Karolina Majorek, Łukasz Kozłowski, Marcin Jąkalski and Janusz M. Bujnicki

21 2.1 Introduction 22

11 12

23 The famous hypothesis formulated by the Nobel Prize laureate Christian Anfinsen states 24 that 25

26 The three-dimensional structure of a native protein $[\ldots]$ is determined by the totality of 27 interatomic interactions and hence by the amino acid sequence, in a given environment. In terms of natural selection through the 'design' of macromolecules during evolution, this idea 28 emphasized the fact that a protein molecule only makes stable, structural sense when it exists 29 under conditions similar to those for which it was selected – the so-called physiological 30 state.1,2 31

32 On the one hand, essentially all globular domains in proteins studied so far appear to 33 conform to this rule. Most proteins (or at least their major fragments) have been found to fold into unique, well-defined, stable three-dimensional structures under very broadly 34 defined 'physiological conditions' that also include 'laboratory conditions' under which 35 36 protein samples are prepared for biophysical and biochemical characterization. In agree-37 ment with Anfinsen's hypothesis, variations in conditions (e.g. change of pH or addition 38 of a ligand) or changes in a sequence (e.g. due to proteolytic cleavage and removal of a 39 sequence fragment) may result in structural changes that are often functionally relevant, e.g. if a protein's function requires opening and closing of a cavity that is used for bind-40 41 ing of another molecule. On the other hand, a growing number of protein sequences (or 42 sequence fragments) have been found to be mostly unstructured (review:³). These 'intrinsically disordered proteins' (IDPs) may assume a defined structure only under very specific 43 44

- 45 Prediction of Protein Structures, Functions, and Interactions Edited by Janusz Bujnicki
- 46 © 2009 John Wiley & Sons, Ltd

conditions, e.g. in the presence of another molecule (e.g. upon binding to another protein
or a ligand). In the absence of a stabilizing factor, these sequences exist as an ensemble
of rapidly interconverting different conformations. Thus, for many IDPs, the Anfinsen's
'stable physiological state for which the protein was selected' is significantly different from
'standard' conditions assumed for other proteins.

Anfinsen's hypothesis implies that the knowledge of amino acid sequence and of a given 6 7 environment should be sufficient to infer the native structure of a protein (or to predict the 8 lack of a stable structure). However, despite seemingly solid theoretical foundations, accurate prediction of protein structure from the primary chemical structure of the polypeptide 9 10 chain remains one of the greatest challenges in biology. Thus far, there have been two major 11 approaches to predict the unknown protein structure from known amino acid sequence: one 12 relies on our knowledge of 'first principles', i.e. the laws of physics, while the other is based on rules inferred from comparative analysis of experimentally solved protein structures. 13 Both approaches had some successes, but neither of them has achieved the final goal. 14

15 The 'physical' approach has been successfully applied already in 1951 by Linus Pauling 16 and Robert Corey, who predicted the existence of two periodic structural motifs formally defined by the pattern of hydrogen bonds that may be formed by the protein backbone: the 17 spiral α -helix with 3.6 amino acid residues per turn⁴ and the flat β -sheet comprising two or 18 more β -strands having an extended zigzag conformation.⁵ These two secondary structure 19 elements (SSEs) now are known to be major features of protein architecture, with >50% of 20 21 residues of an average protein assuming either helical or extended conformation. Helices 22 and strands provide a natural frame for insightful protein structure visualization (with a 23 helix often represented as a tube or a spiral and strand as an arrow), and are widely used to 24 describe protein three-dimensional folds. They are also used by many programs that use 25 simplified protein structure representation (e.g. SSEs instead of individual amino acids) to 26 speed up calculations, e.g. for superposition of protein structures.

27 Secondary structure is much more conserved in the evolution than amino acid sequence; therefore accurate prediction of SSEs from sequence would be of great benefit in structural 28 29 bioinformatics. For instance the knowledge of SSEs can help to guide sequence alignment or improve existing sequence alignment of remotely related sequences with low sequence 30 31 similarity (see Chapter 1 by Kaminska *et al.* in this volume). Secondary structure predic-32 tion is also a good starting point toward elucidating the three-dimensional structure – it 33 serves as an intermediate step in the protein fold-recognition procedure, i.e. identification of templates for comparative modeling (see Chapter 4 by Kosinski et al.) and may provide 34 useful restraints both in comparative modeling and in *de novo* modeling (see Chapter 5 by 35 Gront et al.). However, it has been found that it is quite difficult to predict accurately, which 36 37 type of secondary structure is assumed by each amino acid residue of a protein.⁶ Compu-38 tational simulations of peptide and protein folding based on the 'physics-based' approach 39 have been carried out in attempt to predict both local structure and global conformation (review:⁷). The first applications of force field methods to study peptide conformations 40 date back to calculations performed by Nemethy and Scheraga.⁸ However, due to an ex-41 42 tremely large number of degrees of freedom and very complex calculations of energies, 43 such simulations require extremely large computer resources, such as supercomputers or massively parallel distributed computing.⁹ Alas, despite recent advancement in computer 44 hardware and software, physics-based simulation techniques remain incapable of confi-45 dently predicting structures of even moderately sized proteins (>100 amino acid residues). 46

Definition of Secondary Structure and Its Assignment for Known Protein Structures 41

On the other hand, as soon as the first protein structures solved by X-ray crystallography 1 2 have been determined, it has been observed that secondary structures tend to exhibit regular 3 arrangements of amino acid residues of certain type. The regularities are due to the local 4 periodicity of helical and extended conformations (3.6 and 2 residues per repeated segment, 5 respectively), and the global tendency of a protein to form a well-packed hydrophobic core and to place hydrophilic residues at the surface (at least in water-soluble proteins). For 6 7 example, SSEs buried in the protein core are composed mainly of hydrophobic residues, 8 while SSEs at the protein surface tend to be amphipathic and show an alternating pattern of hydrophobic and hydrophilic residues (usually 010101 for strands and 0011011 or 9 10 0010011 for helices), allowing the respective side chains to be buried in the protein core or exposed to the solvent. Thus, as soon as a sufficient number of protein structures have been 11 determined to make useful statistics, 'knowledge-based' methods have been developed, 12 aiming at predicting SSEs based on calculated tendencies of different residues or peptide 13 fragments to assume particular conformations. The first attempt of predicting secondary 14 15 structure of polypeptides using the 'knowledge-based' approach dates back to the same time 16 as the afore-mentioned physics-based analyses. It was performed in 1965 by Guzzo,¹⁰ who inferred simple rules for preferences of different amino acid residues to form helical and 17 non-helical regions. With the growing number of available protein structures and sequences 18 the statistics have quickly improved, and new algorithms have been developed, yielding 19 methods that are far from perfect, but achieve useful accuracy of about 80%. Currently 20 21 the 'knowledge-based' approach is used by essentially all methods for secondary structure prediction, as well as methods for order/disorder prediction and for inference of other 22 23 simple structural features from the primary sequence. This chapter aims at providing a comprehensive overview of these methods. The 'physical' approach to protein structure 24 25 prediction is beyond the scope of this chapter and will not be reviewed here – instead, it 26 will be referred to in other chapters, in particular Chapter 5 by Gront et al. Because of 27 significant differences between proteins that function as water-soluble and those that are embedded in biological membranes, secondary structure prediction methods for each of 28 two types are different and will be discussed separately. We will also discuss prediction 29 of higher order motifs formed by certain types of SSEs, the so called 'supersecondary' 30 31 structures (e.g. coiled coils or β hairpins), and prediction of contacts between residues that are remote in primary structure. 32

33

34 35

36 2.2 Definition of Secondary Structure and Its Assignment 37 for Known Protein Structures

38

One caveat of the knowledge-based structure prediction is the requirement of unambigu-39 40 ous definition of secondary structure elements or ordered vs. disordered regions. Statistical 41 methods and machine learning methods require the input data to be appropriately clas-42 sified in order to make meaningful predictions. Unfortunately, defining the boundaries between disordered and ordered regions or between helix, sheet, and coil structures is 43 44 arbitrary, and commonly accepted standard assignments do not exist. Therefore, various 45 researchers employed different criteria that in some cases have lead to considerably different assignments. 46

Secondary structure is formally defined by the hydrogen bonds, but the hydrogen bonding 1 2 is correlated with other features, such as dihedral angels generally adopted by particular types of secondary structure, which has given rise to less formal definitions of SSEs. 3 To standardize secondary structure assignment, the Dictionary of Secondary Structure in 4 Proteins (DSSP) was designed.¹¹ It was the first method for protein secondary structure 5 assignment available as a computer program, and has remained most popular until now. 6 7 DSSP classifies each amino acid residue in a protein with known 3D structure into one 8 of 8 types of secondary structure, based on recognition of hydrogen bonding patterns and geometrical features defined in terms of the concepts torsion and curvature of differential 9 10 geometry. The 7 types of SSEs include the α -helix (H), the β -strand (E, for 'extended'), 11 and less frequent types of secondary structure namely non- α helices: 3/10 (G) and π (I), 12 isolated β -bridge (B), highly curved bend (S) and hydrogen-bonded turn (T), while the remaining residues are classified as outside of SSEs.¹¹ 13

Another quite popular method for secondary structure assignment named STRIDE was 14 proposed later by Frishman and Argos.¹² STRIDE is based on the combined use of hy-15 16 drogen bond energy and statistically derived backbone torsional angle information, with parameters of the pattern recognition procedure optimized to improve (compared to DSSP) 17 agreement with manual designations provided by the crystallographers as a standard-18 of-truth. More recently developed methods for secondary structure assignment include 19 P-SEA,¹³ SECSTR,¹⁴ KAKSI,¹⁵ SEGNO,¹⁶ and PALSSE.¹⁷ These methods are based on 20 21 either the hydrogen-bond pattern, geometric features, expert knowledge or their combi-22 nations. However, they often disagree on their assignments, up to 25%. The discrepancy 23 among different methods is caused by nonideal configurations of helices and sheets in 24 experimentally solved structures and by different definitions of helices and strands. Of 25 particular interest is PALSSE, which identifies only two types of SSEs that can be ap-26 proximated by vectors: helix and strand. In contrast to other algorithms, which identify a 27 secondary structure state for every residue in a protein chain, PALSSE attributes residues to SSEs in such a way that consecutive elements may overlap, thus allowing residues located 28 29 at the overlapping region to have more than one secondary structure type. This method is robust to coordinate errors and can be used to define SSEs even in poorly refined and 30 31 low-resolution structures (e.g. if only C- α atoms are available, thus if no hydrogen-bonds are present). PALSSE usually assigns a larger fraction of residues to SSEs as compared to 32 33 other methods, e.g. 80% vs. 53% in the case of DSSP.¹⁷

