [70], and, although the exact fraction of each type of helix in protein native structures depends up to a considerable extent on their definition (in terms of Ramachandran angles, interatomic distances, energy of the hydrogen bonds, etc.), it seems clear that the  $\pi$ -helix is the less common of the three [66]. Now, it is true that, in addition to these helices that have been experimentally confirmed, some others have been proposed. For example, in the same work in which Pauling, Corey and Branson introduce the  $\alpha$ -helix [62], they also describe another candidate: the  $\gamma$ -helix (or  $5.1_{17}$ -helix). Finally, Donohue performed, in 1953, a systematic study of all possible helices and, in addition to the ones already mentioned, he proposed a  $2.2_7$ - and a  $4.3_{14}$ -helix [71]. None of them has been detected in resolved native proteins.

Among the secondary structure elements of proteins, not all regular local patterns are helices: there exist also a variety of repetitive conformations that do not contain strong intra-chain hydrogen bonds and that are less curled than the structures in figure 2.11. For example, the *polyproline II* [72–74], which is

---

<sup>17</sup> Linus Pauling was awarded the Nobel prize in chemistry in 1954 ‘for his research into the nature of the chemical bond and its application to the elucidation of the structure of complex substances’, and Max Perutz shared it with John Kendrew in 1962 ‘for their studies of the structures of globular proteins’.

<sup>18</sup> A conformation that some residues in proteins adopt when an acute turn in the chain is needed.

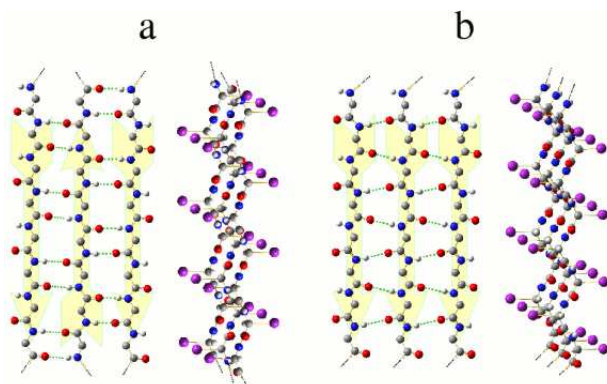


Figure 2.12:  $\beta$ -sheets in the pure (a) antiparallel, and (b) parallel versions. On the left, the top view is shown, with the side chains and the  $\alpha$ -hydrogens omitted for visual convenience and the directions of the strands indicated as yellow arrows. The hydrogen bonds are represented as green dotted lines. On the right, the side view of the sheets is depicted. In this case, the side chains and the  $\alpha$ -hydrogens are included.

thought to be important in the unfolded state of proteins, and, principally, the family of the  $\beta$ -sheets, which are, together with the  $\alpha$ -helices, the most recognizable secondary structure elements in native states of polypeptide chains<sup>19</sup>.

The  $\beta$ -sheets are rather plane structures that are typically formed by several individual  $\beta$ -strands, which align themselves to form stabilizing inter-chain hydrogen bonds with their neighbours. Two pure arrangements of these single threads may be found: the *antiparallel*  $\beta$ -sheets (see figure 2.12a), in which the strands run in opposite directions (read from the amino to the carboxyl terminus); and the *parallel*  $\beta$ -sheets (see figure 2.12b), in which the strands run in the same direction. In both cases, the side chains of neighbouring residues in contiguous strands branch out to the same side of the sheet and may interact. Of course, mixed parallel-antiparallel sheets can also be found.

The next level of protein organization, produced by the assembly of the elements of secondary structure, and also of the chain segments that are devoid of regularity, into a well defined three-dimensional shape, is called *tertiary structure*. The protein folding problem (omitting relevant qualifications that have been partially made and that will be recalled and made more explicit in what follows) may be said to be *the attempt to predict the secondary and the tertiary*

<sup>19</sup> It is probably more correct to define the *secondary structure* as the conformational repetition in *consecutive* residues and, from this point of view, to consider the  $\beta$ -strand as the proper element of secondary structure. In this sense, the assembly of  $\beta$ -strands, the  $\beta$ -sheet, together with some other simple motifs such as the coiled coils made up of two helices, the silk fibroin (made up of stacked  $\beta$ -sheets) or collagen (three coiled threads of a repetitive structure similar to polyproline II), may be said to be elements of *super-secondary structure*, somewhat in between the local secondary structure and the global and more complex tertiary structure (see below).



*structure from the primary structure*, and it will be discussed in the next section.

The *quaternary structure*, which refers to the way in which protein monomers associate to form more complex systems made up of more than one individual chain (such as the ones in figure 1.1), will not be explored in this work.

### 3 The protein folding problem

As we have seen in the previous section and can visually check in figure 1.1, the biologically functional *native* structure of a protein<sup>20</sup> is highly complex. What Kendrew saw in one of the first proteins ever resolved is essentially true for most of them [75]:

*The most striking features of the molecule were its irregularity and its total lack of symmetry.*

Now, since these polypeptide chains are synthesized linearly in the ribosome (i.e., they are not manufactured in the folded conformation), in principle, one may imagine that some specific cellular machinery could be the responsible of the complicated process of folding and, in such a case, the prediction of the native structure could be a daunting task. However, in a series of experiments in the 50s, Christian B. Anfinsen ruled out this scenario and was awarded the Nobel prize for it [76].

The most famous and illuminating experiment that he and his group performed is the refolding of bovine pancreatic ribonuclease (see the scheme in figure 3.1 for reference). They took this protein, which is 124 residues long and has all her eight cysteines forming four disulfide bonds, and added, in a first step, some reducing agent to cleave them. Then, they added urea up to a concentration of 8 M. This substance is known for being a strong denaturing agent (an ‘unfold’er) and produced a ‘scrambled’ form of the protein which is much less compact than the native structure and has no enzymatic activity. From this scrambled state, they took two different experimental paths: in the *positive* one, they removed the urea first and then added some oxidizing agent to reform the disulfide bonds; whereas, in the *negative* path, they poured the oxidizing agent first and removed the urea in a second step.

The resultant species in the two paths are very different. If one removes the urea first and then promotes the formation of disulfide bonds, an homogeneous sample is obtained that is practically indistinguishable from the starting native protein and that keeps full biological activity. The ribonuclease has been ‘unscrambled’! However, if one takes the negative path and let the cysteines form disulfide bonds before removing the denaturing agent, a mixture of products is obtained containing many or all of the possible 105 isomeric disulfide

---

<sup>20</sup> Most native states of proteins are flexible and are comprised not of only one conformation but of a set of closely related structures. This flexibility is essential if they need to perform any biological function. However, to economize words, we will use in what follows the terms *native state*, *native conformation* and *native structure* as interchangeable.

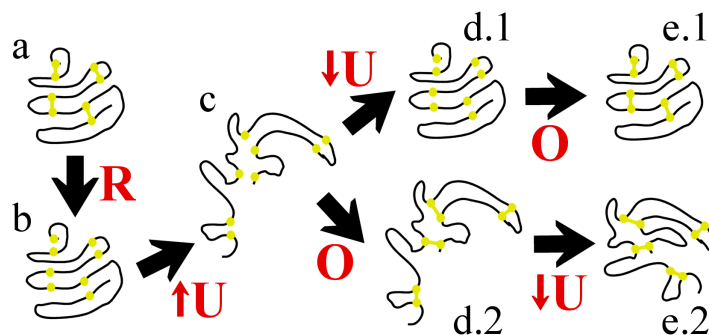


Figure 3.1: Scheme of the refolding of the bovine pancreatic ribonuclease by Anfinsen. The black arrows indicate fundamental steps of the experiment and the red labels next to them designate: **(R)** addition of reducing agent (cleavage of the disulfide bonds), **(O)** addition of oxidizing agent (reformation of the disulfide bonds), **(↑U)** and **(↓U)** increase of the urea concentration up to 8 M and decrease to 0 M respectively. The conformation of the backbone of the protein is schematically depicted by a black line, the cysteines are shown as small yellow circles and the disulfide bonds as line segments connecting them. The different states are labelled: **(a)** starting native enzyme with full activity, **(b)** non-disulfide bonded, folded form, **(c)** representant of the ensemble of inactive ‘scrambled’ ribonuclease, **(d.1)** non-disulfide bonded, folded form, **(e.1)** refolded ribonuclease indistinguishable from (a), **(d.2)** representant of the ensemble of the scrambled, disulfide bonded form, and, finally, **(e.2)** representant of the mixture of the 105 isomeric disulfide bonded forms.

bonded forms<sup>21</sup>. This mixture is essentially inactive, having approximately 1% the activity of the native enzyme.

One of the most clear conclusions that are commonly drawn from this experiment is that *all the information needed to reach the native state is encoded in the sequence of amino acids*. This important statement, which has stood the test of time [24, 77], allows to isolate the system under study (both theoretically and experimentally) and sharply defines the *protein folding problem*, i.e., the prediction of the three-dimensional native structure of proteins from their amino acid sequence (and the laws of physics).

It is true that we nowadays know of the existence of the so-called *molecular chaperones* (see, for example, the GroEL-GroES complex in figure 1.1c), which help the proteins fold in the cellular milieu [78–82]. However, according to the most accepted view [24, 77], these molecular assistants do not add any structural information to the process. Some of them simply prevent accidents related to the *cellular crowding* from happening. Indeed, in the cytoplasm there is not much room: inside a typical bacterium, for example, the total macromolecular

<sup>21</sup> Take an arbitrary cysteine: she can bond to any one of the other seven. From the remaining six, take another one at random: she can bond to five different partners. Take the reasoning to its final and we have  $7 \times 5 \times 3 = 105$  different possibilities.

concentration is approximately 350 mg/ml, whereas a typical protein crystal may contain about 600 mg/ml [77]. This crowding may hinder the correct folding of proteins, since partially folded states (of chains that are either free in the cytoplasm or being synthesized in proximate ribosomes) have more ‘sticky’ hydrophobic surface exposed than the native state, opening the door to aggregation. In order to avoid it, some chaperonins<sup>22</sup> are in charge of providing a shelter in which the proteins can fold alone. Yet another pitfall is that, when the polypeptide chain is being synthesized in the ribosome, it may start to fold incorrectly and get trapped in a non-functional conformation separated by a high energetic barrier from the native state. Again, there exist some chaperones that bind to the nascent chain to prevent this from happening. As we have already pointed out, all this assistance to fold is seen as lacking new structural information and meant only to avoid traps which are not present *in vitro*. It seems as if molecular chaperones’ aim is to make proteins believe that they are not in a messy cell but in Anfinsen’s test tube!

The possibility that this state of affairs opens, the prediction of the three-dimensional native structure of proteins from the only knowledge of the amino acid sequence, is often referred to as ‘the second half of the genetic code’ [83,84]. The reason for such a vehement statement lays in the fact that not all proteins are accessible to the experimental methods of structure resolution (mostly x-ray crystallography and NMR [54,85]) and, for those that can be studied, the process is long and expensive, thus making the databases of known structures grow much more slowly than the databases of known sequences (see figure 1.2 and the related discussion in section 1). To solve ‘the second half of the genetic code’ and bridge this gap is the main objective of the hot scientific field of *protein structure prediction* [24,86,87].

The path that takes to this goal may be walked in two different ways [88,89]: Either at a fast pragmatic pace, using whatever information we have available, increasingly refining the everything-goes prediction procedures by extensive trial-and-error tests and without any need of knowing the details of the physical processes that take place; or at a slow thoughtful pace, starting from first principles and seeking to arrive to the native structure using the same means that Nature uses: the laws of physics.

