

10**Scoring Functions for De Novo Protein Structure Prediction Revisited****Shing-Chung Ngan, Ling-Hong Hung, Tianyun Liu, and Ram Samudrala****Summary**

De novo protein structure prediction methods attempt to predict tertiary structures from sequences based on general principles that govern protein folding energetics and/or statistical tendencies of conformational features that native structures acquire, without the use of explicit templates. A general paradigm for de novo prediction involves sampling the conformational space, guided by scoring functions and other sequence-dependent biases, such that a large set of candidate (“decoy”) structures are generated, and then selecting native-like conformations from those decoys using scoring functions as well as conformer clustering. High-resolution refinement is sometimes used as a final step to fine-tune native-like structures. There are two major classes of scoring functions. Physics-based functions are based on mathematical models describing aspects of the known physics of molecular interaction. Knowledge-based functions are formed with statistical models capturing aspects of the properties of native protein conformations. We discuss the implementation and use of some of the scoring functions from these two classes for de novo structure prediction in this chapter.

Key Words: De novo; physics-based; knowledge-based; potential; protein folding.

1. Introduction

The success of large-scale genome sequencing efforts has spurred structural genomic initiatives, with the goal of determining as many protein folds as possible (*I-4*). At present, structural determination by crystallography and nuclear magnetic resonance (NMR) techniques are still slow and expensive in terms of manpower and resources, despite attempts to automate the

From: *Methods in Molecular Biology*, vol. 413: *Protein Structure Prediction, Second Edition*
Edited by: M. Zaki and C. Bystroff © Humana Press Inc., Totowa, NJ

01 processes. Computational structure prediction algorithms, while not providing
02 the accuracy of the traditional techniques, are extremely quick and inexpensive
03 and can provide useful low-resolution data for structure comparisons (5). Given
04 the immense number of structures that the structural genomic projects are
05 attempting to solve, there would be a considerable gain even if the computa-
06 tional structural prediction approach were applicable only to a subset of proteins.

07 Most current research in protein structure prediction is based on Anfinsen's
08 thermodynamic hypothesis that the native structure of a protein can be deter-
09 mined entirely from its amino acid sequence (6). The two main categories of
10 methods for predicting protein structure from sequence are comparative and de
11 novo modeling. In the comparative modeling category, the methodologies rely
12 on the presence of one or more evolutionarily related template protein structures
13 that are used to construct a model. Traditionally, the evolutionary relationship
14 can be deduced from sequence similarity (7–9) or by “threading” a sequence
15 against a library of structures and selecting the best match (10,11). However,
16 because of the improved sensitivity of the sequence similarity based methods,
17 the threading approach has essentially been supplanted (12,13). In the de novo
18 category, structure prediction methods attempt to predict tertiary structures from
19 sequences based on general principles that govern protein-folding energetics
20 and/or statistical tendencies of conformational features that native structures
21 acquire, without the use of explicit templates (14–16). A general paradigm for de
22 novo structure prediction involves sampling the conformational space, guided
23 with scoring functions and other sequence-dependent biases, such that a large
24 set of candidate (“decoy”) structures are generated, and then selecting native-
25 like conformations from those decoys using scoring functions and conformer
26 clustering as filters (17,18). As a final step, detailed energy potentials are
27 sometimes employed to perform high-resolution refinement on these native-like
28 structures. Although the first papers on protein structure prediction appeared
29 some thirty years ago, de novo structure prediction remains a difficult challenge
30 today (12,13,19–21).

31 Scoring functions are employed in all stages of de novo structure prediction.
32 For the conformational search stage, a selected combination of scoring functions
33 approximates the energy landscape of the protein conformational space.
34 Search methodologies such as Monte Carlo simulated annealing (MCSA) and
35 molecular dynamics (MD) then generate trajectories leading to the minima of
36 the landscape. As the conformational search process needs to evaluate new
37 conformations encountered at every step, it is computationally intensive, and
38 the scoring functions used in this stage need to be computationally efficient.
39 Because none of the existing scoring functions can faithfully reproduce the

Scoring Functions for De Novo Protein Structure Prediction Revisited 243

01 true energy landscape of the conformational space, the search process often
02 leads to many false minima. Thus, one usually repeats the search process many
03 times with many different starting conditions and random seeds and obtains a
04 collection of candidate (“decoy”) structures. Then, a second set of (possibly
05 different) scoring functions are used in the decoy selection stage as filter
06 to eliminate non-native structures and retain the native-like ones. Conformer
07 clustering is often used as an additional step to further refine the collection
08 of the native-like conformations, followed by high-resolution refinement of
09 the few remaining candidate structures. Compared to the functions used in the
10 conformational search stage, the functions employed in the decoy selection
11 stage can be algorithmically more complex and more detailed, because the
12 number of candidate conformations to evaluate is much less than the number of
13 conformations encountered during the search process. Scoring functions used
14 in the high-resolution refinement stage are usually computational expensive
15 functions formulated from detailed mathematical models of short-range interac-
16 tions among atoms, allowing small local perturbations to fine-tune native-like
17 structures.