Discrepancies with structural assignment concern not only algorithms. Protein structures 34 are dynamic objects with some regions more mobile than others. Local conformations near 35 36 the ends of secondary structures vary under native conditions, but may be forced to assume 37 a single conformation in crystals due to packing constraints, hence secondary structure 38 assignments differ by about 5–15 percentage points between different X-ray versions or 39 different NMR models for the same protein⁶ (Figure 2.1). This inherent protein flexibility is the main reason why the theoretical upper limit of secondary structure prediction 40 41 accuracy is about 90%, for a particular SSE assignment method. Recently it has been 42 proposed that instead of relying on single structures, structure assignment methods should 43 be assessed based on the similarity of the secondary structures assigned to established pair-44 wise sequence-alignment benchmarks, where these benchmarks are determined by prior structural alignments of the protein pairs. The use of this criterion has led to identification 45 of STRIDE and KAKSI as the most robust methods (PALSSE was not included in that 46

1 2 3 4 7 8 9



10

11 Figure 2.1 Illustration of secondary structure, conformational variability, and order vs disor-12 der. Three-dimensional structure of small protein SUMO-1 solved by NMR (1a5r in Protein 13 Data Bank), 10 alternative models shown in the 'cartoon' representation. Black spirals indicate 14 α -helices, dark grey arrows indicate β -strands, white coils indicate loops. While the central part 15 shows relatively ordered structure (with only some fluctuations at one end of the helix), the 16 N- and C-terminal regions (left and right, respectively) show 'intrinsic disorder'. Interestingly, a short helical region persists in the disordered N-terminal tail, demonstrating the presence of 17 secondary structure despite the absence of stable tertiary structure 18

19 20

21 comparison), and to development of a consensus of STRIDE, KAKSI, SECSTR, and P-SEA, called SKSP, which is 2–3% higher in agreement with structurally aligned residues 22 than DSSP for three established alignment benchmarks.¹⁸ 23

Summarizing, assignments of secondary structure for one particular protein may vary, 24 25 depending on the method used. Thus, in theoretical protein structure prediction it is important to select, which type of assignment is going to be predicted. Thus far, DSSP has 26 27 been used as the 'golden standard' because of its popularity among crystallographers, but it is likely that methods for secondary structure assignment that are more consistent 28 with the 3D structure alignments (e.g. SKSP) may lead to improved secondary struc-29 ture prediction. Another caveat of secondary structure prediction methods is that they are 30 31 aimed at predicting only the three basic classes of local structure: (H)elix, (E)xtended, and (C)oil; thus, e.g. an 8-letter 'structural alphabet' used in the DSSP notation is reduced to a 32 33 3-letter alphabet. In different algorithms it is done according to different conversion rules, which may yield an apparent increase of accuracy, but cause errors when the predicted 34 secondary structure is used to predict 3D structure.⁶ Karplus and coworkers found that the 35 replacement of simple alphabets of secondary structure with highly informative, detailed 36 37 alphabets can improve detection and alignment of structurally similar, but remotely related proteins.¹⁹ Examples of such alphabets include STR, an enhanced version of DSSP that 38 subdivides DSSP letter E (strand) into six letters, according to properties of a residue's rela-39 tionship to its strand partners (number of partners and their parallel/antiparallel character) 40 or Protein Blocks, a set of overlapping protein backbone fragments of length 5 amino acid 41 residues.²⁰ HMMSTR²¹ implements yet another solution to this problem: in this method 42 protein structure is represented by short structural fragments taken from the database of 43 44 known structures; for secondary structure it uses an alphabet of 11 conformation states, 10 corresponding to Φ - Ψ angle regions and one for cis-peptide bonds. A recently developed 45 method Real-SPINE^{22,23} predicts real values of torsion angles from a given sequence. 46

P1: OTA

3

44 First Steps of Protein Structure Prediction

2.3 Prediction of Secondary Structure and Solvent Accessibility for Water-soluble Proteins

4 The first empirical prediction system aiming at predicting SSEs from protein sequence, 5 based on statistics calculated from structures solved by X-ray crystallography, was developed by Fasman and Chou.^{24,25} This very simple method is based on analysis of the relative 6 7 frequency of each amino acid in helices, strands and coils and on assumption that single 8 residue individually influences secondary structure. Subsequently, a more sophisticated 9 GOR method has been developed, which is based on information theory and Bayesian 10 statistics, and takes into account not only one residue but also adjacent positions in the 11 sequence. These methods have been built-in into commercial software packages for protein 12 sequence analysis and structure modeling and have become very popular among biologists 13 despite their accuracy was only slightly better than random. The main limitation of these 14 early methods was a small amount of known 3D structures, from which parameters could 15 be derived. Besides, these methods did not utilize any evolutionary information and were 16 applicable to single sequences rather than for multiple sequence alignments (MSAs) of homologous sequences.²⁶ Although over the years, both the Chou-Fasman²⁷ and the GOR 17 18 methods²⁸ have been improved, the level of their accuracy is inferior to the best modern 19 methods.

20 A significant improvement in prediction accuracy (>70%) has been achieved by 'sec-21 ond generation' methods such as PHD,²⁹ SAM-T98,³⁰ and PSIPRED,³¹ which utilized 22 MSA-derived information concerning sequence conservation, often combined with ma-23 chine learning techniques such as artificial neural networks (ANNs), nearest-neighbor 24 search (NNS) methods, and support vector machines (SVMs), or advanced statistical 25 methods such as Hidden Markov Models (HMM) (review:³²). These methods were also 26 made available as web servers (Table 2.1). MSA, provided by the user or generated by 27 an internal routine of an algorithm, is usually based on identification of homologs by 28 searches of protein sequence databases (see Chapter 1 by Kaminska et al. in this volume). 29 It is important to note that PSIPRED was the first method, in which iterative PSI-BLAST 30 sequence searches have been introduced, compared to single-pass searches in earlier meth-31 ods. Currently, iterative database searches to obtain the input MSA for prediction methods 32 are considered a standard. Typically, patterns in sequence variability observed in MSA 33 provide information on conservation of core elements (hydrophobic core and regions 34 important for protein function), while the location of insertions and deletions (indels) 35 hints at a position of surface-exposed loops. Incorporated machine learning techniques 36 allow training the methods on known structures to learn characteristic sequence-structure 37 patterns and then use those patterns to predict the secondary structure of the query pro-38 tein. Most of the SSE prediction methods of the above-mentioned generation, or their 39 derivatives developed later, have been associated with predictors of solvent accessibil-40 ity used to identify residues that are buried to different extents in the hydrophobic core 41 (Table 2.2).

⁴² In addition to methods for predicting the three main types of SEEs, there are several ⁴³ methods based on sequence profile analysis for predicting certain types of local structure, ⁴⁴ including various hairpin structures,⁵² or specialized in α -turns,^{71,72} β -turns,^{51,73} γ -turns,⁵⁰ ⁴⁵ and π -turns.⁷⁴ There are also methods to predict conformation of individual residues, e.g. ⁴⁶ the trans/cis state of Pro.^{75,76} To our knowledge, these types of methods have not yet been Prediction of Secondary Structure and Solvent Accessibility for Water-soluble Proteins 45

 Table 2.1
 Software for secondary structure prediction

Program	URL (http://)
	Three-state (Helix/Extended/Coil) prediction
IPSSP ³³ (for single seque	ences) exon.gatech.edu/genemark/ipssp/webIPSSP.cgi
PSIPRED ³¹	bioinf.cs.ucl.ac.uk/psipred/
SSPRO ³⁴	scratch.proteomics.ics.uci.edu/
PHD ²⁹	www.predictprotein.org/
PROFsec ³⁵	www.predictprotein.org/
PRED2ARY ³⁶	alexander.ucsf.edu/~imc/pred2arv/
APSSP2 ³⁷	www.imtech.res.in/raghava/apssp2/
PREDATOR ³⁸	ftp://ftp.ebi.ac.uk/pub/software/unix/predator/
HMMSTR ²¹	www.bioinfo.rpi.edu/~bystrc/hmmstr/
NPREDICT ³⁹	www.cmpharm.ucsf.edu/~nomi/nnpredict.html
PORTER ⁴⁰	distill.ucd.ie/porter/
HYPROSPII ⁴¹	bioinformatics.iis.sinica.edu.tw/HYPROSPII/
SAM-T06 ⁴²	www.soe.ucsc.edu/compbio/SAM_T06/T06-guerv.html
INET ⁴³	www.compbio.dundee.ac.uk/Software/INet/inet.html
SABLE ⁴⁴	sable.cchmc.org/
YASSPP ⁴⁵	glaros.dtc.umn.edu/vasspp/
YASPIN ⁴⁶	ibivu.cs.vu.nl/programs/yaspinwww/
CRNPred ⁴⁷	ftp.bioinformatics.org/pub/crnpred/
JUFO3D ⁴⁸	www.meilerlab.org/index.php
SPINE ²²	sparks.informatics.iupui.edu/SPINE/spine.html
Other types of sec	ondary and supersecondary structure, and other types of local
	conformation
TURNS $(\alpha, \beta, \gamma)^{49,50}$	imtech.res.in/raghava/
β -Turn ⁵¹	serine.umdnj.edu/~zhangq3/betaturn/prediction.htm
TURNPRED ⁵²	www.meilerlab.org/index.php
COILS ⁵³	www.ch.embnet.org/software/COILS_form.html
MARCOIL ⁵⁴	www.isrec.isb-sib.ch/webmarcoil/webmarcoilC1.html
PCOILS ⁵⁵	toolkit.tuebingen.mpg.de/pcoils
PairCoil2 ⁵⁶	groups.csail.mit.edu/cb/paircoil2/paircoil2.html
MultiCoil ⁵⁷	groups.csail.mit.edu/cb/multicoil/cgi-bin/multicoil.cgi
LearnCoil ⁵⁸	groups.csail.mit.edu/cb/learncoil-vmf/cgi-bin/vmf.cgi
۲۰ ^۲ ۸	leta-servers' for secondary structure prediction
JPRED ⁵⁹	www.compbio.dundee.ac.uk/~www-jpred/
NPS@ ⁶⁰	npsa-pbil.ibcp.fr
META-PP ⁶¹	www.predictprotein.org/meta.php
PROTEUS ⁶²	wishart.biology.ualberta.ca/proteus
DISTILL ⁶³	distill.ucd.ie
GeneSilico ⁶⁴	genesilico.pl/meta2/

40

1

41

integrated into metaservers for secondary or tertiary structure prediction and their practical
 utility for protein modeling and function prediction remains to be established.

⁴⁴ Currently the recommended approach to secondary structure prediction involves combin-

⁴⁵ ing the results of different methods; it may involve advanced machine learning approaches,

⁴⁶ such as voting, linear discrimination, neural networks or decision trees⁷⁷ or even simple

1
 Table 2.2
 Software for solvent accessibility prediction
 2 URL (http://) Program 3 4 Jnet43 www.compbio.dundee.ac.uk/Software/JNet/jnet.html 5 PHDacc³⁵ www.predictprotein.org/ 6 PROFacc³⁵ www.predictprotein.org/ SABLE⁶⁵ 7 sable.cchmc.org NETASA⁶⁶ www.netasa.org 8 MLRprdsec⁶⁷ spg.biosci.tsinghua.edu.cn/ 9 WESA⁶⁸ pipe.scs.fsu.edu/wesa.html 10 ACCpro⁶⁹ scratch.proteomics.ics.uci.edu/ 11 SARpred⁷⁰ www.imtech.res.in/raghava/sarpred/ 12 SPINE²² sparks.informatics.iupui.edu/SPINE/spine.html 13 PaleAle⁶³ distill.ucd.ie/paleale/ 14

15 16

consensus.⁷⁸ The idea of combining different prediction methods was first implemented 17 in JPRED,⁵⁹ a consensus meta-server that standardizes input and output requirements of a 18 range of secondary structure prediction algorithms, each representing a different prediction 19 strategy, and computes a consensus of PHD, NNSSP, DSC, and PREDATOR secondary 20 21 structure predictions. In addition, the output of the JPRED server includes predictions of 22 solvent accessibility by the JNET method,⁴³ as well as predictions of coiled-coil regions and transmembrane helices (see below), which however are not directly incorporated in 23 24 the calculation of the secondary structure.