The different protocols belonging to the fast pragmatic way are commonly termed *knowledge-based*, since they take profit from the already resolved structures that are deposited in the PDB [32] or any other empirical information that may be statistically extracted from databases of experimental data. There are basically three pure forms of knowledge-based strategies [90]:

- *Homology modeling* (also called *comparative modeling*) [85,91] is based on the observation that proteins with similar sequences frequently share similar structures [92]. Following this approach, either the whole sequence of the protein that we want to model (the *target*) or some segments of it are *aligned* to a sequence of known structure (the *template*). Then, if some reasonable measure of the *sequence similarity* [93,94] is high enough,

---

<sup>22</sup> A particular subset of the set of molecular chaperones.

the structure of the template is proposed to be the one adopted by the target in the region analysed. Using this strategy, one typically needs more than 50% sequence identity between target and template to achieve high accuracy, and the errors increase rapidly below 30% [87]. Therefore, comparative modeling cannot be used with all sequences, since some recent estimates indicate that  $\sim 40\%$  of genes in newly sequenced genomes do not have significant sequence homology to proteins of known structure [95].

- *Fold recognition* (or *threading*) [86, 96] is based on the fact that, increasingly, new structures deposited in the PDB turn out to fold in shapes that have been seen before, even though conventional sequence searches fail to detect the relationship [97]. Hence, when faced to a sequence that shares low identity with the ones in the PDB, the threading user tries to fit it in each one of the structures in the databases of known folds, selecting the best choices with the help of some scoring function (which may be physics-based or not). Again, fold recognition methods are not flawless and, according to various benchmarks, they fail to select the correct fold from the databases for  $\sim 50\%$  of the cases [86]. Moreover, the fold space is not completely known so, if faced with a novel fold, threading strategies are useless and they may even give false positives. Modern studies estimate that approximately one third of known protein sequences must present folds that have never been seen [98].
- *New fold* (or *de novo prediction*) methods [99, 100] must be used when the protein under study has low sequence identity with known structures and fold recognition strategies fail to fit it in a known fold (because of any of the two reasons discussed in the previous point). The specific strategies used in new fold methods are very heterogeneous, ranging from well-established secondary structure prediction tools or sequence-based identification of sets of possible conformations for short fragments of chains to numerical search methods, such as molecular dynamics, Monte Carlo or genetic algorithms [97].

These knowledge-based strategies may be arbitrarily combined into mixed protocols, and, although the frontiers between them may be sometimes blurry [35], it is clear that the more information available the easier to predict the native structure (see figure 3.2). So that the three types of methods described above turn out to be written in increasing order of difficulty and they essentially coincide with the competing categories of the *CASP experiment*<sup>23</sup> [35, 97]. In this important meeting, held every two years and whose initials loosely stand for *Critical Assessment of techniques for protein Structure Prediction*, experimental structural biologists are asked to release the amino acid sequences of proteins (the *CASP targets*) whose structures are likely to be resolved before the contest starts. Then, the ‘prediction community’ gets on stage and their

---

<sup>23</sup> Since CASP1, people has drifted towards knowledge-based methods and, nowadays, very few groups use pure ab initio approaches [101].

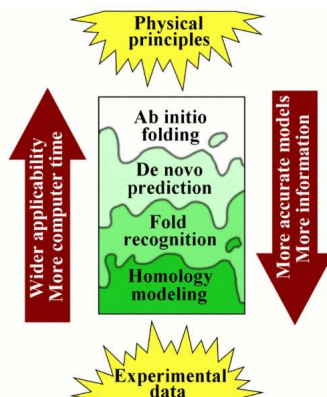


Figure 3.2: Schematic classification of protein structure prediction methods.

members submit the proposed structures (the *models*), which may be found using any chosen method. Finally, a committee of assessors, critically evaluate the predictions, and the results are published, together with some contributions by the best predictors, in a special issue of the journal *Proteins*.

Precisely, in the latest CASP meetings, the expected ordering (based on the available experimental information) of the three aforementioned categories of protein structure prediction has been observed to translate into different qualities of the proposed models (see figure 3.2). Hence, while comparative modeling with high sequence similarity has proved to be the most reliable method to predict the native conformation of proteins (with an accuracy comparable to low-resolution, experimentally determined structures) [87, 102], de novo modeling has been shown to remain still unreliable [35, 88] (although a special remark should be made about the increasingly good results that David Baker and his group are achieving in this field with their program Rosetta [103, 104]).

Opposed to these knowledge-based approaches, the computer simulation of the real physical process of protein folding<sup>24</sup> without using any empirical information and starting from first principles could be termed *ab initio protein folding* or *ab initio protein structure prediction* depending whether the emphasis is laid on the process or on the goal.

Again, the frontier between de novo modeling and ab initio protein folding is not sharply defined and some confusion might arise between the two terms. For example, the potential energy functions included in most empirical force fields such as CHARMM [106, 107] contain parameters extracted from experimental data, while molecular dynamics attempts to fold proteins using these force fields will be considered by most people (including the author) to belong to the ab initio category. As always, the limit cases are clearer, and Baker’s Rosetta [103, 104], which uses statistical data taken from the PDB to bias the secondary

<sup>24</sup> Not a new idea [105].

structure conformational search, may be classified, without any doubt, as a de novo protocol; while, say, a (nowadays unfeasible) simulation of the folding process using quantum mechanics, would be deep in the ab initio region. The situation is further complicated due to the fact that score functions which are based (up to different degrees) on physical principles, are commonly used in conjunction with knowledge-based strategies to prune or refine the candidate models [87, 108]. In the end, the classification of the strategies for finding the native structure of proteins is rather continuous with wriggly, blurry frontiers (see figure 3.2).

It is clear that, despite their obvious practical advantages and the superior results when compared to pure ab initio approaches [86], any knowledge-based features included in the prediction protocols render the assembly mechanisms physically meaningless [109]. If we want to know the real details of protein folding as it happens, for example, to properly study and attack diseases that are related to protein misfolding and aggregation [3], we must resort to pure ab initio strategies. In addition, ab initio folding does not require any experimental information about the protein, apart from its amino acid sequence. Therefore, as new fold strategies, it has a wider range of applicability than homology modeling and fold recognition, and, in contrast with the largely system-oriented protocols developed in the context of knowledge-based methods, most theoretical and computational improvements made while trying to ab initio fold proteins will be perfectly applicable to other macromolecules.

The feedback between strategies is also an important point to stress. Apart from the obvious fact that the knowledge of the whole folding process includes the capability of predicting the native conformation, and the problem of protein structure prediction would be automatically solved if ab initio folding were achieved, the design of accurate energy functions, which is a central part of ab initio strategies (see the next section), would also be very helpful to improve knowledge-based methods that make use of them (such as Rosetta [108]) or to prune and refine the candidate models on a second stage [87]. Additionally, to assign the correct conformation to those chain segments that are devoid of secondary structure (the problem known as *loop modeling*), may be considered as a ‘mini protein folding problem’ [87], and the understanding of the physical behaviour of polypeptide chains would also include a solution to this issue. In the light of all these sweet promises, the long ab initio path to study protein folding constitutes an exciting field of present research.

Before we delve deeper in the details, let us define clearly the playfield in which the match shall take place: Although some details of the protein folding process in vivo are under discussion [110] and many cellular processes are involved in helping and checking the arrival to the correct native structure [81]; although some proteins have been shown to fold cotranslationally [111] (i.e., during their synthesis in the ribosome) and many of them are known to be assisted by molecular chaperones (see the discussion above and references therein); although some proteins contain cis proline peptide bonds or disulfide bonded cysteines in their native structure, and must be in the presence of the respective isomerases in order to fold in a reasonable time (see section 2); although

some residues may be post-translationally changed into side chains that are not included in the standard twenty that are depicted in figure 2.5; and, although some non-peptide molecules may be covalently attached to the protein chain or some cofactor or ion may be needed to reach the native structure, we agree with the words by Alan Fersht [112]:

*We can assume that what we learn about the mechanism of folding of small, fast-folding proteins in vitro will apply to their folding in vivo and, to a large extent, to the folding of individual domains in larger proteins.*

and decided to study those processes that do not include any of the aforementioned complications but that may be rightfully considered as intimately related to the process of folding in the cellular milieu and regarded as a first step on top of which to build a more detailed theory.

Henceforth, we define the *restricted protein folding problem*, as the full description of the physical behaviour, in aqueous solvent and physiological conditions, and (consequently) the prediction of the native structure, of completely synthesized proteins, made up just of the twenty genetically encoded amino acids in figure 2.5, without any molecule covalently attached to them, and needless of molecular chaperones, cofactors, ions, disulfide bonds or cis proline peptide bonds in order to fold properly.

Explicitly mentioned or tacitly assumed, it is this restricted version of the problem the one that is most amenable to physics-based methods and the one that is more commonly tackled in the literature.

## 4 Folding mechanisms and energy functions

After having drawn the boundaries of the problem, we should ask the million-dollar question associated to it: *How does a protein fold into its functional native structure?* In fact, since this feat is typically achieved in a very short time, we must add: *How does a protein fold so fast?* This is the question about the *mechanisms* of protein folding, and, ever since Anfinsen’s experiments, it has been asked once and again and only partially answered [76, 89, 109, 113, 114].

In order to define the theoretical framework that is relevant for the description of the folding process and also to introduce the language that is typically used in the discussions about its mechanisms, let us start with a brief reminder of some important statistical mechanics relations. To do this, we will follow the main ideas in reference [115], although the notation and the assumptions regarding the form of the potential energy, as well as some other minor details, will be different. The presentation will be axiomatic and we will restrict ourselves to the situation in which the macroscopic parameters, such as the temperature  $T$  or, say, the number of water molecules  $N_w$ , do not change. In these conditions, that allow us to drop any multiplicative terms in the partition functions or the probabilities, and also to forget any additive terms to the energies, we



can only focus on the conformational preferences of the system (if, for example, the temperature changed, the neglected terms would be relevant and the expressions that one would need to use would be different). For further details or for the more typical point of view in physics, in which the stress is placed in the variation of the macroscopic thermodynamical parameters, see, for example, reference [116].

The system which we will talk about is the one defined by the *restricted protein folding problem* in the previous section, i.e., *one protein surrounded by  $N_w$  water molecules*<sup>25</sup>; however, one must have in mind that all the subsequent reasoning and the derived expressions are exactly the same for a dilute aqueous solution of a macroscopic number of non-interacting proteins.

Now, if classical mechanics is assumed to be obeyed by our system<sup>26</sup>, then each microscopic state is completely specified by the Euclidean<sup>27</sup> coordinates and momenta of the atoms that belong to the protein (denoted by  $x^\mu$  and  $\pi_\mu$ , respectively, with  $\mu = 1, \dots, N$ ) and those belonging to the water molecules (denoted by  $X^m$  and  $\Pi_m$ , with  $m = N + 1, \dots, N + N_w$ ). The whole set of microscopic states shall be called *phase space* and denoted by  $\Gamma \times \Gamma_w$ , explicitly indicating that it is formed as the direct product of the protein phase space  $\Gamma$  and the water molecules one  $\Gamma_w$ .

The central physical object that determines the time behaviour of the system is the *Hamiltonian* (or *energy*) function

$$H(x^\mu, X^m, \pi_\mu, \Pi_m) = \sum_{\mu} \frac{\pi_\mu^2}{2M_\mu} + \sum_m \frac{\Pi_m^2}{2M_m} + V(x^\mu, X^m), \quad (4.1)$$

where  $M_\mu$  and  $M_m$  denote the atomic masses and  $V(x^\mu, X^m)$  is the *potential energy*.

After equilibrium has been attained at temperature  $T$ , the microscopic details about the time trajectories can be forgot and the average behaviour can be described by the laws of statistical mechanics. In the canonical ensemble, the *partition function* [116] of the system, which is the basic object from which the rest of relevant thermodynamical quantities may be extracted, is given by

$$Z = \frac{1}{h^{N+N_w} N_w!} \int_{\Gamma \times \Gamma_w} \exp[-\beta H(x^\mu, X^m, \pi_\mu, \Pi_m)] dx^\mu dX^m d\pi_\mu d\Pi_m, \quad (4.2)$$

---

<sup>25</sup> At this point of the discussion, the possible presence of non-zero ionic strength is considered to be a secondary issue.