18 There are two broad classes of scoring functions. The first class of functions
19 are largely based on some aspects of the known physics of molecular inter-
20 action, such as the Van der Waals force, electrostatics, and the bending and
21 torsional forces, to determine the energy of a particular conformation (22–27).
22 The second class of functions is knowledge-based. Each of these knowledge-
23 based functions tries to capture some aspects of the properties of protein native
24 conformations, for example, the tendencies of certain residues to form contact
25 with one another or with the solvent. These knowledge-based functions are
26 usually compiled based on the statistics of a database of experimentally deter-
27 mined protein structures (28–34). In essence, the physics-based functions aim at
28 predicting the native structure of a given sequence by mimicking the energetics
29 of protein folding, whereas the knowledge-based functions bypass this inter-
30 mediate step by directly making statistical inferences on what are observed in
31 the database. Thus, the accuracy of the physics-based functions is determined
32 by how realistic the underlying physical models are, whereas the accuracy of
33 the knowledge-based functions is determined by the quality of the database as
34 well as the validity of the statistical assumptions.

35 In an earlier edition, we introduced scoring functions for de novo structure
36 prediction (35). In this chapter, we revisit physics-based and knowledge-based
37 scoring functions in the context of their roles in the current state of the art
38 structure prediction efforts. For the physics-based approach, the often-called
39 Class I force field, which is a common foundation among the widely used

01 molecular modeling force fields such as AMBER, CHARMM, OPLS, and
02 ENCAD, is discussed. Extensions to this force field and the role of modeling
03 solvent effects are also described. For the knowledge-based approach, we
04 study the Bayesian (conditional) probability formalism, using it to derive
05 the all-atom distance-dependent conditional probability discriminatory function
06 (RAPDF) (34). As an additional illustration, we delineate how one can combine
07 the Bayesian probability formalism with the neural network methodology to
08 construct neural network-based scoring functions. Then, a few other novel
09 knowledge-based scoring functions from the recent literature are highlighted.
10 Although it is not strictly a physics- or knowledge-based methodology, we
11 briefly discuss the use of conformer clustering to further enhance decoy
12 selection, as this technique has been shown to be useful in de novo structure
13 prediction. Finally, a sophisticated combined physics- and knowledge-based
14 potential used for high-resolution refinement is described.

15 **2. Theoretical Background and Methods**

16 **2.1. An Overview of Physics-Based Energy Functions**

17
18 Using quantum mechanical techniques, highly accurate energies can be
19 calculated for small organic and inorganic molecules (36,37). However, because
20 of their sizes and flexibility as well as the presence of solvent molecules,
21 proteins are much more difficult systems to model. The polar aqueous
22 environment vastly complicates the calculation of the electrostatic energies. For
23 instance, although there is no dispute that the largest driving force for protein
24 folding is the hydrophobic effect (38,39), which is associated with the decrease
25 of water entropy upon the solvation of non-polar groups, the exact structural
26 configuration of water molecules hydrating the solute remains unknown.

27 Although a full quantum mechanical treatment for a complete protein is not
28 feasible, approximations and simplifications can be made to derive empirical
29 physics-based energies. For example, hydrogen bond geometries that are applic-
30 able to those found in proteins can be determined from quantum mechanical
31 calculations of simple systems (40). Electrostatics calculations can be approx-
32 imated using classical point charges and modifying the dielectric constant to
33 approximate the polarizability of the protein and the solvent. Van der Waals
34 interactions are often approximated by Lennard–Jones potentials. The first use
35 of these approximate functions was in MD simulations, where fast and easily
36 calculated energies were required to determine the force fields. Some proto-
37 types for these types of energies are AMBER (41), CHARMM (42), OPLS (24),
38 and ENCAD (43). Parameters for these energies have been obtained by fitting
39 equations and results of computer simulations to data from experiments and

Scoring Functions for De Novo Protein Structure Prediction Revisited 245

01 from quantum mechanical calculations. These physics-based energies perform
 02 adequately for perturbations around a known native conformation (**44,45**),
 03 because the electrostatic and solvent-dependent information is implicit in the
 04 initial conformation itself. In combination with experimental NMR constraints
 05 (**46,47**), these force fields enable the determination of accurate structures,
 06 so long as there are enough constraints to define the fold. Unfortunately, in
 07 isolation, the solvent and electrostatic modeling is insufficient for full and
 08 reliable simulation of protein folding. As a result, producing accurate protein
 09 folding simulations from physics-based energies alone is still a very challenging
 10 and active area of research.