The most recent class of meta-approaches, exemplified by PROTEUS⁶² exploits the 25 26 observation that if an experimentally determined three-dimensional structure of a closely 27 related protein is known, then copying the secondary structure assignment from the known structure provides a better result than by predicting it de novo. PROTEUS initially carries 28 29 out a sequence similarity search against the PDB database in order to determine if the whole or a part of the query sequence is significantly similar to a known structure, and if 30 31 such a template structure is found, secondary structure mapping is carried out from the 32 template to the query based on a sequence alignment. For the sequence segments that 33 are not covered by template structures, *de novo* secondary structure prediction is carried out with three different, high quality neural network approaches (PSIPRED, JNET and 34 TRANSSEC), whose results are combined into a consensus prediction by the fourth neural 35 36 network. Merging template-based predictions and *de novo* predictions allows PROTEUS 37 to yield a full sequence prediction, regardless of the extent of sequence overlap to a PDB 38 hit (when complete 3D-to-2D mapping is achieved, when only partial coverage is provided 39 and when no homologue with known structure can be found), and to achieve high average accuracy of >80% per residue. A similar approach of merging template-based and *de novo* 40 41 predictions of secondary structure and solvent accessibility has been implemented in the DISTILL suite.⁶³ 42

43 While the early methods of secondary structure prediction were about 60–65% accurate, 44 with accuracy for β -strands only slightly better than random,⁶ the best modern methods 45 reach about 80% accuracy per residue,²² with ~10% lower accuracy for β -strands. The 46 difference between theoretical upper limit of prediction accuracy and actual secondary

inst steps of Hotelin surdetate Hedicuon

Prediction of Secondary Structure for Transmembrane Proteins 47

structure prediction accuracy, and between level of prediction accuracy of α -helices and 1 2 β -strands, is mainly due to difficult to detect long-range interactions that may influence 3 secondary structure formation. It has been shown that the same amino acid sequence of substantial length may fold as α -helix when in one position in primary protein sequence 4 but as β -sheet when in another sequence context.⁷⁹ Besides, during the folding process, 5 a certain fragment of a protein might first adopt a secondary structure preferred by the 6 7 local sequence and later, because of non-local interactions, be transformed to another secondary structure. The latter concern has been addressed in the method 3D-JUFO,48 8 which combines iterative de novo secondary structure prediction using an approach similar 9 10 to PSIPRED with tertiary structure prediction with the ROSETTA method (see Chapter 5 by Gront et al.), followed by re-prediction of SSEs based on local environment of 11 particular residues observed in models of tertiary structure. 3D-JUFO achieves remarkable 12 accuracy of 80%, with notable improvement of accuracy for β -strand prediction (76%) 13 over sequence-only methods. Another interesting recently developed method that brings 14 15 the accuracy of secondary structure prediction close to the theoretical limit combines 16 bioinformatics methodology with experimental techniques of circular dichroism (CD) and Fourier transform infrared (FTIR) spectroscopy for assessing the overall secondary 17 structure content.80 18

19 20

21 2.4 Prediction of Secondary Structure for Transmembrane Proteins

22

23 Membrane proteins are different from water-soluble proteins in that a large fraction of their surface is hydrophobic to enable stability in the environment of a lipid bilayer. They 24 25 constitute about 20–30% of all proteins in the fully sequenced genomes, and are typically involved in cell signaling, molecular pumping and energy transduction. Integral membrane 26 27 proteins consist of one or more transmembrane (TM) segments and can be divided into two structural classes: the α -helical TM proteins and the β -barrel TM proteins, varying 28 in structure, localization and physicochemical features. Typical TM proteins of the more 29 abundant α -helical class are present in all types of biological membranes including outer 30 31 membranes. They comprise one or more hydrophobic α -helical membrane spanning regions 32 separated by hydrophilic loops that are exposed into the solvent (review:⁸¹). TM β -barrel 33 proteins are found only in outer membranes of Gram-negative bacteria, cell wall of Grampositive bacteria, and outer membranes of mitochondria and chloroplasts. They consist of 34 different number of antiparallel, membrane spanning β -strands with a simple up-and-down 35 topology.82 36

37 TM proteins aggregate and precipitate in water and require detergents or nonpolar sol-38 vents for extraction, therefore they are much more difficult to analyze experimentally than their soluble counterparts, in relation to all steps from overexpression to high-resolution 39 structure determination. Although TM proteins represent the most important drug targets, 40 their structure determination has lagged behind that for soluble proteins; currently they 41 represent less than 1% of available crystal structures.⁸³ On the one hand, this situation 42 generates a great deal of pressure to develop effective methods for predicting the structure 43 44 of TM proteins. On the other hand, the paucity of structural data hampers the development of knowledge-based approaches. Nonetheless, for both types of TM proteins specialized 45 structure predictors have been designed, but due to the relatively easily detectable patterns 46

of hydrophobic residues forming α -helical TM segments and much smaller amount of known β -barrel TM proteins structures, the majority of them was focused on the α -helical TM proteins until quite recently (review:⁸⁴).

4 Prediction of TM helices should be intuitively easy due to their hydrophobic nature. However, predictions based solely on hydrophobicity profiles have high error rates. Besides, 5 hydrophobic signal peptides may be confused with TM helices. It is also a consecutive chal-6 7 lenge to predict TM proteins topology. Prediction of the way in which TM segments cross the membrane (inside-out or outside-in) is done mainly by considering the different charge 8 distribution between the inside (cytoplasmic) and outside (extracellular) regions connect-9 ing TM segments, and by application of the so-called 'positive-inside rule'85 based on the 10 11 observation that there is an overrepresentation of positively charged residues in the intra-12 cellular loops of TM proteins. Contemporary approaches usually predict both localization of TM segments and their orientation (topology). The best methods such as PHOBIUS⁸⁶ 13 or MEMSAT3⁸⁷ utilize evolutionary information as well as discriminate against signal 14 peptides. Prediction of TM segments for β -barrel proteins is more difficult, because the 15 16 strands are amphipathic. They contain 10–22 residues with alternating hydrophobic side chains facing the lipid bilayer and hydrophilic side chains facing the internal pore. To 17 predict the β -barrel type of TM proteins a small number of specialized algorithms have 18 been developed based on standard statistical and machine learning techniques including 19 HMMs, ANNs, or SVMs (Table 2.3). 20

21 As in the case of secondary structure prediction for globular soluble proteins, consensus 22 methods perform much better compared to each individual prediction method separately 23 and the recommended strategy for identification membrane spanning segments and their 24 orientation in membranes is to use many different methods and combine results into a consensus prediction. Examples of 'metaservers' that combine the results of several indi-25 vidual methods, providing a more accurate consensus prediction, include BPROMPT,¹⁰⁰ 26 ConPredII,¹¹⁰ and PONGO¹¹¹ for α -helical TM proteins, and ConBBPRED¹¹⁴ for β -barrel 27 TM proteins. The newest trends in TM structure prediction include meta-predictions that 28 utilize predictions of solvent accessibility and secondary structure propensity typical for 29 globular proteins in the form of 'structural profiles'.¹⁰² There have also been attempts 30 to make concurrent prediction of secondary and tertiary structure by simulating folding 31 in lipid membranes, e.g. with modified versions of *de novo* structure prediction methods 32 FRAGFOLD¹¹⁵ and ROSETTA.^{116,117} 33

In addition to predictors specific for TM proteins, a new method MeTaDoR has been recently proposed that predicts membrane-binding peripheral proteins that do not form an integral part of the membrane, but bind to it mostly in a reversible manner and thereby function in various important processes, including cell signaling and membrane trafficking.¹¹³

39

40 **2.5 Prediction of Supersecondary Structure**

41

Individual SSEs may be arranged in simple geometrical shapes forming recurring supersecondary structures. There is a number of well-defined $\alpha - \alpha$, $\beta - \beta$, $\alpha - \beta$ and $\beta - \alpha$ structural motifs that serve as 'building blocks' of tertiary structure. Prediction of supersecondary structures can be an important step toward building a tertiary structure from the specified secondary structure elements.¹¹⁸ The β -hairpin, comprising two adjacent antiparallel

Prediction of Supersecondary Structure 49

Program	URL (http://)			
	α -TM proteins			
HMMTOP ⁸⁸	www.enzim.hu/hmmtop/			
DAS ⁸⁹	www.sbc.su.se/~miklos/DAS/			
PHDhtm ⁹⁰	www.predictprotein.org/			
TMAP ⁹¹	bioinfo4.limbo.ifm.liu.se/tmap/index.html			
TMHMM ⁹²	www.cbs.dtu.dk/services/TMHMM/			
Tmpred ⁹³	www.ch.embnet.org/software/TMPRED_form.html			
MEMSAT3 ⁸⁷	bioinf.cs.ucl.ac.uk/memsat			
TopPred2 ⁹⁴	bioweb.pasteur.fr/seqanal/interfaces/toppred.html			
WHAT ⁹⁵	saier-144-37.ucsd.edu/what.html			
THUMBUP ⁹⁶	sparks.informatics.iupui.edu/Softwares-Services_files/thumbup.htr			
UMDHMM ⁹⁶	sparks.informatics.iupui.edu/Softwares-Services_files/umdhmm.ht			
PRED-TMR ⁹⁷	athina.biol.uoa.gr/PRED-TMR/			
HMM-TM ⁹⁸	biophysics.biol.uoa.gr/HMM-TM/			
ORIENTM ⁹⁹	athina.biol.uoa.gr/orienTM/			
BROMPTIO	www.jenner.ac.uk/BPROMPT			
LOCALIZOME	localizome.org			
PHOBIUS ⁶⁶	phobius.sbc.su.se/			
MINNOU ¹⁰²	minnou.cchmc.org			
DDF: 103	β -Iransmembrane proteins			
BBF*105	www-biology.ucsd.edu/~msaier/transport/software/bbfsource.tar.			
HMM-B21MR**104	gpcr.biocomp.unibo.it/biodec/			
	minnou.cchmc.org			
B2TMPRED ¹⁰³	gpcr.biocomp.unibo.it/cgi/predictors/outer/pred_outercgi.cgi			
PRED-IMBB ¹⁰⁰	bioinformatics.biol.uoa.gr/PRED-IMBB/			
PROFTED 107	cubic.bioc.columbia.edu/services/proftmb/			
IMBEIA-NEI 100	psts.cbrc.jp/tmbeta-net/			
BOWPres	www.bioinfo.no/tools/bomp			
	Weldservers			
ConProdII (α) ¹¹⁰	www.jeimen.ac.uk/brikOmri biginfo si birosaki u ac in/a/ConProd2/			
$PONCO (\alpha)^{111}$	pongo biocomp unibo it/pongo/			
TI IDS $(\alpha)^{112}$	pongo.orocomp.umoo.n/pongo/ sparks informatics jupui adu/Softwaros Sonvices files/tups htm			
$Con BRDPED (R)^{106}$	bioinformatics high upg gr/ConBBPPED/			
CONDUCTION (p)	Membrane-binding peripheral proteins			
112	memorale binding perpileral proteins			

Table 2.3 Software for prediction of TM regions in proteins

*The BBF program is freely available to academic users upon request to the corresponding author.