<sup>26</sup> Although non-relativistic quantum mechanics may be considered to be a much more precise theory to study the problem, the computer simulation of the dynamics of a system with so many particles using a quantum mechanical description lies far in the future. Nevertheless, this more fundamental theory can be used to design better classical potential energy functions (which is one of the main long-term goals of the research performed in our group).

<sup>27</sup> Sometimes, the term *Cartesian* is used instead of *Euclidean*. Here, we prefer to use the latter since it additionally implies the existence of a mass metric tensor that is proportional to the identity matrix, whereas the *Cartesian* label only asks the  $n$ -tuples in the set of coordinates to be bijective with the abstract points of the space [117].

where  $h$  is Planck's constant, we adhere to the standard notation  $\beta := 1/RT$  (per-mole energy units are used all throughout this work, so  $R$  is preferred over  $k_B$ ) and  $N_w!$  is a combinatorial number that accounts for the quantum indistinguishability of the  $N_w$  water molecules. Additionally, as we have anticipated, the multiplicative factor outside the integral sign is a constant that divides out for any observable averages and represents just a change of reference in the Helmholtz free energy. Therefore, we will drop it from the previous expression and the notation  $Z$  will be kept for convenience.

Next, since the principal interest lies on the conformational behaviour of the polypeptide chain, seeking to develop clearer images and, if possible, reduce the computational demands, water coordinates and momenta are customarily *averaged* (or *integrated*) *out* [115, 118], leaving an *effective Hamiltonian*  $H_{\text{eff}}(x^\mu, \pi_\mu; T)$  that depends only on the protein degrees of freedom and on the temperature  $T$ , and whose potential energy (denoted by  $W(x^\mu; T)$ ) is called *potential of mean force* or *effective potential energy*.

This effective Hamiltonian may be either empirically designed from scratch (which is the common practice in the classical force fields typically used to perform molecular dynamics simulations [106, 107, 119–128]) or obtained from the more fundamental, original Hamiltonian  $H(x^\mu, X^m, \pi_\mu, \Pi_m)$  actually performing the averaging out process. In statistical mechanics, the theoretical steps that must be followed if one chooses this second option are very straightforward (at least formally):

The integration over the water momenta  $\Pi_m$  in equation (4.2) yields a  $T$ -dependent factor that includes the masses  $M_m$  and that shall be dropped by the same considerations stated above. On the other hand, the integration of the water coordinates  $X^m$  is not so trivial, and, except in the case of very simple potentials, it can only be performed formally. To do this, we define the *potential of mean force* or *effective potential energy* by

$$W(x^\mu; T) := -RT \ln \left( \int \exp [-\beta V(x^\mu, X^m)] dX^m \right), \quad (4.3)$$

and simply rewrite  $Z$  as

$$Z = \int_{\Gamma} \exp [-\beta H_{\text{eff}}(x^\mu, \pi_\mu; T)] dx^\mu d\pi_\mu, \quad (4.4)$$

with the *effective Hamiltonian* being

$$H_{\text{eff}}(x^\mu, \pi_\mu; T) = \sum_{\mu} \frac{\pi_{\mu}^2}{2M_{\mu}} + W(x^\mu; T). \quad (4.5)$$

At this point, the protein momenta  $\pi_\mu$  may also be averaged out from the expressions. This choice, which is very commonly taken in the literature, largely simplifies the discussion about the mechanisms of protein folding and the images and metaphors typically used in the field. However, to perform this average is not completely harmless, since it brings up a number of technical and

interpretation-related difficulties mostly due to the fact that the marginal probability density in the  $x^\mu$ -space in equation (4.7) is not invariant under a change of coordinates<sup>28</sup> (see appendix A and reference [52] for further details).

Bearing this in mind, the integration over  $\pi_\mu$  produces a new  $T$ -dependent factor, which is dropped as usual, and yields a new form of the partition function, which is the one that will be used from now on in this section:

$$Z = \int_{\Omega} \exp [-\beta W(x^\mu; T)] dx^\mu, \quad (4.6)$$

where  $\Omega$  now denotes the positions part of the protein phase space  $\Gamma$ .

Some remarks may be done at this point: On the one hand, if one further assumes that the original potential energy  $V(x^\mu, X^m)$  separates as a sum of intra-protein, intra-water and water-protein interaction terms, the effective potential energy  $W(x^\mu; T)$  in the equations above may be written as a sum of two parts: a vacuum intra-protein energy and an effective solvation energy [115]. Nevertheless, this simplification is neither justified a priori, nor necessary for the subsequent reasoning about the mechanisms of protein folding; so it will not be assumed herein.

On the other hand, the (in general, non-trivial) dependence of  $W(x^\mu; T)$  on the temperature  $T$  (see equation (4.3)) and the associated fact that it contains the entropy of the water molecules, justifies its alternative denomination of *internal* or *effective free energy*, and also the suggestive notation  $F(x^\mu) := W(x^\mu; T)$  used in some works [130]. Here, however, we prefer to save the name *free energy* for the one that contains some amount of protein conformational entropy and that may be assigned to finite subsets (states) of the conformational space of the chain (see equation (4.10) and the discussion below).

Finally, we will stick to the notational practice of dropping (but remembering) the temperature  $T$  from  $W$  and  $H_{\text{eff}}$ . This is consistent with the situation of constant  $T$  that we wish to investigate and also very natural and common in the literature. In fact, most Hamiltonian functions (and their respective potentials) that are considered to be ‘fundamental’ actually come from the averaging out of degrees of freedom more microscopical than the ones regarded as relevant, and, as a result, the coupling ‘constants’ contained in them are not really constant, but dependent on the temperature  $T$ .

Now, from the *probability density function* (PDF) in the protein conformational space  $\Omega$ , given by,

$$p(x^\mu) = \frac{\exp [-\beta W(x^\mu)]}{Z}, \quad (4.7)$$

we can tell that  $W(x^\mu)$  completely determines the conformational preferences of the polypeptide chain in the thermodynamic equilibrium as a function of each point of  $\Omega$ . On the opposite extreme of the details scale, we may choose to describe the macroscopic state of the system as a whole (like it is normally

---

<sup>28</sup> Note that, if the momenta  $\pi_\mu$  are kept in the integration measure, any canonical transformation leaves the probability density invariant, since its Jacobian determinant is unity [129].

done in physics [116]) and define, for example, the Helmholtz free energy as  $F := -RT \ln Z$ , where no trace of the microscopic details of the system remains.

In protein science, it is also common practice to take a point of view somewhat in the middle of these two limit descriptions, and define *states* that are neither single points of  $\Omega$  nor the whole set, but finite subsets  $\Omega_i \subset \Omega$  comprising many different conformations that are related in some sense. These states must be precisely specified in order to be of any use, and they must fulfill some reasonable conditions, the most important of which is that they must be mutually exclusive, so that  $\Omega_i \cap \Omega_j = \emptyset, \forall i \neq j$  (i.e., no point can lie in two different states at the same time).

Since the two most relevant conceptual constructions used to think about protein folding, the native ( $\mathcal{N}$ ) and the unfolded ( $\mathcal{U}$ ) states, as well as a great part of the language used to talk about protein stability, fit in this formalism, we will now introduce the basic equations associated to it.

To begin with, one can define the partition function of a certain state  $\Omega_i$  as

$$Z_i := \int_{\Omega_i} \exp[-\beta W(x^\mu; T)] dx^\mu, \quad (4.8)$$

so that the probability of  $\Omega_i$  be given by

$$P_i := \frac{Z_i}{Z}. \quad (4.9)$$

The *Helmholtz free energy*  $F_i$  of this state is

$$F_i := -RT \ln Z_i, \quad (4.10)$$

and the following relation for the free energy differences is satisfied:

$$\Delta F_{ij} = F_j - F_i = -RT \ln \frac{Z_j}{Z_i} = -RT \ln \frac{P_j}{P_i} = -RT \ln \frac{[j]}{[i]} = -RT \ln K_{ij}, \quad (4.11)$$

where  $[i]$  denotes the *concentration* (in chemical jargon) of the species  $i$ , and  $K_{ij}$  is the *reaction constant* (using again images borrowed from chemistry) of the  $i \leftrightarrow j$  equilibrium. It is precisely this dependence on the concentrations, together with the approximate equivalence between  $\Delta F$  and  $\Delta G$  at physiological conditions (where the term  $P\Delta V$  is negligible [115]), that renders equation (4.11) very useful and ultimately justifies this point of view based on states, since it relates the quantity that describes protein stability and may be estimated theoretically (the folding free energy at constant temperature and constant pressure  $\Delta G_{\text{fold}} := G_{\mathcal{N}} - G_{\mathcal{U}}$ ) with the observables that are commonly measured in the laboratory (the concentrations  $[\mathcal{N}]$  and  $[\mathcal{U}]$  of the native and unfolded states) [24, 54, 131].

The next step to develop this state-centred formalism is to define the *microscopic PDF in  $\Omega_i$*  as the original one in equation (4.7) conditioned to the knowledge that the conformation  $x^\mu$  lies in  $\Omega_i$ :

$$p_i(x^\mu) := p(x^\mu | x^\mu \in \Omega_i) = \frac{p(x^\mu)}{P_i} = \frac{\exp[-\beta W(x^\mu)]}{Z_i} . \quad (4.12)$$

Now, using this probability measure in  $\Omega_i$ , we may calculate the *internal energy*  $U_i$  as the average potential energy in this state:

$$U_i := \langle W \rangle_i = \int_{\Omega_i} W(x^\mu) p_i(x^\mu) dx^\mu , \quad (4.13)$$

and also define the *entropy* of  $\Omega_i$  as

$$S_i := -R \int_{\Omega_i} p_i(x^\mu) \ln p_i(x^\mu) dx^\mu . \quad (4.14)$$

Finally, ending our statistical mechanics reminder, one can show that the natural thermodynamic relation among the different state functions is recovered:

$$\Delta F_{ij} = \Delta U_{ij} - T \Delta S_{ij} \simeq \Delta G_{ij} = \Delta H_{ij} - T \Delta S_{ij} , \quad (4.15)$$

where  $H$  is the *enthalpy*, whose differences  $\Delta H_{ij}$  may be approximated by  $\Delta U_{ij}$  neglecting the term  $P \Delta V$  again.

Retaking the discussion about the mechanisms of protein folding, we see (again) in equation (4.7) that the potential of mean force  $W(x^\mu)$  completely determines the conformational preferences of the polypeptide chain in the thermodynamic equilibrium. Nevertheless, it is often useful to investigate also the underlying microscopic dynamics. The effective potential energy  $W(x^\mu)$  in equation (4.3) has been simply obtained in the previous paragraphs using the tools of statistical mechanics; the ‘dynamical averaging out’ of the solvent degrees of freedom in order to describe the *time evolution* of the protein subsystem, on the other hand, is a much more complicated (and certainly different) task [132–136]. However, if the relaxation of the solvent is fast compared to the motion of the polypeptide chain, the function  $W(x^\mu)$  turns out to be precisely the effective ‘dynamical’ potential energy that determines the microscopic time evolution of the protein degrees of freedom [133]. Although this condition could be very difficult to check for real cases and it has only been studied in simplified model systems [132, 134, 136], molecular dynamics simulations with classical force fields and explicit water molecules suggest that it may be approximately fulfilled [133, 137, 138]. For the sake of brevity, in the discussion that follows, we will assume that this fast-relaxation actually occurs, so that, when reasoning about the graphical representations (commonly termed *energy landscapes*) of the effective potential energy  $W(x^\mu)$ , we are entitled to switch back and forth from dynamical to statistical concepts.