11 2.1.1. Class I Physics-Based Scoring Function and Its Possible Extensions

12
 13 As we have mentioned, AMBER (**41**), CHARMM (**42**), OPLS (**24**), and
 14 ENCAD (**43**) are some examples of the widely used physics-based force fields
 15 in protein-folding simulation. These force fields share a lot of commonalities
 16 in terms of the underlying physical models used and the mathematical approx-
 17 imations assumed. As an illustration, the AMBER force field, which was first
 18 developed under the direction of Professor Peter Kollman, has the following
 19 form:

$$20 \quad 21 \quad 22 \quad V_{\text{total}} = V_{\text{bond}} + V_{\text{angle}} + V_{\text{torsion}} + V_{\text{non-bond}} \quad (1)$$

23 Here, V_{total} is the total potential energy, V_{bond} is the bond stretching energy,
 24 V_{angle} the angle bending energy, and V_{torsion} the angle torsional energy. Together,
 25 V_{bond} , V_{angle} , and V_{torsion} are denoted as the bonded interactions terms. $V_{\text{non-bond}}$ is
 26 the energy for non-bonded interactions, consisting of a Van der Waals energy
 27 term V_{vdW} and an electrostatics term V_{elec} . Other widely used force fields such
 28 as CHARMM and OPLS employ similar bonded and non-bonded terms in their
 29 formulations, and Eq. 1 is often denoted as the Class I force field.

30 The bond-stretching energy (*see Fig. 1A*) is modeled by treating the bond
 31 as an idealized spring and using a simple quadratic function derivable from the
 32 Hooke's law.

$$33 \quad 34 \quad 35 \quad V_{\text{bond}} = k_{\text{bond}}(r - r_0)^2 \quad (2)$$

36 where k_{bond} is the bond-stretching constant, controlling the stiffness of the bond
 37 spring, and $(r - r_0)$ is the deviation of the bond length from its equilibrium
 38 distance. Unique numerical values for k_{bond} and r_0 are assigned each pair of
 39 atom types.

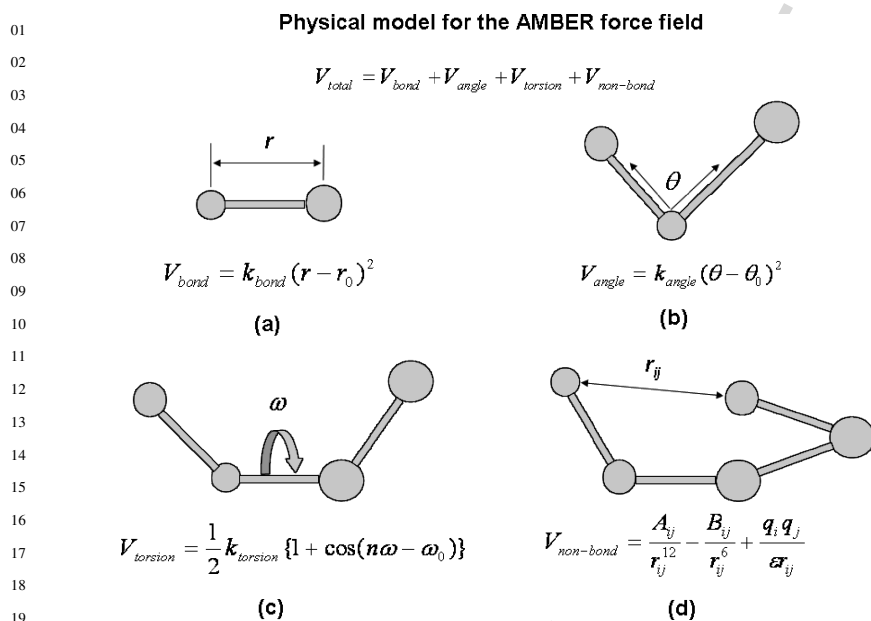


Fig. 1. The physical models for the AMBER molecular mechanics force field. Atoms and bonds are shown. (A) The physical model for bond stretching, (B) the model for angle bending, (C) the model for angle torsional energy, and (D) the model for electrostatics and Van der Waals forces.

The angle bending energy (see Fig. 1B) is similarly modeled by the Hooke's law.

$$V_{angle} = k_{angle} (\theta - \theta_0)^2 \quad (3)$$

where k_{angle} is the angle bending constant, controlling the stiffness of the angle spring. θ is the angle formed by the atom of interest with its two covalently bonded neighbors, and $(\theta - \theta_0)$ is the deviation of the angle from its equilibrium value in radians. Again, unique values for k_{angle} and θ_0 are determined for each bonded triplet of atom types.