38 **HMM-B2TMR is a commercial program, demo version is available.

39 40

1

41 hydrogen bonded β -strands, is an example of the frequently occurring motif for which 42 predictors have been developed. BhairPred¹¹⁹ is an example of a method for discriminating 43 hairpins from non-hairpins; obviously it achieves high accuracy only if the prediction of 44 secondary structure is correct. Coiled coils are another type of super-secondary structure

45 characterized by a bundle of two or more α -helices wrapping around each other. Coiled

46 coil structures have been implicated in inter- and intraprotein interactions, and may be

formed by helices formed by segments of sequence distant in the primary structure or 1 even contributed by different proteins. Thus, coiled coils allow monomeric building blocks 2 to form complex assemblages that can serve as molecular motors and springs (review: 3 ref.¹²⁰). The helices forming coiled coils have a unique pattern of hydrophobicity, which 4 repeats every seven residues (five hydrophobic and two hydrophilic). This sequence pe-5 riodicity has prompted the development of special algorithms to predict the location of 6 α -helices that form coiled coils. According to the recent benchmark, the two best compu-7 tational methods are a HMM-based MARCOIL⁵⁴ and PCOILS,⁵⁵ followed by PairCoil2⁵⁶ 8 (Table 2.2). 9

10

11 12 **2.6 Disorder Prediction**

13

During the past decade, the literature has exploded with reports on intrinsically unstruc-14 15 tured proteins (IDPs). Currently it is estimated that 30–60% of proteins are predicted to 16 contain long stretches of disordered residues. Many of the disordered regions have been confirmed experimentally; they have been often found to be essential for protein function. 17 Interestingly, intrinsic disorder appears to be significantly correlated with certain terms 18 from functional ontologies and with specific functional motifs.¹²¹⁻¹²⁴ In particular, linear 19 motifs¹²⁵ that harbor sites of posttranslational modification, such as phosphorylation, or 20 sites of protein-protein interactions, often fall into regions that are locally disordered or 21 undergo order-disorder transition in different, biologically relevant situations.^{126,127} (see 22 Chapter 1 by Kaminska et al. in this volume). With respect to molecular/biochemical 23 24 function, IDPs have been frequently implicated in protein-nucleic acid interactions as tran-25 scription factors or in protein–protein interactions as e.g. regulators of enzyme activity. 26 With respect to cellular roles, they have been implicated in regulatory processes, in particu-27 lar in regulation of gene expression on the level of transcription and RNA processing, and in cellular signaling. On a more general level, IDPs are crucial for cell survival, proliferation, 28 29 differentiation and apoptosis. Dysfunctions of IDPs may therefore lead to cancer, which makes them particularly important from a biomedical point of view. On the other hand, dis-30 31 ordered regions often prevent crystallization of proteins, or the generation of interpretable 32 NMR data, and in protein bioinformatics – they introduce compositional biases that ham-33 per comparison of sequences of ordered regions. Recognition of disordered regions in a protein is therefore important for delineating boundaries of stably folded protein domains 34 for structural and functional studies and for reducing bias in sequence similarity analyses 35 by avoiding alignment of disordered regions against ordered ones (reviews:^{128,129}). De-36 37 tection of disordered regions may also facilitate identification of domains (see Chapter 1 by Kaminska et al.). A very important resource for disorder is the DISPROT database 38 (http://www.disprot.org).¹³⁰ It links structure and function information for proteins that 39 contains at least one experimentally determined disordered region. 40

The relatively frequent occurrence of IDPs and their importance in understanding protein structure-function relationships and cellular processes make it worthwhile to develop predictors of protein disordered regions. Since the SEG algorithm for identification of lowcomplexity regions that are typically associated with molecular disorder was developed in 1994,¹³¹ an increasing number of groups have been developing such methods. However, as with secondary structure, it is not immediately clear how to unambiguously define

Prediction of Long-range Contacts between Amino Acid Residues 51

'disorder'. The lack of stable structure and conformational heterogeneity can manifest 1 2 itself either at the secondary or tertiary level, and may include sites with varying extent of residual secondary structure and conformational mobility: molten globules, pre-molten 3 globules, liquid-like collapsed-disordered state, or gas-like extended-disordered state.¹³² 4 Various researchers employed different criteria for defining disorder, resulting in numer-5 ous predictors that attempt to identify different features. Thus, depending on a research 6 question being asked, using a single disorder predictor may be insufficient to achieve a 7 meaningful prediction. Ferron *et al.*¹²⁸ presented an informative review of a number of 8 methods published until 2006, highlighting their advantages and drawbacks. Table 2.4 9 10 presents succinct descriptions of disorder predictors, taking into account also the most recently published methods. 11

According to our own benchmark focused at accuracy of predicting regions of short dis-12 order (using the criterion employed in CASP-7, i.e. the absence of resolved coordinates in 13 crystal structures¹⁵³), the best methods include POODLE,¹⁴⁷ DisPSSMP,¹³⁸ and iPDA.¹⁴² 14 We have developed a meta-predictor that reports the results of two primary coiled-coil pre-15 16 dictors (COILS and Marcoil; see above and Table 2.2), and 10 primary disorder predictors (DISOPRED2, GlobPlot, Spritz, DISPROT (VSL2), IUPred, POODLE-L, POODLE-S, 17 iPDA, PrDOS, and DisPSSMP, see Table 2.4). It also calculates a consensus prediction. 18 The disorder meta-predictor is available via the gateway of the GeneSilico metaserver⁶⁴ at 19 http://genesilico.pl/meta2/. 20

21 22

P1: OTA

23 2.7 Prediction of Long-range Contacts between Amino Acid Residues

24

In addition to predicting local structure, a number of methods have been developed to predict 25 contacts between residues that are remote in primary structure. This type of information is 26 27 of particular interest, because it has been shown that it is possible to directly infer threedimensional protein structures, if a sufficiently large number of contacts are known with 28 sufficient accuracy. It has been estimated that as few as one contact on average per seven 29 residues may be sufficient.¹⁵⁴ Various measures of distance and various thresholds may be 30 used to define a contact between two residues (see e.g. ref.¹⁵⁵), however the most common 31 definition of contact used in prediction methods is as a C β -C β pair (C α in the case of Gly 32 residue) less than or equal to 8 Å apart.¹⁵⁶ According to the recent benchmark within the 33 framework of the CASP7 experiment, the best contact predictor is an ANN associated with 34 the SAM-T06 structure prediction server.¹⁵⁷ Other well-performing programs (according 35 to the CASP7 benchmark or to other tests published by their authors) that are available as 36 37 web-servers have been summarized in Table 2.5. Special kinds of methods for long-range contact prediction are those for identification 38

of Cys residues involved in disulfide bond formation (review: ref.¹⁵⁸). Disulfide bonds are 39 primary covalent cross-links between two Cys residues in proteins that play critical roles 40 in stabilizing the protein structures. They can impose a substantial distance and angular 41 42 constraint on the backbone of protein, thus making a large contribution to the stabilization of protein tertiary structures. A number of proposed algorithms for prediction of disulfide 43 44 bonding states of Cys (involved in disulfide formation or not), as well as prediction of disulfide connectivity patterns (with the prior knowledge of disulfide bonding states) have 45 been implemented as freely available web servers (Table 2.5). Most of these methods 46

Program	URL (http://)	Short description
DisEMBL ^{TM133}	Dis.embl.de/	ANN trained to predict classic loop (DSSP), flexible loops with high B-factors, missing coordinates in X-ray structures, regions of low-complexity and prone to aggregation.
DISOPRED2 ¹³⁴	bioinf.cs.ucl.ac.uk/disopred/ disopred.html	SVM trained to predict residues wit missing coordinates. Standalone version available.
DISpro ¹³⁵	www.ics.uci.edu/~baldig/ dispro.html	Recursive neural networks (RNNs) trained to predict missing coordinates.
DISPROT ^{136,137}	www.ist.temple.edu/disprot/ predictor.php	VL2 (least-squares regression) and VL3 (ANN) predict long disorder, VSL2 predicts both short and lon disorder. Standalone version available.
DisPSSMP ¹³⁸	biominer.bime.ntu.edu.tw/ dispssmp/	Radial Basis Function Network (RBFN) trained to predict missing coordinates.
DRIP-PRED ¹³⁹	sbcweb.pdc.kth.se/cgi-bin/ maccallr/disorder/submit.pl	Self-organizing maps (SOMs) trainer to predict missing coordinates.
FoldIndex© ¹⁴⁰	Bip.weizmann.ac.il/fldbin/findex	Simple method to predict whether a given protein will fold or not, bas on average hydrophobicity and n charge.
FoldUnfold ¹⁴¹	skuld.protres.ru/~mlobanov/ ogu/ogu.cgi	A statistical method to predict regic of weak packing density (less tha 8 Å between heavy atoms of non-adjacent residues).
GlobPlot2 ¹³³	globplot.embl.de/	A simple method based on several hydrophobicity scales to predict regions of missing coordinates ar loops with high B-factors.
iPDA ¹⁴²	biominer.cse.yzu.edu.tw/ipda	A successor of DisPSSMP. Incorporates information about sequence conservation, predicted secondary structure, sequence complexity and hydrophobic clusters.
IUPred ¹⁴³	iupred.enzim.hu/	Estimates pairwise interaction energies using a statistical potent Disordered regions tend to exhib poor inter-residue contact capaci
NORSp ¹⁴⁴	www.rostlab.org/services/NORSp/	Predicts long regions exposed to th solvent, with no regular seconda structure.
PONDR ^{®145}	www.pondr.com/	A commercial package containing several predictors based on FFNN

Table 2.4 Software for disorder prediction

Summary 53

P1: OTA

Program	URL (http://)	Short description
POODLE (S,L,W) ^{146,147}	mbs.cbrc.jp/poodle/poodle.html	L predicts long disorder using an SVM. S adds analysis of PSSMs generated by PSI-BLAST to detect short disorder. W uses Joachims' spectral graph transducer (SGT) to classify entire proteins as either disordered or ordered.
PrDOS ¹⁴⁸	prdos.hgc.jp	Predicts missing coordinates using SVM and PSSMs from PSI-BLAST.
PreLink ¹⁴⁹	genomics.eu.org/prelink/	Identifies regions with biased composition and poor in hydrophobic clusters to predict regions with missing coordinates.
RONN ¹⁵⁰	www.strubi.ox.ac.uk/RONN	Predicts missing coordinates using an ANN.
SEG ¹³¹	mendel.imp.ac.at/METHODS/ seg.server.html	Ancient precursor of modern disorder predictors, identifies regions of low sequence complexity.
SPRITZ ¹⁵¹	distill.ucd.ie/spritz/	Predicts long and short disorder (missing coordinates) using two separate SVMs. Utilizes secondary structure predicted by PORTER.
Grishin Lab Disorder Predictor ¹⁵²	prodata.swmed.edu/disorder/ disorder_prediction/predict.cgi	Predicts missing coordinates based on a PSSM and optimized propensities of amino acid residues toward disorder.
GeneSilico ⁶⁴	genesilico.pl/meta2/	A metaserver that predicts different types of disorder using weighted consensus of several methods.