Now, just after noting that  $F(x^\mu)$  is the central physical object needed to tackle the elucidation of the folding mechanisms, we realize that the number of degrees of freedom  $N$  in an average-length polypeptide chain is large enough for the size of the conformational space (which is exponential on  $N$ ) to be astronomically astronomical. This fact was, for years, regarded as a problem, and is normally called *Levinthal’s paradox* [139]. Although it belongs to the set of

paradoxes that (like Zeno’s or Epimenides’) are called so without actually being problematic<sup>29</sup>, thinking about it and using the language and the images related to it have dominated the views on folding mechanisms for a long time [89]. The paradox itself was first stated in a talk entitled ‘How to fold gracefully’ given by Cyrus Levinthal in 1969 [140] and it essentially says that, if, in the course of folding, a protein is required to sample all possible conformations (a hypothesis that ignores completely the laws of dynamics and statistical mechanics) and the conformation of a given residue is independent of the conformations of the rest (which is also false), then the protein will never fold to its native structure.

For example, let us assume that each one of the 124 residues in Anfinsen’s ribonuclease (see section 3) can take up any of the six different discrete backbone conformations in table 1 (side chain degrees of freedom are not relevant in this qualitative discussion, since they only affect the structure locally). This makes a total of  $6^{124} \simeq 10^{96}$  different conformations for the chain. If they were visited in the shortest possible time (say,  $\sim 10^{-12}$  s, approximately the time required for a single molecular vibration [141]), the protein would need about  $10^{76}$  years to sample the whole conformational space. Of course, this argument is just a *reductio ad absurdum* proof (since proteins do fold!) of the a priori evident statement that protein folding cannot be a completely random trial-and-error process (i.e., a random walk in conformational space). The *golf-course* energy landscape in figure 4.1a represents this non-realistic, paradoxical situation: the point describing the conformation of the chain wanders aimlessly on the enormous denatured plateau until it suddenly finds the native well by pure chance.

Levinthal himself argued that a solution to his paradox could be that the folding process occurs along well-defined *pathways* that take every protein, like an ordered column of ants, from the *unfolded state* to the native structure, visiting partially folded intermediates en route [114, 142]. The *ant-trail* energy landscape in figure 4.1b is a graphical depiction of the pathway image.

This view, which is typically referred to as the *old view* of folding [133, 139, 143], is largely influenced by the situation in simple chemical reactions, where the barriers surrounding the minimum energy paths that connect the different local minima are very steep compared to  $RT$ , and the dynamical trajectories are, consequently, well defined. In protein folding, however, due to the fact that the principal driving forces are much weaker than those relevant for chemical reactions and comparable to  $RT$ , short-lived transient interactions may form randomly among different residues in the chain and the system describes stochastic trajectories that are never the same. Henceforth, since the native state may be reached in many ways, it is unlikely that a single minimum energy path dominates over the rest of them [133].

In the late 80s, a *new view* of folding mechanisms began to emerge based on these facts and inspired on the statistical mechanics of spin glasses [133, 139, 141, 144–146]. According to it, when a large number of identical proteins (from  $10^{15}$

---

<sup>29</sup> In fact, Levinthal did not use the word ‘paradox’ and, just after stating the problem, he proposed a possible solution to it.

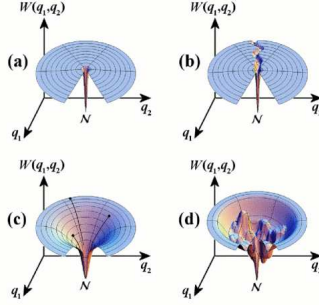


Figure 4.1: Possible energy landscapes of a protein. The conformational space is assumed to be two-dimensional, the degrees of freedom being  $q_1$  and  $q_2$ . The degrees of freedom of the solvent have been *integrated out* (see the text), and the effective potential energy  $W(q_1, q_2)$  is a function of these two variables, which are internal degrees of freedom of the molecule.  $\mathcal{N}$  stands for native state and it is assumed here to be the global minimum. (a) *Flat golf course*: the energy landscape as it would be if Levinthal’s paradox were a real problem. (b) *Ant trail*: the old-view pathway solution to Levinthal’s paradox. (c) *New-view smooth funnel*. (d) *More realistic partially rugged funnel*. (Figures taken from reference [130] with kind permission and somewhat modified.)

to  $10^{18}$  [147]) are introduced in a test tube in the conditions of the *restricted protein folding problem* defined in section 3, a conformational equilibrium is attained between the native ensemble of states  $\mathcal{N}$  and the ensemble made up of the rest of (non-functional) conformations (the unfolded state  $\mathcal{U}$ ). At the same time, what is happening at the microscopic level is that each single molecule is following a partially stochastic trajectory determined by the intrinsic energetics of the system (given by  $W(x^\mu)$ ) and subject to random fluctuations due to the thermal noise. Of course, all trajectories are different, some towards the native state and some towards the unfolded state, but, if we focus on a single molecule at an arbitrary time, the probability that she is wandering in the native basin is very high (typically more than 99%) and, in the rare case that we happen to choose a protein that is presently unfolded, we will most certainly watch a very fast race towards the native state.

In order for this to happen, we need that the energy landscape be *funneled* towards the native state, like in figures 4.1c and 4.1d, so that any microscopic trajectory has more probability to evolve in the native direction than in the opposite one at every point of the conformational space (the ‘ruggedness’ of the funnel must also be small in order to avoid getting trapped in deep local minima during the course of folding). In this way, the solution to Levinthal’s paradox could be said to be ‘funnels, not tunnels’ [148], and the deterministic pathway image is changed by a statistical treatment in which folding is a heterogeneous reaction involving broad ensembles of structures [149], the kinetic intermediates



that are sometimes observed experimentally being simply more or less deep wells in the walls of the funnel. Anyway, although this new view has been validated both experimentally [150] and theoretically [147], and it is widely accepted as correct by the scientific community, one must note that it is not contradictory with the old view, since the latter is only a particular case of the former in which the funnel presents a deep canyon through which most of the individual proteins roll downwards. In fact, in some studied cases, one may find a single pathway that dominates statistically [138, 147].

A marginal issue that arises both in the old and new views, is whether the native state is the global minimum of the effective potential energy  $W(x^\mu)$  of the protein (in which case the folding process is said to be *thermodynamically controlled*) or it is just the lowest-lying kinetically-accessible local minimum (in which case we talk about *kinetic control*) [115]. This question was raised by Anfinsen [76], who assumed the first case to be the correct answer and called the assumption the *thermodynamic hypothesis*. Although Levinthal pointed out a few years later that this was not necessary and that kinetic control was perfectly possible [140], and also despite some indications against it [151, 152], it is now widely accepted that the thermodynamic hypothesis is fulfilled most of the times, and almost always for small single-domain proteins [24, 77, 81, 115]. Of course, nothing fundamental changes in the overall picture if the energy landscape is funneled towards a local minimum of  $W(x^\mu)$  instead of being funneled to a global one, however, from the computational point of view there is a difference: In the latter case, the prediction of the native state may be tackled both dynamically and by simple minimization<sup>30</sup> of the function  $W(x^\mu)$  (for example, using *simulated annealing* [153, 154] or similar schemes), whereas, if the thermodynamic hypothesis is broken, the native structure may still be found performing molecular dynamics simulations, but minimization procedures could be misleading and technically problematic. This is so because, although local minima may also be found and described, the knowledge about towards which one of them the protein trajectories converge depends on kinetic information, which is absent from the typical minimization algorithms.

Now, even though a funneled energy function provides the only consistent image that accounts for all the experimental facts about protein folding, one must still explain the fact that the landscape is just like that. If one looks at a protein as if it were the first time, one sees that it is a heteropolymer made up of twenty different types of amino acid monomers (see section 2). Such a system, due to its many degrees of freedom, the constraints imposed by chain connectivity and the different affinities that the monomers show for their neighbours and for the environment, presents a large degree of *frustration*, that is, there is not a single conformation of the chain which optimizes all the interactions at the same time<sup>31</sup>. For the vast majority of the sequences, this

<sup>30</sup> See appendix A for some technical but relevant remarks about the minimization of the effective potential energy function.

<sup>31</sup> In order to be entitled to give such a simple definition, we need that the effective potential energy of the system separates as a sum of terms with the minima at different points (either because it is split in few-body terms, or because it is split in different ‘types of interactions’,



would lead to a rugged energy landscape with many low-energy states, high barriers, strong traps, etc.; up to a certain degree, a landscape similar to that of spin glasses. A landscape in which fast-folding to a unique three-dimensional structure is impossible!

However, a protein is not a random heteropolymer. Its sequence has been selected and improved along thousands of millions of years by natural selection<sup>32</sup>, and the score function that decided the contest, the fitness that drove the process, is just its ability to fold into a well-defined native structure in a biologically reasonable time<sup>33</sup>. Henceforth, the energy landscape of a protein is not like the majority of them, proteins are a selected minority of heteropolymers for which there exists a privileged structure (the native one) so that, in every point of the conformational space, it is more stabilizing, on average, to form ‘native contacts’ than to form ‘non-native’ ones (an image radically implemented by Gō-type models [156]). Bryngelson and Wolynes [146] have termed this fewer conflicting interactions than typically expected the *principle of minimal frustration*, and this takes us to a natural definition of a *protein* (opposed to a general *polypeptide*): a *protein* is a polypeptide chain whose sequence has been naturally selected to satisfy the principle of minimal frustration.

Now, we should note that this funneled shape emerges from a very delicate balance. Proteins are only marginally stable in solution, with an unfolding free energy  $\Delta G_{\text{unfold}}$  typically in the 5 – 15 kcal/mol range. However, if we split this relatively small value into its enthalpic and entropic contributions, using equation (4.15) and the already mentioned fact that the term  $P\Delta V$  is negligible at physiological conditions [115],

$$\Delta G_{\text{unfold}} = \Delta H_{\text{unfold}} - T\Delta S_{\text{unfold}} , \quad (4.16)$$

we find that it is made up of the difference between two quantities ( $\Delta H_{\text{unfold}}$  and  $T\Delta S_{\text{unfold}}$ ) that are typically an order of magnitude larger than  $\Delta G_{\text{unfold}}$  itself [115,157], i.e., the native state is enthalpically favoured by hundreds of kilocalories per mole and entropically penalized by approximately the same amount.

In addition, both quantities are strongly dependent on the details of the effective potential energy  $W(x^\mu)$  (see equations (4.13) and (4.14)), which could be imagined to be made up of the sum of thousands of non-covalent terms each one of a size comparable to  $\Delta G_{\text{unfold}}$ . This very fine tuning that has been achieved after thousands of millions of years of natural selection is easily destroyed by a

---

such as van der Waals, Coulomb, hydrogen-bonds, etc.). This is a classical image which is rigorously wrong but approximately true (and very useful to think). If one does not want to assume the existence of ‘interactions’ or few-body terms that may conflict with one another, one may jump directly to the conclusion, noting that the energy landscape of a random heteropolymer is glassy but without introducing the concept of frustration.

<sup>32</sup> The problem of finding the protein needle in the astronomical haystack of all possible sequences and its solution are presented as another paradox, *the blind watchmaker paradox*, and inspiring discussed by Richard Dawkins in reference [155].

<sup>33</sup> One may argue that the ability to perform a catalytic function also enters the fitness criterium. While this is true, it is probably a less important factor than the folding skill, since the active site of enzymes is generally localized in a small region of the surface of the protein and it could be, in principle, assembled on top of many different folds.

single-residue mutation or by slightly altering the temperature, the  $pH$  or the concentration of certain substances in the environment (parameters on which  $W(x^\mu)$  implicitly depends).