The torsional energy (see Fig. 1C) is represented by an n -fold periodic function:

$$V_{torsion} = \frac{1}{2} k_{torsion} [1 + \cos(n\omega - \omega_0)] \quad (4)$$

Here, the torsional angle ω is the dihedral angle defined by a quartet of bonded atoms, and ω_0 is the reference angle. $k_{torsion}$ is a constant for the

Scoring Functions for De Novo Protein Structure Prediction Revisited 247

01 n -fold periodic interaction. n represents the periodicity of the torsional barrier,
 02 reflecting the intrinsic symmetry in the dihedral angle for the quartet of the
 03 bonded atoms. Unique values of k_{torsion} , n , and ω_0 are assigned to each bonded
 04 quartet of atom types. In practice, parameterization of torsional energies also
 05 corrects for bonding energy terms unaccounted for by the simple bending and
 06 stretching models. Additional torsional energy terms (denoted as “improper
 07 torsions” in the literature) can be added to ensure that subtle properties such as
 08 chirality and planarity are preserved.

09 For the non-bonded interactions, AMBER and other commonly used force
 10 fields employ a 6–12 Lennard–Jones potential to represent the Van der Waals
 11 interactions between two non-bonded atoms, and the Coulomb’s law to model
 12 the interactions of two charged atoms (see **Fig. 1D**):

$$V_{\text{non-bond}} = \left(\frac{A_{ij}}{r_{ij}^{12}} - \frac{B_{ij}}{r_{ij}^6} \right) + \left(\frac{q_i q_j}{\epsilon r_{ij}} \right) \quad (5)$$

16 The Van der Waals interaction consists of two components, a short-range
 17 attractive force that quickly vanishes when the distance between the interacting
 18 atoms, r_{ij} , is greater than a few Angstrom and an even shorter-range repulsive
 19 force that dominates when r_{ij} is less than the sum of their individual atomic
 20 radii. B_{ij} and A_{ij} in Eq. 5 control the attractive and the repulsive compo-
 21 nents of the steric potential. A_{ij} can be calculated from quantum mechanics
 22 considerations or measured from atomic polarizability experiments, and B_{ij}
 23 can be calculated from crystallographic data. For the electrostatics, interacting
 24 atoms are treated as point charges of q_i and q_j . The value of the dielectric
 25 constant ϵ accounts for the attenuation of electrostatic interaction by the polar
 26 environment. In more sophisticated solvent models, which are discussed later,
 27 the constant ϵ is replaced by a function dependent on r_{ij} . Earlier versions of
 28 AMBER had an explicit term to take into account hydrogen bonding. The latest
 29 versions incorporate hydrogen-bonding effects into the parameterization of the
 30 electrostatic and van der Waals terms, as these two terms are found to be able to
 31 sufficiently represent the distance and angle dependencies of hydrogen bonds
 32 in molecular mechanics modeling (48).

33 Currently, except in the high-resolution refinement stage, idealized backbone
 34 and side-chain bond lengths and angles are often used in de novo structure
 35 prediction. Hence, the energy associated with the bonded interactions terms
 36 V_{bond} , V_{angle} , and V_{torsion} can be regarded as constant. Improvement in structure
 37 prediction can conceivably be achieved by enhancing the physical models for
 38 the non-bonded terms. For example, one can replace the Van der Waals terms
 39 in Eq. 5 by a buffered 14–7 potential (49,50), by the Morse function (51),

01 or by the Buckingham–Fowler potential (52). The goal is to reduce the Pauli
02 exclusion barrier so as to allow sufficient sampling of conformations in the
03 neighborhood of the native structure during molecular mechanics or Monte
04 Carlo simulations.

05 For the electrostatic term, the physical model of fixed charges at atom
06 centers is found to be insufficient to describe charge polarization in the aqueous
07 environment. Examples of the more sophisticated electrostatics models involve
08 generalizing the point charge model with multi-center multi-pole expansion.
09 This can be done through the cumulative atomic multi-pole moment method,
10 the distributed multi-pole analysis, or an atoms-in-molecules-based multi-pole
11 moment method (53–55). Even though these types of model improvement
12 are computationally expensive, several groups have been making significant
13 progress in incorporating polarizable force fields for MD simulation of proteins.
14 For example, see refs. 56–58.