30 31

employ combinations of various machine-learning techniques and utilize information from 32 multiple sequence alignments (e.g. to identify correlated Cys pairs) and predicted secondary 33 structure. Their reported prediction accuracy reaches 80%, but different methods have not 34 been compared directly with each other on the same test set. 35

37 2.8 Summary 38

39

36

To obtain better quality of secondary structure prediction, when no related structures are 40 41 known, it is advisable to follow some general rules:

42 First, it is important to use multiple sequence information, but if target sequence shows 43 high similarity to none or to only a few other proteins it is worth trying to search different databases (e.g. not only the non-redundant database at the NCBI, but also protein 44 sequences deduced from unfinished genomes and environmental sequencing projects) to 45 46 find moderately divergent sequences that can be used to build MSA (see also Chapter 1

chap02 JWBK331-Bujnicki October 1, 2008 8:21 Printer: Yet to come

	URL (http://)	
	SAM-T06 ¹⁵⁷	www.soe.ucsc.edu/research/compbio/SAM_T06/T06-guery.htm
	GPCPred ¹⁵⁹	sbcweb.pdc.kth.se/cgi-bin/maccallr/gpcpred/submit.pl
	PROFcon ¹⁶⁰	www.predictprotein.org/submit_profcon.html
	CONpro ¹⁶¹	www.ics.uci.edu/~baldig/scratch/
	SVMcon ¹⁶²	www.bioinfotool.org/symcon.html
	CRNPRED ⁴⁷	ftp.bioinformatics.org/pub/crnpred/
0	CMAPpro ⁶⁹	scratch.proteomics.ics.uci.edu/
1	Distill ¹⁵¹	distill.ucd.ie/distill/
י ר	PoCM ¹⁶³	foo.maths.uq.edu.au/~nick/Protein/contact.html
2	CMA ¹⁶⁴	ligin.weizmann.ac.il/cma/
3	HMMSTR-CM ¹⁶⁵	www.bioinfo.rpi.edu/~bystrc/hmmstr/server.php
4	BETApro $(\beta)^{166}$	www.ics.uci.edu/~baldig/betasheet.html
5	DIANNA (C-C) ¹⁶⁷	clavius.bc.edu/~clotelab/DiANNA/
6	Dlpro 2.0 (C-C) ¹⁶⁸	contact.ics.uci.edu/bridge.html
7	DISULFIND (C-C) ¹⁶⁹	cassandra.dsi.unifi.it/disulfind/index.php
8	GDAP (C-C) ¹⁷⁰	www.doe-mbi.ucla.edu/~boconnor/GDAP/
9	DCON (C-C) ¹⁷¹	gpcr.biocomp.unibo.it/cgi/predictors/cys-cys/pred_dconcgi.cgi
0	DISULFIDE (C-C) ¹⁷²	foo.maths.uq.edu.au/~huber/disulfide/

Table 2.5 Software for prediction of long-range contacts and disulfide bonds

22 23

P1: OTA

1

24 by Kaminska et al.). If the sequence is a true 'ORFan' with no homologs, a specialized method IPSSP33 (Table 2.1) may be used. Alignments that include remotely related se-25 quences should be inspected in the most divergent regions, and sequences that cannot 26 27 be aligned with confidence should be removed. In case of secondary structure prediction algorithms that do not accept an MSA as an input (but e.g. construct one from scratch by 28 29 themselves), secondary structure may be predicted independently for a few homologous 30 sequences and checked for mutual consistency. Correctly aligned positions should display 31 similar structure; therefore regions of low sequence similarity with different predictions 32 should be checked for possible errors in MSA.

Second, we recommend using meta-servers for disorder and secondary structure pre-33 diction, because combining results of several best prediction methods into a consensus 34 prediction is more reliable than relying on any individual method alone. Agreement be-35 tween methods usually indicates confident prediction, while disagreement may indicate 36 37 various things: different peculiarities of methods used, poor alignment in the input data, 38 and/or non-standard type of secondary structure, such as surface-exposed β -strands with 39 bulges that are often mispredicted as helices due to their irregular pattern of hydrophobic 40 and hydrophilic residues. It is also important to remember that specialized methods are 41 usually better for predicting particular types of structure than general-purpose methods for 42 secondary structure prediction. Therefore, it may be useful to use methods for prediction 43 of TM regions to pre-screen sequence for non-globular elements and 'mask' them before 44 considering regular secondary structure prediction.

Third, when selecting 'best' methods for consensus prediction it is important to remember that many authors use different benchmarks to assess their methods and that many

References 55

published accuracies have been shown to be overestimated, when these methods were as-1 sessed in rigorous blind tests on standard benchmarks, such as within CASP¹⁷³ or EVA.¹⁷⁴ 2 Although secondary structure predictions are no longer assessed in CASP, the EVA website 3 4 (http://cubic.bioc.columbia.edu/eva/) is updated automatically each week, to cope with the 5 large number of existing prediction servers and the constant changes in the prediction methods. EVA currently assesses servers for secondary structure prediction, contact prediction, 6 7 comparative protein structure modeling and threading/fold recognition. The identity of the test set assures that the competition is fair, while a large sample of targets assures that 8 methods are compared reliably. 9

10 Fourth, we recommend making simultaneous predictions of secondary structure, solvent accessibility, and disorder, as they usually reinforce each other (e.g. regions of disorder 11 usually exhibit little tendency to form secondary structure and their residues are predicted 12 to be largely solvent-exposed). However, discrepancies in this regard (i.e. presence of 13 confidently predicted secondary structure and/or buried residues within regions of disorder) 14 15 may indicate interesting structural and functional elements, such as partially folded molten 16 globule-like structures or candidates for linear motifs (see Chapter 1 by Kaminska et al. in this volume). Thus, again we recommend using meta-servers for making predictions 17 on the level of primary and secondary structure, in particular if they are going to be used 18 as restraints for modeling of protein tertiary structure (see articles by Kosinski et al. and 19 Gront et al. in this volume). 20

21

22

23 Acknowledgements

24

We thank present and former members of the Bujnicki lab in IIMCB and at the UAM for stimulating discussions and collaboration in development of some of the methods mentioned in this article. The authors acknowledge the support from past and current grants for the development of structure prediction methods from Polish Ministry of Science, NIH, Framework Programme of the EU, EMBO, and HHMI.

30 31

³² References

- 33
- ³⁴
 ³⁵
 C.J. Epstein, R.F. Goldberger, and C.B. Anfinsen, The genetic control of tertiary protein structure. Model systems. *Cold Spring Harb Symp Quant Biol*, 28, 439–449 (1963).
- C.B. Anfinsen, Studies on the principles that govern the folding of protein chains. Nobel lecture, December 11, 1972, in *Nobel Lectures, Chemistry 1971–1980*, T. Frängsmyr (Ed), World Scientific Publishing Co., Singapore, 1993.
- 39
 31. H.J. Dyson, and P.E. Wright, Intrinsically unstructured proteins and their functions, *Nat Rev Mol Cell Biol*, 6, 197–208 (2005).
- 40
 4. L. Pauling, R.B. Corey, and H.R. Branson, The structure of proteins; Two hydrogen-bonded
 41
 41 helical configurations of the polypeptide chain, *Proc Natl Acad Sci U S A*, 37, 205–211 (1951).
- 42 5. L. Pauling, and R.B. Corey, The pleated sheet, a new layer configuration of polypeptide chains,
 43 *Proc Natl Acad Sci U S A*, 37, 251–256 (1951).
- 6. B. Rost, Review: Protein secondary structure prediction continues to rise, J Struct Biol, 134, 204–218 (2001).
- 45
 7. G.A. Chasse, A.M. Rodriguez, M.L. Mak, *et al.*, Peptide and protein folding, *J Mol Struct THEOCHEM*, **537**, 319–361 (2001).

1

2

3

4

5

8

56 First Steps of Protein Structure Prediction

- 8. G. Nemethy, and H.A. Scheraga, Theoretical determination of sterically allowed conformations of a polypeptide chain by a computer method, *Biopolymers*, **3**, 155–181 (1965).
- H.A. Scheraga, M. Khalili, and A. Liwo, Protein-folding dynamics: Overview of molecular simulation techniques, *Annu Rev Phys Chem*, 58, 57–83 (2007).
- 10. A.V. Guzzo, The influence of amino-acid sequence on protein structure, *Biophys J*, **5**, 809–822 (1965).
- W. Kabsch, and C. Sander, Dictionary of protein secondary structure: Pattern recognition of
 hydrogen-bonded and geometrical features, *Biopolymers*, 22, 2577–2637 (1983).
 - D. Frishman, and P. Argos, Knowledge-based protein secondary structure assignment, *Proteins*, 23, 566–579 (1995).
- ⁹
 ¹³. G. Labesse, N. Colloc'h, J. Pothier, and J.P. Mornon, P-Sea: A new efficient assignment of secondary structure from C alpha trace of proteins, *Comput Appl Biosci*, **13**, 291–295 (1997).
- 14. M.N. Fodje, and S. Al-Karadaghi, Occurrence, conformational features and amino acid propen sities for the Pi-helix, *Protein Eng*, **15**, 353–358 (2002).
- 15. J. Martin, G. Letellier, A. Marin, J.F. Taly, A.G. de Brevern, and J.F. Gibrat, Protein secondary structure assignment revisited: A detailed analysis of different assignment methods, *BMC Struct Biol*, 5, 17 (2005).
- 16. M.V. Cubellis, F. Cailliez, and S.C. Lovell, Secondary structure assignment that accurately reflects physical and evolutionary characteristics, *BMC Bioinformatics*, 6 Suppl 4, S8 (2005).
- ¹⁸ 17. I. Majumdar, S.S. Krishna, and N.V. Grishin, Palsse: A program to delineate linear secondary structural elements from protein structures, *BMC Bioinformatics*, **6**, 202 (2005).
- ¹⁹
 ¹⁸ W. Zhang, A.K. Dunker, and Y. Zhou, Assessing secondary structure assignment of protein structures by using pairwise sequence-alignment benchmarks, *Proteins*, (2007).
- R. Karchin, M. Cline, Y. Mandel-Gutfreund, and K. Karplus, Hidden Markov models that use predicted local structure for fold recognition: Alphabets of backbone geometry, *Proteins*, 51, 504–514 (2003).
- 20. C. Etchebest, C. Benros, S. Hazout, and A.G. de Brevern, A structural alphabet for local protein structures: Improved prediction methods, *Proteins*, (2005).
- 25 21. C. Bystroff, V. Thorsson, and D. Baker, Hmmstr: A hidden Markov model for local sequence 26 structure correlations in proteins, *J Mol Biol*, **301**, 173–190 (2000).
- 23. B. Xue, O. Dor, E. Faraggi, and Y. Zhou, Real-value prediction of backbone torsion angles, *Proteins*, (2008).
- 24. P.Y. Chou, and G.D. Fasman, Prediction of protein conformation, *Biochemistry*, 13, 222–245
 (1974).
- 25. P.Y. Chou, and G.D. Fasman, Conformational parameters for amino acids in helical, beta-sheet, and random coil regions calculated from proteins, *Biochemistry*, 13, 211–222 (1974).
- 26. J.M. Levin, S. Pascarella, P. Argos, and J. Garnier, Quantification of secondary structure prediction improvement using multiple alignments, *Protein Eng*, 6, 849–854 (1993).
- 27. H. Chen, F. Gu, and Z. Huang, Improved Chou-Fasman method for protein secondary structure
 prediction, *BMC Bioinformatics*, **7 Suppl 4**, S14 (2006).
- 28. A. Kloczkowski, K.L. Ting, R.L. Jernigan, and J. Garnier, Combining the Gor V algorithm with
 evolutionary information for protein secondary structure prediction from amino acid sequence,
 Proteins, 49, 154–166 (2002).
- ⁵⁹ 29. B. Rost, C. Sander, and R. Schneider, Phd: An automatic mail server for protein secondary structure prediction, *Comput Appl Biosci*, **10**, 53–60 (1994).
- 30. J. Park, K. Karplus, C. Barrett, *et al.*, Sequence comparisons using multiple sequences detect three times as many remote homologues as pairwise methods, *J Mol Biol*, **284**, 1201–1210 (1998).
- 31. D.T. Jones, Protein secondary structure prediction based on position-specific scoring matrices, *J Mol Biol*, **292**, 195–202 (1999).
- 45 32. J. Heringa, Computational methods for protein secondary structure prediction using multiple
 46 sequence alignments, *Curr Protein Pept Sci*, 1, 273–301 (2000).