For the same reasons, if the folding process is intended to be simulated theoretically, the chances of missing the native state and (what is even worse) of producing a non-funneled landscape, which is very difficult to explore using conventional molecular dynamics or minimization algorithms, are very high if poor energy functions are used [144, 158, 159]. Therefore, it is not surprising that current force fields [106, 107, 119–128], which include a number of strong assumptions (additivity of the ‘interactions’, mostly pairwise terms, simple functional forms, etc.), are widely recognized to be incapable of folding proteins [24, 86, 100, 102, 160–163].

The improvement of the effective potential energy functions describing polypeptides, with the long-term goal of reliable *ab initio* folding, is one of the main objectives pursued in our group, and probably one of the central issues that must be solved before the wider framework of the protein folding problem can be tackled. The enormous mathematical and computational complexity that the study of these topics entails, renders the incorporation of the physicists community essential for the future advances in molecular biology. That the boundaries of what is normally considered ‘physics’ are expanding is obvious, and so it is that the investigation of the behaviour of biological macromolecules is a very appealing part of the new territory to explore.

## Acknowledgments

I wish to thank J. L. Alonso, J. Sancho and I. Calvo for illuminating discussions and for the invaluable help to perform the transition mentioned in the title of this work.

This work has been supported by the research projects E24/3 and PM048 (Aragón Government), MEC (Spain) FIS2006-12781-C02-01 and MCyT (Spain) FIS2004-05073-C04-01. P. Echenique and is supported by a BIFI research contract.

## A Probability density functions

Let us define a *stochastic* or *random variable*<sup>34</sup> as a pair  $(X, p)$ , with  $X$  a subset of  $\mathbb{R}^n$  for some  $n$  and  $p$  a function that takes  $n$ -tuples  $x \equiv (x_1, \dots, x_n) \in X$  to positive real numbers,

$$\begin{aligned} p : X &\longrightarrow [0, \infty) \\ x &\longmapsto p(x) \end{aligned}$$

Then,  $X$  is called *range*, *sample space* or *phase space*, and  $p$  is termed *probability distribution* or *probability density function* (PDF). The phase space can be

---

<sup>34</sup> See Van Kampen [164] for a more complete introduction to probability theory.

discrete, a case with which we shall not deal here, or continuous, so that  $p(x) dx$  (with  $dx := dx_1 \cdots dx_n$ ) represents the probability of occurrence of some  $n$ -tuple in the set defined by  $(x, x + dx) := (x_1, x_1 + dx_1) \times \cdots \times (x_n, x_n + dx_n)$ , and the following normalization condition is satisfied:

$$\int_X p(x) dx = 1 . \quad (\text{A.1})$$

It is precisely in the continuous case where the interpretation of the function  $p(x)$  alone is a bit problematic, and playing intuitively with the concepts derived from it becomes dangerous. On one side, it is obvious that  $p(x)$  is not the probability of the value  $x$  happening, since the probability of any specific point in a continuous space must be zero (what is the probability of selecting a random number between 3 and 4 and obtaining *exactly*  $\pi$ ?). In fact, the correct way of using  $p(x)$  to assign probabilities to the  $n$ -tuples in  $X$  is ‘to multiply it by differentials’ and say that it is the probability that any point in a differentially small interval occurs (as we have done in the paragraph above equation (A.1)). The reason for this may be expressed in many ways: one may say that  $p(x)$  is an object that only makes sense under an integral sign (like a Dirac delta), or one may realize that only probabilities of finite subsets of  $X$  can have any meaning. In fact, it is this last statement the one that focuses the attention on the fact that, if we decide to reparameterize  $X$  and perform a change of variables  $x'(x)$ , what should not change are the integrals over finite subsets of  $X$ , and, therefore,  $p(x)$  cannot transform as a scalar quantity (i.e., satisfying  $p'(x') = p(x(x'))$ ), but according to a different rule.

If we denote the *Jacobian matrix* of the change of variables by  $\partial x / \partial x'$ , we must have that

$$p'(x') = \left| \det \left( \frac{\partial x}{\partial x'} \right) \right| p(x(x')) , \quad (\text{A.2})$$

so that, for any finite set  $Y \subset X$  (with its image by the transformation denoted by  $Y'$ ), and indicating the probability of a set with a capital  $P$ , we have the necessary property

$$P(Y) := \int_Y p(x) dx = \int_{Y'} p'(x') dx' =: P'(Y') . \quad (\text{A.3})$$

All in all, the object that has meaning content is  $P$  and not  $p$ . If one needs to talk about things such as the *most probable regions*, or *the most probable states*, or *the most probable points*, or if one needs to compare in any other way the relative probabilities of different parts of the phase space  $X$ , an *arbitrary* partition of  $X$  into finite subsets  $(X_1, \dots, X_i, \dots)$  must be defined<sup>35</sup>. These  $X_i$  should be considered more useful *states* than the individual points  $x \in X$  and their probabilities  $P(X_i)$ , which, contrarily to  $p(x)$ , do not depend on the

---

<sup>35</sup> Two additional reasonable properties should be asked to such a partition: (i) the sets in it must be exclusive, i.e.,  $X_i \cap X_j = \emptyset, \forall i \neq j$ , and (ii) they must fill the phase space,  $\bigcup_i X_i = X$

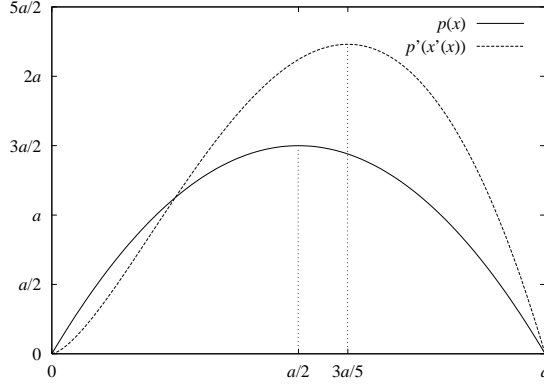


Figure A.1: Probability density functions  $p(x)$  and  $p'(x'(x))$  in equations (A.4) and (A.5) respectively. In the axes, the quantities  $x$  and  $p(x)$  are shown for convenience. Note that the area enclosed by the two curves is different; this is because  $p'(x'(x))$  is normalized with the measure  $dx'$  and not with  $dx$ , which is the one implicitly assumed in this representation.

coordinates chosen, should be used as the meaningful quantities about which to make well-defined probabilistic statements.

To illustrate this, let us see an example: suppose we have a 1-dimensional PDF

$$p(x) = \frac{6}{a^3} x(a-x) . \quad (\text{A.4})$$

The maximum of  $p(x)$  is at  $x = a/2$ , however, it would not be very clever to declare that  $x = a/2$  is the most probable value of  $x$ , since one may choose to describe the problem with a different but perfectly legitimate variable  $x'$ , whose relation to  $x$  is, say,  $x = x'^2$ , and find the PDF in terms of  $x'$  using equation (A.2):

$$p'(x') = \frac{12}{a^3} x'^3 (a - x'^2) . \quad (\text{A.5})$$

Now, insisting on the mistake, we may find the maximum of  $p'(x')$ , which lies at  $x' = (3a/5)^{1/2}$  (see figure A.1), and declare it the most probable value of  $x'$ . But, according to the change of variables given by  $x = x'^2$ , the point  $x' = (3a/5)^{1/2}$  corresponds to  $x = 3a/5$  and, certainly, it is not possible that  $x = a/2$  and  $x = 3a/5$  are the most probable values of  $x$  at the same time!

To sum up, only finite regions of continuous phase spaces can be termed *states* and meaningfully assigned a probability that do not depend on the coordinates chosen. In order to do that, an *arbitrary* partition of the phase space must be defined.

Far from being an academic remark, this is relevant in the study of the equilibrium of proteins, where, very commonly, Anfinsen's *thermodynamic hypothesis* is invoked (see section 4). Loosely speaking, it says that *the functional native state of proteins lies at the minimum of the effective potential energy* (i.e., the maximum of the associated Boltzmann PDF, proportional to  $e^{-\beta W}$ , in equation (4.7)), but, according to the properties of PDFs described in the previous paragraphs, much more qualifying is needed.

First, one must note that all complications arise from the choice of integrating out the momenta (for example, in equation (4.6)) to describe the equilibrium distribution of the system with a PDF dependent only on the potential energy. If the momenta were kept and the PDF expressed in terms of the complete Hamiltonian as  $p(q^\mu, \pi_\mu) = e^{-\beta H}/Z$ , then, it would be invariant under canonical changes of coordinates (which are the physically allowed ones), since the Jacobian determinant that appears in equation (A.2) equals unity in such a case. If we now look, using this complete description in terms of  $H$ , for the *most probable point*  $(q^\mu, \pi_\mu)$  in the whole dynamical phase space, the answer does not depend on the coordinates chosen: It is the point with all momenta  $\pi_\mu$  set to zero (since the kinetic energy is a positive defined quadratic form on the  $\pi_\mu$ ), and the positions  $q^\mu$  set to those that minimize the potential energy  $V(q^\mu)$ , denoted by  $q_{\min}^\mu$ . If we now perform a point transformation, which is a particular case of the larger group of canonical transformations [165],

$$q^\mu \rightarrow q'^\mu(q^\mu) \quad \text{and} \quad \pi_\mu \rightarrow \pi'_\mu = \frac{\partial q^\nu}{\partial q'^\mu} \pi_\nu, \quad (\text{A.6})$$

the *most probable point* in the new coordinates turns out to be ‘the same one’, i.e., the point  $(q'^\mu, \pi'_\mu) = (q'^\mu(q_{\min}^\mu), 0)$ , and all the insights about the problem are consistent.

However, if one decides to integrate out the momenta, the marginal PDF on the positions that remains has a more complicated meaning than the joint one on the whole phase space and lacks the reasonable properties discussed above. The central issue is that the marginal  $p(q^\mu)$  (for example, the one in equation (4.7)) quantifies the probability that the positions of the system be in the interval  $(q^\mu, q^\mu + dq^\mu)$  *without any knowledge about the momenta*, or, otherwise stated, *for any value of the momenta*.

In Euclidean coordinates, the volume in momenta space does not depend on the positions, however, in general curvilinear coordinates, the accessible momenta volume is different from point to point, and one can say the same about the *kinetic entropy* (see reference [52]) associated with the removed  $\pi_\mu$ , which, apart from the potential energy, also enters the coordinate PDF.

If, despite these inconveniences, the description in terms of only the positions  $q^\mu$  is chosen to be kept (which is typically recommendable from the computational point of view), two different approaches may be followed to assure the meaningfulness of the statements made: Either some partition of the conformational space into finite subsets must be defined, as it is described in the beginning of this appendix and as it is done in reference [118], or the position-dependent kinetic entropies that appear when curvilinear coordinates are used and that are

introduced in reference [52] must be included in the effective potential energy function.