16 2.1.2. Protein Structures in Aqueous Environment

17
18 Protein structures are formed in the presence of aqueous environment, and
19 therefore, in order for the search of energy-minimized protein conformation
20 to be accurate, the effect of the solvent must be taken into account. Explicit
21 solvent models that simulate individual water molecules [for example, TIPS
22 (59,60), SPC (61), and F3C (62)] are too slow to be practicable for protein
23 structure prediction. Truncation of the non-bonded potentials such that interac-
24 tions beyond a fixed cutoff distance are ignored can improve speed. However,
25 it often leads to undesirable artifacts and reduced accuracy (63). Combining
26 Ewald’s approach with fast Fourier transform, Darden and his colleagues have
27 developed the particle mesh Ewald method to describe long-range interac-
28 tions more efficiently (64). However, direct simulation with explicit water is
29 still highly computational expensive even with this and other advances. On
30 the contrary, the effect of solvation can be modeled implicitly by averaging
31 solvent-solute interaction using mean field formulation and by decomposing
32 the solvation energy into an electrostatic component and a so-called non-polar
33 component, which accounts for everything else. For electrostatics, Poisson–
34 Boltzmann (65,66) models extend the simple Coulombic potential by allowing
35 charge distributions within the solute and having separate dielectrics for the
36 solvent and solute. Unfortunately, there are no general analytical solutions
37 for the Poisson–Boltzmann equation for irregular protein shapes and precise
38 numerical solutions (for example, by finite differences using GRASP/Delphi
39 (67)) can be very computationally expensive. Faster solutions can be obtained

Scoring Functions for De Novo Protein Structure Prediction Revisited 249

01 using generalized-Born (GB) approximations (68), which have been incorpo-
02 rated into MD simulations. For the non-polar term, which includes hydrophobic
03 interactions, the energy is usually modeled as a simple linear function of
04 solvent accessible area. The resulting generalized-Born/surface-area (GBSA)
05 models are more accurate than the simple non-bonded interaction terms and
06 can rival knowledge-based functions for scoring small loops in accuracy (69).
07 However, the amount of parameterization involved in GBSA models also rivals
08 that of knowledge-based energies. Recently, other approximate methods for
09 solving the Poisson–Boltzman equation may prove to be as or more accurate
10 with less parameterization (70). Besides the Poisson–Boltzmann and gener-
11 alized Born-type approaches, another category of implicit models describes the
12 solvent effect in terms of the dielectric screening of electrostatic interaction
13 within the protein molecule. For example, this can be done by defining the
14 dielectric coefficient as a simple function of distance (71,72) and as a more
15 detailed function involving solvent-excluded volume (73), the distance of a
16 charge from the protein surface, and the degree of exposure of a charge point
17 to the solvent (74).

18 In summary, the implicit solvent models are computationally much more
19 efficient than the explicit models. The tradeoff is the inability to represent
20 the detailed interaction structures between the solvent and the solute, which
21 can be essential in determining the overall energy landscape. Furthermore, the
22 lack of polarizability in the continuum solvent treatments precludes a flexible
23 description of charge distributions in the aqueous environment.

24 **2.2. An Overview of the Knowledge-Based Scoring Functions**

25
26 The physics-based functions are formulated from underlying approximate
27 physical models. In contrast, knowledge-based functions are derivable directly
28 from properties observed in known folded proteins (75). Although the basis of
29 the knowledge-based propensities is still physical, the statistical “black-box”
30 approach to the weighting of physical effects has proved to be more effective
31 than explicitly specifying the form and calculating the coefficients in traditional
32 physics-based energies. As a result, almost all of the most successful de novo
33 structure prediction techniques have both physics-based and knowledge-based
34 components.

35 The hydrophobic moment (76) is an example of a simple heuristic energy
36 function. It is analogous to the physical moment of inertia except that the
37 mass term is replaced by a measure of the hydrophobicity of the residue.
38 Minimization of this function leads to compact structures with hydrophobic
39 residues in the core. In general, any property that is differentially observed in

01 folded proteins and unfolded proteins can be converted into an energy function.
 02 Hidden Markov models (HMM), neural nets, support vector machines (SVM),
 03 and trial and error have been used to find such properties. A particularly
 04 useful class of knowledge-based functions is the pairwise distance preferences
 05 (**11,34,77**), which reflect proper packing. Consequently, the pairwise distance
 06 preference scoring functions can be found in many of the top-performing de
 07 novo methods, for example, ROSETTA (**16**), FRAGFOLD (**78**), TASSER (**79**),
 08 CABS (**80**), and PROTINFO (**81**).