References 57

1

2

9

- 33. Z. Aydin, Y. Altunbasak, and M. Borodovsky, Protein secondary structure prediction for a single-sequence using hidden semi-Markov models, *BMC Bioinformatics*, **7**, 178 (2006).
- 34. G. Pollastri, D. Przybylski, B. Rost, and P. Baldi, Improving the prediction of protein secondary structure in three and eight classes using recurrent neural networks and profiles, *Proteins*, 47, 228–235 (2002).
- 5 35. B. Rost, G. Yachdav, and J. Liu, The Predictprotein Server, *Nucleic Acids Res*, 32, W321–326
 (2004).
- J.M. Chandonia, and M. Karplus, Neural networks for secondary structure and structural class
 predictions, *Protein Sci*, 4, 275–285 (1995).
 - 37. G.P.S. Raghava. Apssp2: A combination method for protein secondary structure prediction based on neural network and example based learning. Casp5 Online Abstract. A-132, 2002.
- 38. D. Frishman, and P. Argos, Seventy-five percent accuracy in protein secondary structure pre diction, *Proteins*, 27, 329–335 (1997).
- 39. D.G. Kneller, F.E. Cohen, and R. Langridge, Improvements in protein secondary structure prediction by an enhanced neural network, *J Mol Biol*, **214**, 171–182 (1990).
- 40. G. Pollastri, and A. McLysaght, Porter: A new, accurate server for protein secondary structure prediction, *Bioinformatics*, **21**, 1719–1720 (2005).
- 41. H.N. Lin, J.M. Chang, K.P. Wu, T.Y. Sung, and W.L. Hsu, Hyprosp Ii-a knowledge-based
 hybrid method for protein secondary structure prediction based on local prediction confidence,
 Bioinformatics, 21, 3227–3233 (2005).
- 42. K. Karplus, S. Katzman, G. Shackleford, *et al.*, Sam-T04: What is new in protein-structure prediction for Casp6, *Proteins*, **61 Suppl 7**, 135–142 (2005).
- ¹⁹
 ⁴³. J.A. Cuff, and G.J. Barton, Application of multiple sequence alignment profiles to improve protein secondary structure prediction, *Proteins*, **40**, 502–511 (2000).
- 44. R. Adamczak, A. Porollo, and J. Meller, Combining prediction of secondary structure and solvent accessibility in proteins, *Proteins*, 59, 467–475 (2005).
- 46. K. Lin, V.A. Simossis, W.R. Taylor, and J. Heringa, A simple and fast secondary structure prediction method using hidden neural networks, *Bioinformatics*, 21, 152–159 (2005).
- 47. A.R. Kinjo, and K. Nishikawa, Crnpred: Highly accurate prediction of one-dimensional
 protein structures by large-scale critical random networks, *BMC Bioinformatics*, 7, 401 (2006).
- 48. J. Meiler, and D. Baker, Coupled prediction of protein secondary and tertiary structure, *Proc* Natl Acad Sci U S A, 100, 12105–12110 (2003).
- 49. H. Kaur, and G.P. Raghava, Prediction of beta-turns in proteins from multiple alignment using
 neural network, *Protein Sci*, 12, 627–634 (2003).
- 50. H. Kaur, and G.P. Raghava, A neural-network based method for prediction of gamma-turns in proteins from multiple sequence alignment, *Protein Sci*, 12, 923–929 (2003).
- ⁵⁵ 51. Q. Zhang, S. Yoon, and W.J. Welsh, Improved method for predicting {beta}-turn using support vector machine, *Bioinformatics*, (2005).
- 52. M. Kuhn, J. Meiler, and D. Baker, Strand-loop-strand motifs: prediction of hairpins and diverg ing turns in proteins, *Proteins*, 54, 282–288 (2004).
- 53. A. Lupas, M. Van Dyke, and J. Stock, Predicting coiled coils from protein sequences, *Science*,
 252, 1162–1164 (1991).
- 50 54. M. Delorenzi, and T. Speed, An Hmm model for coiled-coil domains and a comparison with Pssm-based predictions, *Bioinformatics*, 18, 617–625 (2002).
- 55. M. Gruber, J. Soding, and A.N. Lupas, Comparative analysis of coiled-coil prediction methods,
 J Struct Biol, 155, 140–145 (2006).
- 56. A.V. McDonnell, T. Jiang, A.E. Keating, and B. Berger, Paircoil2: Improved prediction of coiled coils from sequence, *Bioinformatics*, 22, 356–358 (2006).
- 57. E. Wolf, P.S. Kim, and B. Berger, Multicoil: A program for predicting two- and three-stranded coiled coils, *Protein Sci*, 6, 1179–1189 (1997).
- 58. M. Singh, B. Berger, and P.S. Kim, Learncoil-Vmf: Computational evidence for coiled-coil-like
 motifs in many viral membrane-fusion proteins, *J Mol Biol*, 290, 1031–1041 (1999).

1

7

8

58 First Steps of Protein Structure Prediction

- 59. J.A. Cuff, M.E. Clamp, A.S. Siddiqui, M. Finlay, and G.J. Barton, Jpred: A consensus secondary structure prediction server, Bioinformatics, 14, 892-893 (1998). 2
- Y. Guermeur, C. Geourjon, P. Gallinari, and G. Deleage, Improved performance in protein sec-3 ondary structure prediction by inhomogeneous score combination, Bioinformatics, 15, 413-421 4 (1999).
- 5 61. V.A. Eyrich, and B. Rost, Meta-Pp: Single interface to crucial prediction servers, Nucleic Acids 6 Res, **31**, 3308–3310 (2003).
 - 62. S. Montgomerie, S. Sundararaj, W.J. Gallin, and D.S. Wishart, Improving the accuracy of protein secondary structure prediction using structural alignment, BMC Bioinformatics, 7, 301 (2006).
- 9 63. G. Pollastri, A.J. Martin, C. Mooney, and A. Vullo, Accurate prediction of protein secondary 10 structure and solvent accessibility by consensus combiners of sequence and structure information, BMC Bioinformatics, 8, 201 (2007). 11
- 64. M.A. Kurowski, and J.M. Bujnicki, Genesilico protein structure prediction meta-server, Nucleic 12 Acids Res, 31, 3305–3307 (2003). 13
- 65. R. Adamczak, A. Porollo, and J. Meller, Accurate prediction of solvent accessibility using 14 neural networks-based regression, Proteins, 56, 753-767 (2004).
- 15 66. S. Ahmad, and M.M. Gromiha, Netasa: Neural network based prediction of solvent accessibility, 16 Bioinformatics, 18, 819-824 (2002).
- 67. S. Qin, Y. He, and X.M. Pan, Predicting protein secondary structure and solvent accessibility 17 with an improved multiple linear regression method, *Proteins*, **61**, 473–480 (2005). 18
- 68. H. Chen, and H.X. Zhou, Prediction of solvent accessibility and sites of deleterious mutations 19 from protein sequence, Nucleic Acids Res, 33, 3193-3199 (2005).
- 20 J. Cheng, A.Z. Randall, M.J. Sweredoski, and P. Baldi, Scratch: A protein structure and structural 21 feature prediction server, Nucleic Acids Res, 33, W72-76 (2005).
- 22 70. A. Garg, H. Kaur, and G.P. Raghava, Real value prediction of solvent accessibility in proteins using multiple sequence alignment and secondary structure, Proteins, 61, 318-324 (2005). 23
- 71. H. Kaur, and G.P. Raghava, Prediction of alpha-turns in proteins using Psi-blast profiles and 24 secondary structure information, Proteins, 55, 83-90 (2004).
- 25 72. Y. Wang, Z. Xue, and J. Xu, Better prediction of the location of alpha-turns in proteins with 26 support vector machine, Proteins, 65, 49-54 (2006).
- 73. H. Kaur, and G.P. Raghava, A neural network method for prediction of beta-turn types in 27 proteins using evolutionary information, *Bioinformatics*, 20, 2751–2758 (2004). 28
- 74. Y. Wang, Z.D. Xue, X.H. Shi, and J. Xu, Prediction of Pi-turns in proteins using Psi-blast 29 profiles and secondary structure information, Biochem Biophys Res Commun, 347, 574-580 30 (2006).
- 31 75. M.L. Wang, W.J. Li, M.L. Wang, and W.B. Xu, Support vector machines for prediction of peptidyl prolyl Cis/trans isomerization, J Pept Res, 63, 23-28 (2004). 32
- 76. J. Song, K. Burrage, Z. Yuan, and T. Huber, Prediction of Cis/trans isomerization in proteins 33 using Psi-blast profiles and secondary structure information, BMC Bioinformatics, 7, 124 34 (2006).
- 35 77. R.D. King, M. Ouali, A.T. Strong, et al., Is it better to combine predictions?, Protein Eng, 13, 36 15-19 (2000).
- 78. M. Albrecht, S.C. Tosatto, T. Lengauer, and G. Valle, Simple consensus procedures are effective 37 and sufficient in secondary structure prediction, Protein Eng, 16, 459-462 (2003). 38
- 79. D.L. Minor, Jr., and P.S. Kim, Context-dependent secondary structure formation of a designed 39 protein sequence, Nature, 380, 730-734 (1996).
- 40 80. J.G. Lees, and R.W. Janes, Combining sequence-based prediction methods and circular dichro-41 ism and infrared spectroscopic data to improve protein secondary structure determinations, BMC Bioinformatics, 9, 24 (2008). 42
- 81. J.L. Popot, and D.M. Engelman, Helical membrane protein folding, stability, and evolution, 43 Annu Rev Biochem, 69, 881-922 (2000). 44
- 82. G.E. Schulz, Transmembrane beta-barrel proteins, Adv Protein Chem, 63, 47-70 (2003).
- 45 83. K. Lundstrom, Structural genomics for membrane proteins, Cell Mol Life Sci, 63, 2597–2607 46 (2006).