## References

- [1] A. M. LESK, *Introduction to Protein Architecture*, Oxford University Press, Oxford, 2001.
- [2] Y. CHO, S. GORINA, P. D. JEFFREY, and N. P. PAVLETICH, *Crystal structure of a p53 tumor suppressor-DNA complex: Understanding tumorigenic mutations*, Science **265** (1994) 346–355.
- [3] C. M. DOBSON, *Protein-misfolding diseases: Getting out of shape*, Nature **729** (2002) 729.
- [4] J. W. KELLY, *Towards and understanding of amyloidogenesis*, Nat. Struct. Biol. **9** (2002) 323.
- [5] E. H. KOO, P. T. LANSBURY JR., and J. W. KELLY, *Amyloid diseases: Abnormal protein aggregation in neurodegeneration*, Proc. Natl. Acad. Sci. USA **96** (1999) 9989–9990.
- [6] M. H. NIELSEN, F. S. PEDERSEN, and J. KJEMS, *Molecular strategies to inhibit HIV-1 replication*, Retrovirology **2** (2005) 10.
- [7] H. OHTAKA and E. FREIRE, *Adaptive inhibitors of the HIV-1 protease*, Prog. Biophys. Mol. Biol. **88** (2005) 193–208.
- [8] U. BACHA, J. BARRILA, A. VELÁZQUEZ-CAMPOY, S. LEAVITT, and E. FREIRE, *Identification of novel inhibitors of the SARS associated coronavirus main proteinase 3CLpro*, Biochemistry **43** (2004) 4906–4912.
- [9] S. VENKATRAMAN, F. G. NJORGE, V. M. GIRIJAVALLABHAN, V. S. MADISON, N. H. YAO, A. J. PRONGAY, N. BUTKIEWICZ, and J. PICHARDO, *Design and synthesis of depeptidized macrocyclic inhibitors of Hepatitis C NS3-4A protease using structure-based drug design*, J. Med. Chem. **48** (2005) 5088–5091.
- [10] J. POEHLGAARD and S. DOUTHWAITE, *The bacterial ribosome as a target for antibiotics*, Nat. Rev. Microbiol. **3** (2005) 870–881.
- [11] C. SMITH, *Drug target validation: Hitting the target*, Nature **422** (2003) 341–347.
- [12] F. VOLLRATH and D. PORTER, *Spider silk as archetypal protein elastomer*, Soft Matter **2** (2006) 377–385.
- [13] J. SCHELLER, K.-H. GÜHRS, F. GROSSE, and U. CONRAD, *Production of spider silk proteins in tobacco and potato*, Nat. Biotech. **19** (2001) 573–577.
- [14] C. M. BELLINGHAM and F. W. KEELEY, *Self-ordered polymerization of elastin-based biomaterials*, Curr. Opin. Sol. State Mat. Sci. **8** (2004) 135–139.

- [15] A. Y. WANG, X. MO, C. S. CHEN, and S. M. YU, *Facile modification of collagen directed by collagen mimetic peptides*, J. Am. Chem. Soc. **127** (2005) 4130–4131.
- [16] S. A. MASKARINEC and D. A. TIRRELL, *Protein engineering approaches to biomaterials design*, Curr. Opin. Biotech. **16** (2005) 422–426.
- [17] M. STRONG, *Protein nanomachines*, PLoS Biology **2** (2004) 0305.
- [18] MANY AUTHORS, *What don't we know?*, Science **309** (2005) 78–102.
- [19] L. D. STEIN, *Human genome: End of the beginning*, Nature **431** (2004) 915–916.
- [20] A. WOOLFE, M. GOODSON, D. K. GOODE, P. SNELL, G. K. MCEWEN, T. VAVOURI, S. F. SMITH, P. NORTH, H. CALLAWAY, K. KELLY, K. WALTER, I. ABNIZOVA, W. GILKS, Y. J. K. EDWARDS, J. E. COOKE, and G. ELGAR, *Highly conserved non-coding sequences are associated with vertebrate development*, PLoS Biology **3** (2005) 0116.
- [21] M. A. NOBREGA, I. OVCHARENKO, V. AFZAL, and E. M. RUBIN, *Scanning human gene deserts for long-range enhancers*, Science **302** (2003) 413.
- [22] W. W. GIBBS, *The unseen genome: Gems among the junk*, Scientific American **289** (2003) 46–53.
- [23] G. GIBSON and S. V. MUSE, *A Primer of Genome Science*, Sinauer, Sunderland, 2nd edition, 2004.
- [24] C. GÓMEZ-MORENO CALERA and J. SANCHO SANZ, editors, *Estructura de Proteínas*, Ariel ciencia, Barcelona, 2003.
- [25] B. BOECKMANN, A. BAIROCH, R. APWEILER, M.-C. BLATTER, A. ESTREICHER, E. GASTEIGER, M. J. MARTIN, K. MICHOD, C. O'DONOVAN, I. PHAN, S. PILBOUT, and M. SCHNEIDER, *The SWISS-PROT protein knowledgebase and its supplement TrEMBL in 2003*, Nucleic Acids Research **31** (2003) 365–370.
- [26] S. ORCHARD, H. HERMJAKOB, and R. APWEILER, *Annotating the human proteome*, Mol. Cell. Proteomics **4** (2005) 435–440.
- [27] J. T. PELTON and L. R. MCLEAN, *Spectroscopic methods for analysis of protein secondary structure*, Anal. Biochem. **277** (2000) 167–176.
- [28] D. A. CASE, H. J. DYSON, and P. E. WRIGHT, *Use of chemical shifts and coupling constants in nuclear magnetic resonance structural studies on peptides and proteins*, Methods Enzymol. **239** (1994) 392–416.
- [29] G. J. KLEYWEGT, *Validation of protein crystal structures*, Acta Crystallogr. D **56** (2000) 249–265.
- [30] J. DRENTH, *Principles of Protein X-Ray Crystallography*, Springer-Verlag, New York, 1999.
- [31] J. P. GLUSKER, *X-ray crystallography of proteins*, Methods Biochem. Anal. **37** (1994) 1–72.



- [32] H. M. BERMAN, J. WESTBROOK, Z. FENG, G. GILLILAND, T. N. BHAT, H. WEISSIG, I. N. SHINDYALOV, and P. E. BOURNE, *The Protein Data Bank*, Nucleic Acids Research **28** (2000) 235–242.
- [33] Millennium Ecosystem Assessment. *Ecosystems and Human Well-being: Biodiversity Synthesis*, Washington, DC, 2005, <http://www.millenniumassessment.org/en/Products.Synthesis.aspx>.
- [34] T. CASTRIGIANÒ, P. D’ONORIO DE MEO, D. COZZETTO, I. G. TALAMO, and A. TRAMONTANO, *The PMDB Protein Model Database*, Nucleic Acids Research **34** (2006) 306–309.
- [35] A. TRAMONTANO, *Of men and machines*, Nat. Struct. Biol. **10** (2003) 87–90.
- [36] M. D. S. KUMAR, K. A. BAVA, M. M. GROMIHA, P. PRABAKARAN, K. KITAJIMA, H. UEDAIRA, and A. SARAI, *ProTherm and ProNIT: Thermodynamic databases for proteins and protein-nucleic acid interactions*, Nucleic Acids Research **34** (2006) 204–206.
- [37] T.-Y. LEE, H.-D. HUANG, J.-H. HUNG, H.-Y. HUANG, Y.-S. YANG, and T.-H. WANG, *dbPTM: An information repository of protein post-translational modification*, Nucleic Acids Research **34** (2006) 622–627.
- [38] M. D. S. KUMAR and M. M. GROMIHA, *PINT: Protein-protein Interactions Thermodynamic Database*, Nucleic Acids Research **34** (2006) 195–198.
- [39] R. DAWKINS, *River Out of Eden: A Darwinian View of Life*, BasicBooks, New York, 1995.
- [40] M. J. FRISCH, G. W. TRUCKS, H. B. SCHLEGEL, G. E. SCUSERIA, M. A. ROBB, J. R. CHEESEMAN, J. A. MONTGOMERY, JR., T. VREVEN, K. N. KUDIN, J. C. BURANT, J. M. MILLAM, S. S. IYENGAR, J. TOMASI, V. BARONE, B. MENNUECCI, M. COSSI, G. SCALMANI, N. REGA, G. A. PETERSSON, H. NAKATSUJI, M. HADA, M. EHARA, K. TOYOTA, R. FUKUDA, J. HASEGAWA, M. ISHIDA, T. NAKAJIMA, Y. HONDA, O. KITAO, H. NAKAI, M. KLENE, X. LI, J. E. KNOX, H. P. HRATCHIAN, J. B. CROSS, V. BAKKEN, C. ADAMO, J. JARAMILLO, R. GOMPERTS, R. E. STRATMANN, O. YAZYEV, A. J. AUSTIN, R. CAMMI, C. POMELLI, J. W. OCHTERSKI, P. Y. AYALA, K. MOROKUMA, G. A. VOTH, P. SALVADOR, J. J. DANNENBERG, V. G. ZAKRZEWSKI, S. DAPPRICH, A. D. DANIELS, M. C. STRAIN, O. FARKAS, D. K. MALICK, A. D. RABUCK, K. RAGHAVACHARI, J. B. FORESMAN, J. V. ORTIZ, Q. CUI, A. G. BABOUL, S. CLIFFORD, J. CIOŚLOWSKI, B. B. STEFANOV, G. LIU, A. LIASHENKO, P. PISKORZ, I. KOMAROMI, R. L. MARTIN, D. J. FOX, T. KEITH, M. A. AL-LAHAM, C. Y. PENG, A. NANAYAKKARA, M. CHALLACOMBE, P. M. W. GILL, B. JOHNSON, W. CHEN, M. W. WONG, C. GONZALEZ, and J. A. POPLE, *Gaussian 03, Revision C.02*, 2004, Gaussian, Inc., Wallingford, CT.
- [41] R. S. CAHN, S. C. INGOLD, and V. PRELOG, *Specification of molecular chirality*, Angew. Chem. Int. Ed. **5** (1966) 385–415.
- [42] M. KLUSMANN, H. IWAMURA, S. P. MATHEW, D. H. WELLS JR., U. PANDYA, A. ARMSTRONG, and D. G. BLACKMOND, *Thermodynamic control of asymmetric amplification in amino acid catalysis*, Nature **441** (2006) 621–623.

- [43] T. HEAD-GORDON, M. HEAD-GORDON, M. J. FRISCH, C. BROOKS III, and J. POPLE, *A theoretical study of alanine dipeptide and analogs*, Intl. J. Quant. Chem. **16** (1989) 311-322.
- [44] R. F. FREY, J. COFFIN, S. Q. NEWTON, M. RAMEK, V. K. W. CHENG, F. A. MOMANY, and L. SCHÄFER, *Importance of correlation-gradient geometry optimization for molecular conformational analyses*, J. Am. Chem. Soc. **114** (1992) 5369-5377.
- [45] I. R. GOULD, W. D. CORNELL, and I. H. HILLIER, *A quantum mechanical investigation of the conformational energetics of the alanine and glycine dipeptides in the gas phase and in aqueous solution*, J. Am. Chem. Soc. **116** (1994) 9250-9256.
- [46] M. BEACHY, D. CHASMAN, R. MURPHY, T. HALGREN, and R. FRIESNER, *Accurate ab initio quantum chemical determination of the relative energetics of peptide conformations and assessment of empirical force fields*, J. Am. Chem. Soc. **119** (1997) 5908-5920.
- [47] C.-H. YU, M. A. NORMAN, L. SCHÄFER, M. RAMEK, A. PEETERS, and C. VAN ALSENOY, *Ab initio conformational analysis of N-formyl L-alanine amide including electron correlation*, J. Mol. Struct. **567-568** (2001) 361-374.
- [48] R. VARGAS, J. GARZA, B. P. HAY, and D. A. DIXON, *Conformational study of the alanine dipeptide at the MP2 and DFT levels*, J. Phys. Chem. A **106** (2002) 3213-3218.
- [49] A. PERCZEL, Ö. FARKAS, I. JÁKLI, I. A. TOPOL, and I. G. CSIZMADIA, *Peptide models. XXXIII. Extrapolation of low-level Hartree-Fock data of peptide conformation to large basis set SCF, MP2, DFT and CCSD(T) results. The Ramachandran surface of alanine dipeptide computed at various levels of theory*, J. Comp. Chem. **24** (2003) 1026-1042.
- [50] Z.-X. WANG and Y. DUAN, *Solvation effects on alanine dipeptide: A MP2/cc-pVTZ//MP2/6-31G\*\* study of  $(\Phi, \Psi)$  energy maps and conformers in the gas phase, ether and water*, J. Comp. Chem. **25** (2004) 1699-1716.
- [51] P. ECHENIQUE and J. L. ALONSO, *Definition of Systematic, Approximately Separable and Modular Internal Coordinates (SASMIC) for macromolecular simulation*, J. Comp. Chem. **27** (2006) 1076-1087.
- [52] P. ECHENIQUE, I. CALVO, and J. L. ALONSO, *Quantum mechanical calculation of the effects of stiff and rigid constraints in the conformational equilibrium of the Alanine dipeptide*, J. Comp. Chem. **27** (2006) 1748-1755.
- [53] P. ECHENIQUE and J. L. ALONSO, *Efficient model chemistries for peptides. I. Split-valence Gaussian basis sets and the heterolevel approximation in RHF and MP2*, Submitted, 2007.
- [54] T. E. CREIGHTON, *Proteins: Structures and Molecular Properties*, Freeman, W. H., New York, 2nd edition, 1992.
- [55] M. DI GIULIO, *The origin of the genetic code: Theories and their relationships, a review*, Biosystems **80** (2005) 175-184.