10 2.2.1. Deriving Knowledge-Based Scoring Functions from the Bayesian 11 Probability Formalism

12 A majority of the knowledge-based scoring functions have their theoretical
 13 foundations rooted in the Bayesian (conditional) probability formalism. In such
 14 a formalism, we view a given set of conformations for a protein sequence as
 15 comprising a subset of correct conformations $\{C\}$ and a subset of incorrect
 16 conformations $\{I\}$. Furthermore, we consider a set of conformational properties,
 17 which can be any feature of protein structure that differs significantly between
 18 the subset of incorrect conformations and the subset of correct conformations.
 19 Examples are the preferences of some amino acid subsequences to exhibit
 20 certain torsion angles, to form contacts with other amino acid types, and so on.
 21 In this subheading, for the purpose of illustration, we focus on the set of inter-
 22 atomic distances within a structure $\{d_{ab}^{ij}\}$, where d_{ab}^{ij} is the distance between
 23 atoms numbers i and j , of type a and b . We want to determine $P(C|\{d_{ab}^{ij}\})$, the
 24 probability that the structure is a member of the “correct” subset, given that
 25 it contains the distances $\{d_{ab}^{ij}\}$. A standard way to achieve this is to express
 26 $P(C|\{d_{ab}^{ij}\})$ in terms of probabilities derivable from experimental structures,
 27 through the Bayes’ theorem:

$$28 \quad 29 \quad 30 \quad 31 \quad P(C|\{d_{ab}^{ij}\}) = P(C) \times \frac{P(\{d_{ab}^{ij}\}|C)}{P(\{d_{ab}^{ij}\})} \quad (6)$$

32 Here, $P(\{d_{ab}^{ij}\}|C)$ is the probability of observing the set of distances $\{d_{ab}^{ij}\}$
 33 in a correct structure. $P(\{d_{ab}^{ij}\})$ is the probability of observing such a set of
 34 distances in any correct or incorrect structure, and $P(C)$ is the probability that
 35 any structure picked at random belongs to the correct subset. $P(\{d_{ab}^{ij}\}|C)$ is
 36 regarded as a posterior probability in the sense that the underlying population
 37 for the probability distribution consists of structures that are already known
 38 to belong to the “correct” subset. On the contrary, $P(\{d_{ab}^{ij}\})$ is regarded as a
 39 prior probability in the sense that its underlying population is composed of

Scoring Functions for De Novo Protein Structure Prediction Revisited 251

01 structures whose class memberships have not yet been determined. We should
 02 note that both $P(\{d_{ab}^{ij}\}|C)$ and $P(\{d_{ab}^{ij}\})$ are highly difficult to compute, because
 03 the input arguments to these probability functions are the multitude of distance
 04 variables. A full model capturing the dependency among these variables would
 05 be extremely complex and would require a huge amount of training data to
 06 determine all the implicit parameters. Hence, to ensure computational feasibility
 07 of Eq. 6, one often makes the simplifying, albeit not strictly correct, assumption
 08 that the distances are statistically independent of one another, that is:

$$09 \quad P(\{d_{ab}^{ij}\}|C) = \prod_{i,j} P(d_{ab}^{ij}|C); P(\{d_{ab}^{ij}\}) = \prod_{i,j} P(d_{ab}^{ij}) \quad (7)$$

12 Then, combining Eqs. 6 and 7 gives us

$$13 \quad P(C|\{d_{ab}^{ij}\}) = P(C) \prod_{i,j} \frac{P(d_{ab}^{ij}|C)}{P(d_{ab}^{ij})} \quad (8)$$

14 For a given protein sequence, $P(C)$ is a constant independent of conformation
 15 and therefore can be omitted because we are only interested in selecting native-
 16 like conformations among decoys for a fixed protein sequence. Equation 8
 17 suggests a scoring function S , which is proportional to the negative log
 18 conditional probability that the given structure is correct, given a set of
 19 distances. Equation 8 suggests a scoring function S , which is proportional to the negative log
 20 conditional probability that the given structure is correct, given a set of
 21 distances.

$$22 \quad S(\{d_{ab}^{ij}\}) = \sum_{i,j} s(d_{ab}^{ij}); s(d_{ab}^{ij}) = -\log \left(\frac{P(d_{ab}^{ij}|C)}{P(d_{ab}^{ij})} \right) \quad (9)$$

23 An advantage of using Eq. 9 instead of Eq. 8 as a scoring function is that
 24 in the logarithm form, the pitfall of repeated multiplication of small numbers
 25 is eliminated, and therefore, it is easier to be implemented on the computer.