References 59

P1: OTA

1 2

3

4

5

6

7

8

9

84.	A. Elofsson	n, and G.	von Heijn	e, Membrane	e protein	structure:	Prediction	versus re	eality,	Ann	и
	Rev Biocher	m, 76 , 12	5-140 (200	07).							
05	T 0'	10	11 ·· D	11	. 1	C 1				-	

- L. Sipos, and G. von Heijne, Predicting the topology of eukaryotic membrane proteins, *Eur J Biochem*, 213, 1333–1340 (1993).
- 86. L. Kall, A. Krogh, and E.L. Sonnhammer, A combined transmembrane topology and signal peptide prediction method, *J Mol Biol*, **338**, 1027–1036 (2004).
- 87. D.T. Jones, Improving the accuracy of transmembrane protein topology prediction using evolutionary information, *Bioinformatics*, **23**, 538–544 (2007).
- G.E. Tusnady, and I. Simon, The Hmmtop transmembrane topology prediction server, *Bioin-formatics*, 17, 849–850 (2001).
- 89. M. Cserzo, E. Wallin, I. Simon, G. von Heijne, and A. Elofsson, Prediction of transmembrane
 alpha-helices in prokaryotic membrane proteins: The dense alignment surface method, *Protein Eng*, **10**, 673–676 (1997).
- 90. B. Rost, P. Fariselli, and R. Casadio, Topology prediction for helical transmembrane proteins at 86% accuracy, *Protein Sci.*, 5, 1704–1718 (1996).
- 91. F. Milpetz, P. Argos, and B. Persson, Tmap: A new email and WWW service for membrane protein structural predictions, *Trends Biochem Sci*, 20, 204–205 (1995).
- 92. E.L. Sonnhammer, G. von Heijne, and A. Krogh, A hidden Markov model for predicting
 transmembrane helices in protein sequences, *Proc Int Conf Intell Syst Mol Biol*, 6, 175–182
 (1998).
- 93. K. Hofmann, and W. Stoffel, Tmbase a database of membrane spanning proteins segments, Biol Chem Hoppe Seyler, 347, 166 (1993).
- ¹⁹ 94. G. von Heijne, Membrane protein structure prediction. Hydrophobicity analysis and the positive-inside rule, *J Mol Biol*, 225, 487–494 (1992).
- 95. Y. Zhai, and M.H. Saier, Jr., A web-based program (What) for the simultaneous prediction of hydropathy, amphipathicity, secondary structure and transmembrane topology for a single protein sequence, *J Mol Microbiol Biotechnol*, 3, 501–502 (2001).
- 96. H. Zhou, and Y. Zhou, Predicting the topology of transmembrane helical proteins using mean burial propensity and a hidden-Markov-model-based method, *Protein Sci*, **12**, 1547–1555 (2003).
- 97. C. Pasquier, and S.J. Hamodrakas, An hierarchical artificial neural network system for the
 classification of transmembrane proteins, *Protein Eng*, **12**, 631–634 (1999).
- 98. P.G. Bagos, T.D. Liakopoulos, and S.J. Hamodrakas, Algorithms for incorporating prior topological information in Hmms: Application to transmembrane proteins, *BMC Bioinformatics*, 7, 189 (2006).
- 99. T.D. Liakopoulos, C. Pasquier, and S.J. Hamodrakas, A novel tool for the prediction of trans membrane protein topology based on a statistical analysis of the Swissprot database: The
 Orientm algorithm, *Protein Eng*, 14, 387–390 (2001).
- ³³ 100. P.D. Taylor, T.K. Attwood, and D.R. Flower, Bprompt: A consensus server for membrane protein prediction, *Nucleic Acids Res*, **31**, 3698–3700 (2003).
- 101. S. Lee, B. Lee, I. Jang, S. Kim, and J. Bhak, Localizome: A server for identifying transmembrane
 topologies and Tm helices of eukaryotic proteins utilizing domain information, *Nucleic Acids Res*, 34, W99–103 (2006).
- B. Cao, A. Porollo, R. Adamczak, M. Jarrell, and J. Meller, Enhanced recognition of protein transmembrane domains with prediction-based structural profiles, *Bioinformatics*, 22, 303–309 (2006).
 102. W. Zhao, and M. Land, and J. Meller, Enhanced recognition of protein transmembrane domains with prediction-based structural profiles, *Bioinformatics*, 22, 303–309
- ³⁹ 103. Y. Zhai, and M.H. Saier, Jr., The Beta-Barrel Finder (Bbf) Program, Allowing identification
 of outer membrane beta-barrel proteins encoded within prokaryotic genomes, *Protein Sci*, 11,
 2196–2207 (2002).
- Horizon 104. P.L. Martelli, P. Fariselli, A. Krogh, and R. Casadio, A sequence-profile-based Hmm for predicting and discriminating beta barrel membrane proteins, *Bioinformatics*, 18 Suppl 1, S46–53 (2002).
- I. Jacoboni, P.L. Martelli, P. Fariselli, V. De Pinto, and R. Casadio, Prediction of the transmem brane regions of beta-barrel membrane proteins with a neural network-based predictor, *Protein*
- 46 Sci, **10**, 779–787 (2001).

- 106. P.G. Bagos, T.D. Liakopoulos, I.C. Spyropoulos, and S.J. Hamodrakas, Pred-Tmbb: A web server for predicting the topology of beta-barrel outer membrane proteins, *Nucleic Acids Res*, 32, W400–404 (2004).
- H.R. Bigelow, D.S. Petrey, J. Liu, D. Przybylski, and B. Rost, Predicting transmembrane beta-barrels in proteomes, *Nucleic Acids Res*, 32, 2566–2577 (2004).
- M.M. Gromiha, S. Ahmad, and M. Suwa, Neural network-based prediction of transmembrane
 beta-strand segments in outer membrane proteins, *J Comput Chem*, 25, 762–767 (2004).
- 7 109. F.S. Berven, K. Flikka, H.B. Jensen, and I. Eidhammer, Bomp: A program to predict integral beta-barrel outer membrane proteins encoded within genomes of gram-negative bacteria, *Nucleic Acids Res*, 32, W394–399 (2004).
- ⁹
 ¹¹⁰. M. Arai, H. Mitsuke, M. Ikeda, *et al.*, Conpred Ii: A consensus prediction method for obtaining transmembrane topology models with high reliability, *Nucleic Acids Res*, **32**, W390–393 (2004).
- 111. M. Amico, M. Finelli, I. Rossi, *et al.*, Pongo: A web server for multiple predictions of all-alpha transmembrane proteins, *Nucleic Acids Res*, 34, W169–172 (2006).
- 112. H. Zhou, C. Zhang, S. Liu, and Y. Zhou, Web-based toolkits for topology prediction of transmembrane helical proteins, fold recognition, structure and binding scoring, folding-kinetics analysis and comparative analysis of domain combinations, *Nucleic Acids Res*, 33, W193–197 (2005).
- 113. N. Bhardwaj, R.V. Stahelin, R.E. Langlois, W. Cho, and H. Lu, Structural bioinformatics
 prediction of membrane-binding proteins, *J Mol Biol*, **359**, 486–495 (2006).
- 18
 114. P.G. Bagos, T.D. Liakopoulos, and S.J. Hamodrakas, Evaluation of methods for predicting the topology of beta-barrel outer membrane proteins and a consensus prediction method, *BMC Bioinformatics*, 6, 7 (2005).
- 115. M. Pellegrini-Calace, A. Carotti, and D.T. Jones, Folding in lipid membranes (film): A novel
 method for the prediction of small membrane protein 3D structures, *Proteins*, 50, 537–545
 (2003).
- 116. V. Yarov-Yarovoy, J. Schonbrun, and D. Baker, Multipass membrane protein structure prediction using Rosetta, *Proteins*, 62, 1010–1025 (2006).
 117. D. Bath, J. Schonbrun, and D. Baker, Toward high modulities and decises of tensor.
- ²⁴ 117. P. Barth, J. Schonbrun, and D. Baker, Toward high-resolution prediction and design of transmembrane helical protein structures, *Proc Natl Acad Sci U S A*, **104**, 15682–15687 (2007).
- 118. X. de la Cruz, E.G. Hutchinson, A. Shepherd, and J.M. Thornton, Toward predicting protein topology: An approach to identifying beta hairpins, *Proc Natl Acad Sci U S A*, 99, 11157–11162 (2002).
 110. M. Dinging and D. Dinging and Dinging and Dinging and D. Dinging and D. Dinging and D. Ding
- 119. M. Kumar, M. Bhasin, N.K. Natt, and G.P. Raghava, Bhairpred: Prediction of beta-hairpins in a protein from multiple alignment information using Ann and Svm techniques, *Nucleic Acids Res*, 33, W154–159 (2005).
- 120. A.N. Lupas, and M. Gruber, The structure of alpha-helical coiled coils, *Adv Protein Chem*, **70**,
 37–78 (2005).
- 121. A. Lobley, M.B. Swindells, C.A. Orengo, and D.T. Jones, Inferring function using patterns of native disorder in proteins, *PLoS Comput Biol*, 3, e162 (2007).
- H. Xie, S. Vucetic, L.M. Iakoucheva, *et al.*, Functional anthology of intrinsic disorder. 3.
 Ligands, post-translational modifications, and diseases associated with intrinsically disordered
 proteins, *J Proteome Res*, 6, 1917–1932 (2007).
- S. Vucetic, H. Xie, L.M. Iakoucheva, *et al.*, Functional anthology of intrinsic disorder. 2.
 Cellular components, domains, technical terms, developmental processes, and coding sequence diversities correlated with long disordered regions, *J Proteome Res*, 6, 1899–1916 (2007).
- H. Xie, S. Vucetic, L.M. Iakoucheva, *et al.*, Functional anthology of intrinsic disorder. 1.
 Biological processes and functions of proteins with long disordered regions, *J Proteome Res*,
 6, 1882–1898 (2007).
- 41 6, 1882–1898 (2007).
 42 125. V. Neduva, and R.B. Russell, Linear motifs: Evolutionary interaction switches, *FEBS Lett*, 579, 3342–3345 (2005).
- 126. M. Fuxreiter, P. Tompa, and I. Simon, Local structural disorder imparts plasticity on linear motifs, *Bioinformatics*, 23, 950–956 (2007).
- 45 127. Y. Zhang, B. Stec, and A. Godzik, Between order and disorder in protein structures: Analysis
 46 of 'dual personality' fragments in proteins, *Structure*, **15**, 1141–1147 (2007).