- [56] R. D. KNIGHT, S. J. FREELAND, and L. F. LANDWEBER, *Selection, history and chemistry: The three faces of the genetic code*, Trends Biochem. Sci. **24** (1999) 241–247.
- [57] M. S. WEISS, A. JABS, and R. HILGENFELD, *Peptide bonds revisited*, Nat. Struct. Biol. **5** (1998) 676.
- [58] J. F. SWAIN and L. M. GIERASH, *A new twist for an Hsp70 chaperone*, Nat. Struct. Biol. **9** (2002) 406–408.
- [59] V. I. LIM and A. S. SPIRIN, *Stereochemical analysis of ribosomal transpeptidation conformation of nascent peptide*, J. Mol. Biol. **188** (1986) 565–574.
- [60] G. N. RAMACHANDRAN and C. RAMAKRISHNAN, *Stereochemistry of polypeptide chain configurations*, J. Mol. Biol. **7** (1963) 95–99.
- [61] L. BRAGG, J. C. KENDREW, and M. F. PERUTZ, *Polypeptide chain configurations in crystalline proteins*, Proc. Roy. Soc. London Ser. A **203** (1950) 321–357.
- [62] L. PAULING, R. B. COREY, and H. R. BRANSON, *The structure of proteins: Two hydrogen-bonded helical configurations of the polypeptide chain*, Proc. Natl. Acad. Sci. USA **37** (1951) 205–211.
- [63] D. EISENBERG, *The discovery of the  $\alpha$ -helix and  $\beta$ -sheet, the principal structural features of proteins*, Proc. Natl. Acad. Sci. USA **100** (2003) 11207–11210.
- [64] J. C. KENDREW, G. BODO, H. M. DINTZIS, R. G. PARRISH, H. WYCKOFF, and D. C. PHILLIPS, *A three-dimensional model of the myoglobin molecule obtained by x-ray analysis*, Nature **181** (1958) 662–666.
- [65] M. F. PERUTZ, M. G. ROSSMAN, A. F. CULLIS, H. MUIRHEAD, G. WILL, and A. C. T. NORTH, *Structure of haemoglobin: A three-dimensional Fourier synthesis at 5.5 Å resolution, obtained by x-ray analysis*, Nature **185** (1960) 416–422.
- [66] M. N. FODJE and S. AL-KARADAGHI, *Occurrence, conformational features and amino acid propensities for the  $\pi$ -helix*, Protein Eng. **15** (2002) 353–358.
- [67] M. E. KARPEN, P. L. DE-HASETH, and K. E. NEET, *Differences in the amino acid distributions of  $3_{10}$ -helices and  $\alpha$ -helices*, Prot. Sci. **1** (1992) 1333–1342.
- [68] E. N. BAKER and R. E. HUBBARD, *Hydrogen bonding in globular proteins*, Prog. Biophys. Mol. Biol. **44** (1984) 97–179.
- [69] M. L. HIGGINS, *The structure of fibrous proteins*, Chem. Rev. **32** (1943) 195–218.
- [70] B. W. LOW and R. B. BAYBUTT, *The  $\pi$  helix – A hydrogen bonded configuration of the polypeptide chain*, J. Am. Chem. Soc. **74** (1952) 5806–5807.
- [71] J. DONOHUE, *Hydrogen bonded helical configurations of the polypeptide chain*, Proc. Natl. Acad. Sci. USA **39** (1953) 470–478.

- [72] B. ZAGROVIC, J. LIPFERT, E. J. SORIN, I. S. MILLETT, W. F. VAN GUNSTEREN, S. DONIACH, and V. S. PANDE, *Unusual compactness of a polyproline type II structure*, Proc. Natl. Acad. Sci. USA **102** (2005) 11698–11703.
- [73] R. V. PAPPU and G. D. ROSE, *A simple model for polyproline II structure in unfolded states of alanine-based peptides*, Prot. Sci. **11** (2002) 2437–2455.
- [74] B. J. STAPLEY and T. P. CREAMER, *A survey of left-handed polyproline II helices*, Prot. Sci. **8** (1999) 587–595.
- [75] J. C. KENDREW, *Myoglobin and the structure of proteins*, Nobel Lecture, 1962, [http://nobelprize.org/nobel\\_prizes/chemistry/laureates/1962/](http://nobelprize.org/nobel_prizes/chemistry/laureates/1962/).
- [76] C. B. ANFINSSEN, *Principles that govern the folding of protein chains*, Science **181** (1973) 223–230.
- [77] C. M. DOBSON, *The nature and significance of protein folding*, in *Mechanism of Protein Folding*, edited by R. H. PAIN, pp. 1–33, Oxford University Press, New York, 2000.
- [78] F. U. HARTL and M. HAYER-HARTL, *Molecular chaperones in the cytosol: from nascent chain to folded protein*, Science **295** (2002) 1852–1858.
- [79] F. U. HARTL, *Molecular chaperones in cellular protein folding*, Nature **381** (2002) 571–580.
- [80] J. ELLIS, *Proteins as molecular chaperones*, Nature **328** (1987) 378–379.
- [81] C. M. DOBSON, *Protein folding and misfolding*, Nature **426** (2003) 884–890.
- [82] A. L. HORWICH, E. U. WEBER-BAN, and D. FINLEY, *Chaperone rings in protein folding and degradation*, Proc. Natl. Acad. Sci. USA **96** (1999) 11033–11040.
- [83] C. HARDIN, T. V. POGORELOV, and Z. LUTHEY-SCHULTEN, *Ab initio protein structure prediction*, Curr. Opin. Struct. Biol. **12** (2002) 176–181.
- [84] Y. DUAN and P. A. KOLLMAN, *Computational protein folding: From lattice to all-atom*, IBM Systems Journal **40** (2001) 297–309.
- [85] E. KRIEGER, S. B. NABUURS, and G. VRIEND, *Homology modeling*, in *Structural Bioinformatics*, edited by P. E. BOURNE and H. WEISSIG, pp. 507–521, Wiley-Liss, 2003.
- [86] K. GINALSKI, N. V. GRISHIN, A. GODZIK, and L. RYCHLEWSKI, *Practical lessons from protein structure prediction*, Nucleic Acids Research **33** (2005) 1874–1891.
- [87] M. JACOBSON and A. SALI, *Comparative protein structure modeling and its applications to drug discovery*, Ann. Rep. Med. Chem. **39** (2004) 259–276.
- [88] V. DAGGETT and A. FERSHT, *The present view of the mechanism of protein folding*, Nat. Rev. Mol. Cell Biol. **4** (2003) 497.
- [89] B. HONIG, *Protein folding: From the Levinthal paradox to structure prediction*, J. Mol. Biol. **293** (1999) 283–293.

- [90] D. BAKER and A. SALI, *Protein structure prediction and structural genomics*, Science **294** (2001) 93–96.
- [91] M. A. MARTI-RENO, A. C. STUART, A. FISER, R. SÁNCHEZ, F. MELO, and A. SALI, *Comparative protein structure modeling of genes and genomes*, Annu. Rev. Biophys. Biomol. Struct. **29** (2000) 291–325.
- [92] C. CHOTHIA and A. M. LESK, *The relation between the divergence of sequence and structure in proteins*, EMBO J. **5** (1986) 823–826.
- [93] S. F. ALTSCHUL, T. L. MADDEN, A. A. SCHÄFFER, J. ZHANG, Z. ZHANG, W. MILLER, and L. D. J., *Gapped BLAST and PSI-BLAST: A new generation of protein database search programs*, Nucleic Acids Research **25** (1997) 3389–33402.
- [94] S. HENIKOFF, *Scores for sequence searches and alignments*, Curr. Opin. Struct. Biol. **6** (1996) 352–360.
- [95] U. PIEPER, N. ESWAR, F. P. DAVIS, H. BRABERG, M. S. MADHUSUDHAN, A. ROSSI, M. MARTI-RENO, R. KARCHIN, B. M. WEBB, D. ERAMIAN, M.-Y. SHEN, L. KELLY, F. MELO, and A. SALI, *MODBASE: A database of annotated comparative protein structure models and associated resources*, Nucleic Acids Research **34** (2006) 291–295.
- [96] J. U. BOWIE, R. LUTHY, and D. EISENBERG, *A method to identify protein sequences that fold into a known three-dimensional structure*, Science **253** (1991) 164–170.
- [97] J. MOULT, K. FIDELIS, B. ROST, T. HUBBARD, and A. TRAMONTANO, *Critical Assessment of Methods of Protein Structure Prediction (CASP)—Round 6*, PROTEINS: Struct. Funct. Bioinf. **7** (2005) 3–7.
- [98] C. A. ORENGO and J. M. THORNTON, *Protein families and their evolution—A structural perspective*, Annu. Rev. Biochem. **74** (2005) 867–900.
- [99] P. BRADLEY, K. M. S. MISURA, and D. BAKER, *Toward high-resolution de novo structure prediction for small proteins*, Science **309** (2005) 1868–1871.
- [100] O. SCHUELER-FURMAN, C. WANG, P. BRADLEY, K. MISURA, and D. BAKER, *Progress in modeling of protein structures and interactions*, Science **310** (2005) 638–642.
- [101] J. MOULT, *A decade of CASP: Progress, bottlenecks and prognosis in protein structure prediction*, Curr. Opin. Struct. Biol. **15** (2005) 285–289.
- [102] R. BONNEAU and D. BAKER, *Ab initio protein structure prediction: Progress and prospects*, Annu. Rev. Biophys. Biomol. Struct. **30** (2001) 173–189.
- [103] C. A. ROHL, C. E. STRAUSS, K. M. MISURA, and D. BAKER, *Protein structure prediction using Rossetta*, Methods Enzymol. **383** (2004) 66–93.
- [104] C. A. ROHL, C. E. STRAUSS, D. CHIVIAN, and D. BAKER, *Modeling structurally variable regions in homologous proteins with Rossetta*, PROTEINS: Struct. Funct. Bioinf. **55** (2004) 656–677.