26 One can replace the set of distances $\{d_{ab}^{ij}\}$ with another type of conforma-
 27 tional property, say for example $\{m_a^i\}$, where m_a^i represents the value of that
 28 conformational property attained by residue number i of amino acid type a .
 29 This leads to another scoring function:

$$30 \quad S(\{m_k\}) = -\sum_k \log \left(\frac{P(m_k|C)}{P(m_k)} \right) \quad (10)$$

31 To gain an intuitive understanding of the scoring function, we note that if the
 32 chosen conformational property does not differ significantly between the subset
 33 of incorrect conformations and the subset of correct conformations, then the
 34
 35
 36
 37
 38
 39

01 values of $P(m_k|C)$ and $P(m_k)$ will tend to be close to each other. The resulting
 02 score S will always be close to 0 and is not an informative measure for decoy
 03 discrimination. On the contrary, if the conformational property is well chosen,
 04 that is, it differs significantly between incorrect and correct conformations, then
 05 for a native-like structure, $P(m_k|C)$ will tend to dominate $P(m_k)$, yielding a
 06 negative (good) score for S . On the contrary, for a non-native structure, the
 07 opposite occurs, yielding a positive (bad) score.

08 2.2.2. Compilation of the Probabilities

10 Before one can use Eq. 9 as a scoring function, the statistics for the posterior
 11 probability $P(d_{ab}^{ij}|C)$ and the prior probability $P(d_{ab}^{ij})$ need to be compiled.
 12 To compile the statistics for $P(d_{ab}^{ij}|C)$, we can tabulate the intra-molecular
 13 distances observed in a database of experimentally determined conformations.
 14 Such a database is usually extracted from the Protein Data Bank (PDB) (82,83).
 15 For example, one can proceed to select all the proteins from the PDB that also
 16 appear in the e-value filtered ASTRAL SCOP genetic domain sequence subset
 17 list with the threshold e-value set at 10^{-4} (84). Such an e-value is chosen,
 18 so that sampling bias (i.e., including too many homologous proteins) can be
 19 avoided. We then evaluate the quantity

$$20 \quad P(d_{ab}^{ij}|C) \equiv \frac{N(d_{ab})}{\sum_d N(d_{ab})} \quad (11)$$

23 where $N(d_{ab})$ is the number of occurrences of atom types a and b in a distance
 24 bin d in the database.

25 To compile the statistics of the prior probability $P(d_{ab}^{ij})$, we apply a formula
 26 similar to Eq. 11. But the question is: What would be an appropriate database
 27 from which to tabulate the counts? Samudrala and Moult (34) argued that
 28 methods employed for structure prediction usually produce compact models,
 29 whether the result is topologically correct or not. Thus, they consider a good
 30 choice of prior distribution to be found in the set of possible compact conformations
 31 and assume that averaging over different atom types in experimental
 32 conformations is an adequate representation of random arrangements of these
 33 atom types in any compact conformation. The probability $P(d_{ab})$ of finding
 34 atom types a and b in a distance bin d in any native-like or non-native compact
 35 conformation is thus approximated by:

$$36 \quad P(d_{ab}) = \frac{\sum_{ab} N(d_{ab})}{\sum_d \sum_{ab} N(d_{ab})} \quad (12)$$

Scoring Functions for De Novo Protein Structure Prediction Revisited 253

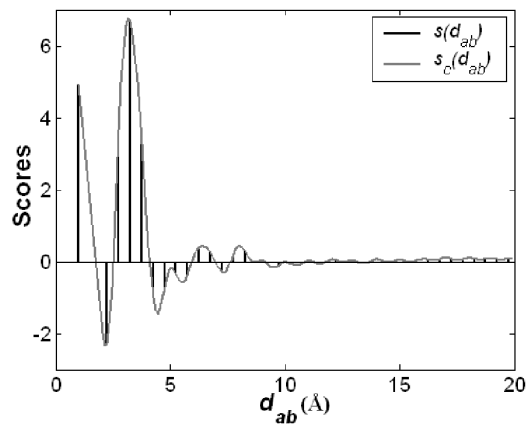
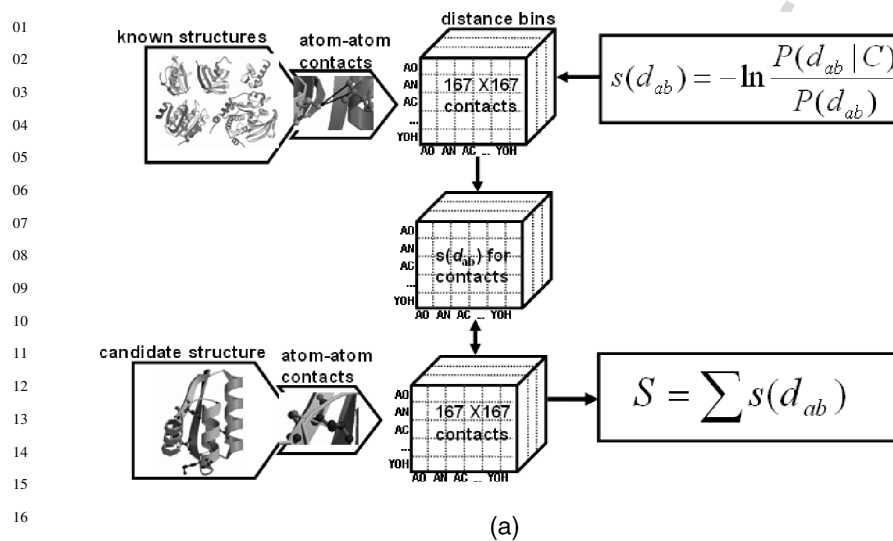
01 where $\sum_{ab} N(d_{ab})$ is the total number of contacts between all pairs of atom
02 types in a particular distance bin d , and the denominator is the total number of
03 contacts between all pairs of atom types summed over the distance bins d . The
04 pairwise distance preference function described in **Subheading 2.2.1.**, Eq. 9,
05 together with Eq. 11 and the prior distribution assumption of Eq. 12, is termed
06 the RAPDF in (34). **Figure 2A** highlights the essential components of this
07 scoring function.