References 61

- 128. F. Ferron, S. Longhi, B. Canard, and D. Karlin, A practical overview of protein disorder prediction methods, *Proteins*, 65, 1–14 (2006).
- 3 129. Z. Dosztanyi, M. Sandor, P. Tompa, and I. Simon, Prediction of protein disorder at the domain level, *Curr Protein Pept Sci*, 8, 161–171 (2007).
- ⁴ 130. M. Sickmeier, J.A. Hamilton, T. LeGall, *et al.*, Disprot: The database of disordered proteins, *Nucleic Acids Res*, **35**, D786–793 (2007).
- I.31. J.C. Wootton, Non-globular domains in protein sequences: Automated segmentation using complexity measures, *Comput.Chem.*, 18, 269–285 (1994).
- ⁹ 133. R. Linding, L.J. Jensen, F. Diella, P. Bork, T.J. Gibson, and R.B. Russell, Protein disorder
 ¹⁰ prediction: Implications for structural proteomics, *Structure*, **11**, 1453–1459 (2003).
- 134. J.J. Ward, J.S. Sodhi, L.J. McGuffin, B.F. Buxton, and D.T. Jones, Prediction and functional analysis of native disorder in proteins from the three kingdoms of life, *J Mol Biol*, 337, 635–645 (2004).
 135. L.G. R. L. L. L. L. L. D. D. Life the second second
- 135. J. Cheng, M. Sweredoski, and P. Baldi, Accurate prediction of protein disordered regions by mining protein structure data, *Data Mining and Knowledge Discovery*, **11**, 213–222 (2005).
- 16 136. S. Vucetic, C.J. Brown, A.K. Dunker, and Z. Obradovic, Flavors of protein disorder, *Proteins*,
 52, 573–584 (2003).
- 18
 137. Z. Obradovic, K. Peng, S. Vucetic, P. Radivojac, and A.K. Dunker, Exploiting heterogeneous sequence properties improves prediction of protein disorder, *Proteins*, 61 Suppl 7, 176–182 (2005).
- 138. C.T. Su, C.Y. Chen, and Y.Y. Ou, Protein disorder prediction by condensed Pssm considering
 propensity for order or disorder, *BMC Bioinformatics*, 7, 319 (2006).
- 139. R.M. MacCallum. Order/disorder prediction with self organising maps. CASP 6 meeting,
 Online paper. http://www.forcasp.org/paper2127.html
- 140. J. Prilusky, C.E. Felder, T. Zeev-Ben-Mordehai, *et al.*, Foldindex: A simple tool to predict whether a given protein sequence is intrinsically unfolded, *Bioinformatics*, 21, 3435–3438 (2005).
- 141. O.V. Galzitskaya, S.O. Garbuzynskiy, and M.Y. Lobanov, Foldunfold: Web server for
 the prediction of disordered regions in protein chain, *Bioinformatics*, 22, 2948–2949
 (2006).
- 142. C.T. Su, C.Y. Chen, and C.M. Hsu, Ipda: Integrated protein disorder analyzer, *Nucleic Acids Res*, 35, W465–472 (2007).
- 143. Z. Dosztanyi, V. Csizmok, P. Tompa, and I. Simon, Iupred: Web server for the prediction of
 intrinsically unstructured regions of proteins based on estimated energy content, *Bioinformatics*,
 21, 3433–3434 (2005).
- 144. J. Liu, and B. Rost, Norsp: Predictions of long regions without regular secondary structure, Nucleic Acids Res, 31, 3833–3835 (2003).
- ³⁴
 ³⁵
 ³⁶ P. Romero, Z. Obradovic, X. Li, E.C. Garner, C.J. Brown, and A.K. Dunker, Sequence complexity of disordered protein, *Proteins*, 42, 38–48 (2001).
- 146. K. Shimizu, S. Hirose, and T. Noguchi, Poodle-S: Web application for predicting protein
 disorder by using physicochemical features and reduced amino acid set of a position-specific
 scoring matrix, *Bioinformatics*, 23, 2337–2338 (2007).
- 147. S. Hirose, K. Shimizu, S. Kanai, Y. Kuroda, and T. Noguchi, Poodle-L: A two-level Svm prediction system for reliably predicting long disordered regions, *Bioinformatics*, 23, 2046–2053 (2007).
- 148. T. Ishida, and K. Kinoshita, Prdos: Prediction of disordered protein regions from amino acid
 sequence, *Nucleic Acids Res*, 35, W460–464 (2007).
- 43
 449. K. Coeytaux, and A. Poupon, Prediction of unfolded segments in a protein sequence based on amino acid composition, *Bioinformatics*, 21, 1891–1900 (2005).
- tailing acta composition, *Distrigormatics*, 22, 1091 1966 (2000).
 150. Z.R. Yang, R. Thomson, P. McNeil, and R.M. Esnouf, Ronn: The bio-basis function neural network technique applied to the detection of natively disordered regions in proteins, *Bioinfor*-
- 46 *matics*, **21**, 3369–3376 (2005).

- 62 First Steps of Protein Structure Prediction
- 151. A. Vullo, O. Bortolami, G. Pollastri, and S.C. Tosatto, Spritz: A server for the prediction of intrinsically disordered regions in protein sequences using kernel machines, *Nucleic Acids Res*, 34, W164–168 (2006).
- 152. N.B. Holladay, L.N. Kinch, and N.V. Grishin, Optimization of linear disorder predictors yields tight association between crystallographic disorder and hydrophobicity, *Protein Sci*, 16, 2140–2152 (2007).
- Is3. L. Bordoli, F. Kiefer, and T. Schwede, Assessment of disorder predictions in Casp7, *Proteins*,
 69 Suppl 8, 129–136 (2007).
- 8
 154. J. Skolnick, A. Kolinski, and A.R. Ortiz, Monsster: A method for folding globular proteins with a small number of distance restraints, *J Mol Biol*, 265, 217–241 (1997).
- ⁹ 155. M.J. Pietal, I. Tuszynska, and J.M. Bujnicki, Protmap2d: Visualization, comparison, and anal ¹⁰ ysis of 2D maps of protein structure, *Bioinformatics*, (2007).
- 11 156. J.M. Izarzugaza, O. Grana, M.L. Tress, A. Valencia, and N.D. Clarke, Assessment of intramolec ular contact predictions for Casp7, *Proteins*, 69 Suppl 8, 152–158 (2007).
- 157. G. Shackelford, and K. Karplus, Contact prediction using mutual information and neural nets, *Proteins*, 69 Suppl 8, 159–164 (2007).
 158. C.H. Tasi, C.H. Chan, B.L. Chan, C.Y. Kao, H.L. Liu, and I.B. Hay, Bioinformation and neural nets.
- ¹⁴ 158. C.H. Tsai, C.H. Chan, B.J. Chen, C.Y. Kao, H.L. Liu, and J.P. Hsu, Bioinformatics approaches
 ¹⁵ for disulfide connectivity prediction, *Curr Protein Pept Sci*, **8**, 243–260 (2007).
- 16 159. R.M. MacCallum, Striped Sheets and Protein Contact Prediction, *Bioinformatics*, 20 Suppl 1,
 i224–231 (2004).
- 160. M. Punta, and B. Rost, Profcon: Novel prediction of long-range contacts, *Bioinformatics*, 21, 2960–2968 (2005).
 161. C. Ballactri, D. Baldi, D. Farigalli, and P. Casadia. Prediction of accordination number and relative.
- ¹⁹ 161. G. Pollastri, P. Baldi, P. Fariselli, and R. Casadio, Prediction of coordination number and relative
 ²⁰ solvent accessibility in proteins, *Proteins*, 47, 142–153 (2002).
- 21 162. J. Cheng, and P. Baldi, Improved residue contact prediction using support vector machines and
 22 a large feature set, *BMC Bioinformatics*, 8, 113 (2007).
- 163. N. Hamilton, K. Burrage, M.A. Ragan, and T. Huber, Protein contact prediction using patterns of correlation, *Proteins*, 56, 679–684 (2004).
 164. V. Sakalan, F. Faral, S. Carran, et al. Space, A mits of table for participation and diction and distinguishing the second sec
- ²⁴ 164. V. Sobolev, E. Eyal, S. Gerzon, *et al.* Space: A suite of tools for protein structure prediction and analysis based on complementarity and environment, *Nucleic Acids Res*, **33**, W39–43 (2005).
- 165. Y. Shao, and C. Bystroff, Predicting interresidue contacts using templates and pathways, *Pro- teins*, **53 Suppl 6**, 497–502 (2003).
- 166. J. Cheng, and P. Baldi, Three-stage prediction of protein beta-sheets by neural networks, alignments and graph algorithms, *Bioinformatics*, 21 Suppl 1, i75–84 (2005).
- ²⁹ 167. F. Ferre, and P. Clote, Dianna 1.1: An extension of the Dianna web server for ternary cysteine classification, *Nucleic Acids Res*, **34**, W182–185 (2006).
- 168. J. Cheng, H. Saigo, and P. Baldi, Large-scale prediction of disulphide bridges using kernel
 methods, two-dimensional recursive neural networks, and weighted graph matching, *Proteins*,
 62, 617–629 (2006).
- 169. A. Ceroni, A. Passerini, A. Vullo, and P. Frasconi, Disulfind: A disulfide bonding state and cysteine connectivity prediction server, *Nucleic Acids Res*, 34, W177–181 (2006).
- ³⁵ 170. B.D. O'Connor, and T.O. Yeates, Gdap: A web tool for genome-wide protein disulfide bond
 prediction, *Nucleic Acids Res*, 32, W360–364 (2004).
- P. Fariselli, P. Martelli, and R. Casadio, A neural network based method for predicting the disulfide connectivity in proteins, in *Knowledge Based Intelligent Information Engineering Systems and Allied Technologies*, E. Damiani (ed.), International Conference on Knowledge-Based Intelligent Engineering (Kes 2002), IOS Press, Amsterdam, 2002.
- In Song, Z. Yuan, H. Tan, T. Huber, and K. Burrage, Predicting disulfide connectivity from protein sequence using multiple sequence feature vectors and secondary structure, *Bioinformatics*, 23, 3147–3154 (2007).
- 43 173. P. Aloy, A. Stark, C. Hadley, and R.B. Russell, Predictions without templates: New folds, secondary structure, and contacts in Casp5, *Proteins*, **53 Suppl 6**, 436–456 (2003).
- 44
 45
 46
 474. I.Y. Koh, V.A. Eyrich, M.A. Marti-Renom, *et al.* Eva: Evaluation of protein structure prediction servers, *Nucleic Acids Res*, **31**, 3311–3315 (2003).
- 46