- [105] M. LEVITT and A. WARSHEL, *Computer simulation of protein folding*, Nature **253** (1975) 694–698.
- [106] A. D. MAC KERELL JR., B. BROOKS, C. L. BROOKS III, L. NILSSON, B. ROUX, Y. WON, and M. KARPLUS, *CHARMM: The energy function and its parameterization with an overview of the program*, in *The Encyclopedia of Computational Chemistry*, edited by P. v. R. SCHLEYER et al., pp. 217–277, John Wiley & Sons, Chichester, 1998.
- [107] B. R. BROOKS, R. E. BRUCCOLERI, B. D. OLAFSON, D. J. STATES, S. SWAMINATHAN, and M. KARPLUS, *CHARMM: A program for macromolecular energy, minimization, and dynamics calculations*, J. Comp. Chem. **4** (1983) 187–217.
- [108] D. BAKER, *Prediction and design of macromolecular structures and interactions*, Phil. Trans. R. Soc. London B Biol. Sci. **361** (2006) 459–463.
- [109] J. SKOLNICK, *Putting the pathway back into protein folding*, Proc. Natl. Acad. Sci. USA **102** (2005) 2265–2266.
- [110] M. A. BASHAROV, *Protein folding*, J. Cell. Mol. Med. **7** (2003) 223–237.
- [111] B. HARDESTY and G. KRAMER, *Folding of a nascent peptide on the ribosome*, Prog. Nucleic Acid Res. Mol. Biol. **66** (2001) 41–66.
- [112] A. R. FERSHT and V. DAGGETT, *Protein folding and unfolding at atomic resolution*, Cell **108** (2002) 573–582.
- [113] V. DAGGETT and A. R. FERSHT, *Is there a unifying mechanism for protein folding?*, Trends Biochem. Sci. **28** (2003) 18–25.
- [114] C. LEVINTHAL, *Are there pathways for protein folding?*, J. Chim. Phys. **65** (1968) 44–45.
- [115] T. LAZARIDIS and M. KARPLUS, *Thermodynamics of protein folding: a microscopic view*, Biophys. Chem. **100** (2003) 367–395.
- [116] W. GREINER, H. STOCKER, and L. NEISE, *Thermodynamics and Statistical Mechanics*, Classical Theoretical Physics, Springer, New York, 2004.
- [117] B. A. DUBROVIN, A. T. FOMENKO, and S. P. NOVIKOV, *Modern Geometry — Methods and Applications*, volume I. The Geometry of Surfaces, Transformation Groups and Fields, Springer, New York, 1984.
- [118] T. LAZARIDIS and M. KARPLUS, *Effective energy function for proteins in solution*, PROTEINS: Struct. Funct. Gen. **35** (1999) 133–152.
- [119] W. F. VAN GUNSTEREN and M. KARPLUS, *Effects of constraints on the dynamics of macromolecules*, Macromolecules **15** (1982) 1528–1544.
- [120] W. D. CORNELL, P. CIEPLAK, C. I. BAYLY, I. R. GOULD, J. MERZ, K. M., D. M. FERGUSON, D. C. SPELLMEYER, T. FOX, J. W. CALDWELL, and P. A. KOLLMAN, *A second generation force field for the simulation of proteins, nucleic acids, and organic molecules*, J. Am. Chem. Soc. **117** (1995) 5179–5197.

- [121] D. A. PEARLMAN, D. A. CASE, J. W. CALDWELL, W. R. ROSS, T. E. CHEATHAM III, S. DEBOLT, D. FERGUSON, G. SEIBEL, and P. KOLLMAN, *AMBER, a computer program for applying molecular mechanics, normal mode analysis, molecular dynamics and free energy calculations to elucidate the structures and energies of molecules*, Comp. Phys. Commun. **91** (1995) 1–41.
- [122] W. L. JORGENSEN and J. TIRADO-RIVES, *The OPLS potential functions for proteins. Energy minimization for crystals of cyclic peptides and Crambin*, J. Am. Chem. Soc. **110** (1988) 1657–1666.
- [123] W. L. JORGENSEN, D. S. MAXWELL, and J. TIRADO-RIVES, *Development and testing of the OPLS all-atom force field on conformational energetics and properties of organic liquids*, J. Am. Chem. Soc. **118** (1996) 11225–11236.
- [124] T. A. HALGREN, *Merck Molecular Force Field. I. Basis, form, scope, parametrization, and performance of MMFF94*, J. Comp. Chem. **17** (1996) 490–519.
- [125] T. A. HALGREN, *Merck Molecular Force Field. II. MMFF94 van der Waals and electrostatic parameters for intermolecular interactions*, J. Comp. Chem. **17** (1996) 520–552.
- [126] T. A. HALGREN, *Merck Molecular Force Field. III. Molecular geometrics and vibrational frequencies for MMFF94*, J. Comp. Chem. **17** (1996) 553–586.
- [127] T. A. HALGREN, *Merck Molecular Force Field. IV. Conformational Energies and Geometries for MMFF94*, J. Comp. Chem. **17** (1996) 587–615.
- [128] T. A. HALGREN, *Merck Molecular Force Field. V. Extension of MMFF94 using experimental data, additional computational data, and empirical rules*, J. Comp. Chem. **17** (1996) 616–641.
- [129] V. I. ARNOLD, *Mathematical Methods of Classical Mechanics*, Graduate Texts in Mathematics, Springer, New York, 2nd edition, 1989.
- [130] K. A. DILL and H. S. CHAN, *From Levinthal to pathways to funnels: The “new view” of protein folding kinetics*, Nat. Struct. Biol. **4** (1997) 10–19.
- [131] A. FERSHT, *Structure and Mechanism in Protein Science: A Guide to Enzyme Catalysis and Protein Folding*, Freeman, W. H., New York, 1998.
- [132] S. REICH, *Smoothed Langevin dynamics of highly oscillatory systems*, Physica D **118** (2000) 210–224.
- [133] C. M. DOBSON, A. ŠALI, and M. KARPLUS, *Protein folding: A perspective from theory and experiment*, Angew. Chem. Int. Ed. **37** (1998) 868–893.
- [134] D. PERCHAK, J. SKOLNICK, and R. YARIS, *Dynamics of rigid and flexible constraints for polymers. Effect of the Firman potential*, Macromolecules **18** (1985) 519–525.
- [135] M. R. PEAR and J. H. WEINER, *Brownian dynamics study of a polymer chain of linked rigid bodies*, J. Chem. Phys. **71** (1979) 212.



- [136] E. HELFAND, *Flexible vs. rigid constraints in Statistical Mechanics*, J. Chem. Phys. **71** (1979) 5000.
- [137] T. LAZARIDIS and M. KARPLUS, *Discrimination of the native from misfolded protein models with an energy function including implicit solvation*, J. Mol. Biol. **288** (1999) 477–487.
- [138] T. LAZARIDIS and M. KARPLUS, “New view” of protein folding reconciled with the old through multiple unfolding simulations, Science **278** (1997) 1928–1931.
- [139] K. A. DILL, *Polymer principles and protein folding*, Prot. Sci. **8** (1999) 1166–1180.
- [140] C. LEVINthal, *How to fold gracefully*, in *Mossbauer Spectroscopy in Biological Systems*, edited by J. T. P. DEBRUNNER and E. MUNCK, pp. 22–24, Allerton House, Monticello, Illinois, 1969, University of Illinois Press.
- [141] J. D. BRYNGELSON, J. N. ONUCHIC, N. D. SOCCI, and P. G. WOLYNES, *Funnels, pathways, and the energy landscape of protein folding: A synthesis*, Proteins **21** (1995) 167–195.
- [142] H. S. CHAN and K. A. DILL, *Protein folding in the landscape perspective: Chevron plots and non-Arrhenius kinetics*, PROTEINS: Struct. Funct. Gen. **30** (1998) 2–33.
- [143] R. L. BALDWIN, *The nature of protein folding pathways: The classical versus the new view*, J. Biomol. NMR **5** (1995) 103–109.
- [144] J. N. ONUCHIC and P. G. WOLYNES, *Theory of protein folding*, Curr. Opin. Struct. Biol. **14** (2004) 70–75.
- [145] S. S. PLOTKIN and J. ONUCHIC, *Understanding protein folding with energy landscape theory. Part I: Basic concepts*, Quart. Rev. Biophys. **35** (2002) 111–167.
- [146] J. D. BRYNGELSON and P. G. WOLYNES, *Spin-glasses and the statistical-mechanics of protein folding*, Proc. Natl. Acad. Sci. USA **84** (1987) 7524–7528.
- [147] R. DAY and V. DAGGETT, *Ensemble versus single-molecule protein unfolding*, Proc. Natl. Acad. Sci. USA **102** (2005) 13445–1450.
- [148] S. E. RADFORD and C. M. DOBSON, *Insights into protein folding using physical techniques: Studies of lysozyme and alpha-lactalbumin*, Phil. Trans. R. Soc. London B Biol. Sci. **348** (1995) 17–25.
- [149] C. D. SNOW, H. NGUYEN, V. S. PANDE, and M. GRUEBELE, *Absolute comparison of simulated and experimental protein-folding dynamics*, Nature **420** (2002) 102–106.
- [150] A. A. DENIZ, T. A. LAURENCE, G. S. BELIGERE, M. DAHAN, A. B. MARTIN, D. S. CHEMLA, P. E. DAWSON, P. G. SCHULTZ, and S. WEISS, *Single-molecule protein folding: Diffusion fluorescence resonance energy transfer studies of the denaturation of chymotrypsin inhibitor 2*, Proc. Natl. Acad. Sci. USA **97** (2000) 5179–5184.

- [151] D. BAKER, *Metastable states and folding free energy barriers*, Nat. Struct. Biol. **5** (1998) 1021–1034.
- [152] J. L. SOHL, S. S. JASWAL, and D. A. AGARD, *Unfolded conformations of  $\alpha$ -lytic protease are more stable than its native state*, Nature **395** (1998) 817–819.
- [153] S. KIRKPATRICK, C. D. GELATT, and M. P. VECCHI, *Optimization by simulated annealing*, Science **220** (1983) 671–680.
- [154] V. CERNY, *A thermodynamical approach to the travelling salesman problem: An efficient simulation algorithm*, J. Optimiz. Theory App. **45** (1985) 41–51.
- [155] R. DAWKINS, *The Blind Watchmaker: Why the Evidence of Evolution Reveals a Universe without Design*, W. W. Norton & Company, New York, 1987.
- [156] N. GŌ and H. TAKETOMI, *Respective roles of short- and long-range interactions in protein folding*, Proc. Natl. Acad. Sci. USA **75** (1978) 559–563.
- [157] G. I. MAKHATADZE and P. L. PRIVALOV, *Energetics of protein structure*, Adv. Prot. Chem. **47** (1995) 307–425.
- [158] P. DERREUMAUX, *Ab initio polypeptide structure prediction*, Theo. Chem. Acc. **104** (2000) 1–6.
- [159] R. A. ABAGYAN, *Protein structure prediction by global energy optimization*, in *Computer Simulations of Biomolecular Systems*, edited by W. F. VAN GUNSTEREN, volume 3, Kluwer academic publishing, Dordrecht, 1997.
- [160] C. D. SNOW, E. J. SORIN, Y. M. RHEE, and V. S. PANDE, *How well can simulation predict protein folding kinetics and thermodynamics?*, Annu. Rev. Biophys. Biomol. Struct. **34** (2005) 43–69.
- [161] A. R. MACKERELL JR., M. FEIG, and C. L. BROOKS III, *Extending the treatment of backbone energetics in protein force fields: Limitations of gas-phase quantum mechanics in reproducing protein conformational distributions in molecular dynamics simulations*, J. Comp. Chem. **25** (2004) 1400–1415.
- [162] A. V. MOROZOV, T. KORTemme, K. TSEMEKHMAN, and D. BAKER, *Close agreement between the orientation dependence of hydrogen bonds observed in protein structures and quantum mechanical calculations*, Proc. Natl. Acad. Sci. USA **101** (2004) 6946–6951.
- [163] M. KARPLUS and J. A. MCCAMMON, *Molecular dynamics simulations of biomolecules*, Nat. Struct. Biol. **9** (2002) 646–652.
- [164] N. G. VAN KAMPEN, *Stochastic processes in Physics and Chemistry*, North-Holland, Amsterdam, 1981.
- [165] H. GOLDSTEIN, C. POOLE, and J. SAFKO, *Classical Mechanics*, Addison-Wesley, 3rd edition, 2002.