08 Besides the above method of estimating prior distributions, various other
09 approaches have also been suggested. Subramaniam et al. (85) assumed that all
10 distances are equally probable, and Avbelj and Moulton (86) considered the set of
11 distances observed in some random coil model as appropriate. Lu and Skolnick
12 (87) employed a quasi-chemical approximation. Alternatively, Zhou and Zhou
13 (88) assumed that the residues follow uniform distribution everywhere in the
14 protein and developed a new reference state termed “distance-scaled, finite
15 ideal-gas reference state.”
16

2.2.3. A Pairwise Distance Scoring Function in Continuous Form

17
18
19 The RAPDF scoring function uses discrete distance bins to compile the
20 probability scores. Specifically, contact distances between 0 and 3 Å are
21 grouped into bin 1, 3 and 4 Å into bin 2, 4 and 5 Å into bin 3, and so on up to
22 the 20 Å cutoff. As a result, the score for observing any distance within a bin
23 width is the same for a given pair of atom types. However, the distance prefer-
24 ences between atom types should vary in a continuous manner as the distances
25 between the contacts vary. We can seek a function to interpolate between the
26 scores across the discrete bins such that the score for a given distance can be
27 uniquely defined. Several methods for interpolating discrete points, including
28 linear, polynomial, cubic spline, and band-limited interpolations, have been
29 tested for their efficacy to improve the discriminatory power of RAPDF. The
30 best among the tested methods is band-limited interpolation, derivable from
31 the Fourier Theorems. It assumes that the variation of the log-likelihood scores
32 fluctuates slowly enough such that the scores for any given distance can be
33 exactly reconstructed from the scores across the discrete bins.

34 Given a pair of atom types a and b at a particular distance, a “continuous” log-
35 likelihood score $s_c(d_{ab})$ can be calculated by interpolating between the scores
36 across the discrete bins of $s(d_{ab})$ through the Shannon’s sampling theorem,
37 resulting in a smooth curve (89). (see **Fig. 2B** for illustration.) Given an amino
38 acid sequence in a particular conformation, $s_c(d_{ab})$ of all contacts between pairs
39 of atom types at any distance within the 20 Å cutoff is summed to yield the total



(b)

33
34
35
36
37
38
39

Fig. 2. The all-atom distance-dependent conditional probability discriminatory function (RAPDF) and its extension, the interpolated RAPDF function. (A) The essential feature of the RAPDF scoring function. A matrix giving the log-likelihood scores for pairwise contact among different atom types at various discrete distance bins is computed using a database of known experimental structures. Then, given a candidate (“decoy”) structure, appropriate entries in the matrix can be extracted and summed to give a log-likelihood score for the structure. (B) The application of band-limited

Scoring Functions for De Novo Protein Structure Prediction Revisited 255

01 log-likelihood score to evaluate whether the conformation is native-like or not.
02 The interpolated RAPDF (IRAPDF) has been evaluated by various decoy sets.
03 Comparison between the IRAPDF and the RAPDF shows that the band-limited
04 interpolation leads to an improved discriminatory power.

05 **2.3. Neural Network Knowledge-Based Scoring Functions**

06
07 Rather than predicting whether an entire structure is native-like or not, neural
08 network algorithms are often used to predict the likelihood of occurrence of a
09 certain conformational property for each residue along a given protein sequence.
10 Examples of the properties are the tendencies of an amino acid to be exposed
11 or buried relative to the solvent (90–92), to be part of the helix, strand, or
12 coil local structures (93–95), the expected number of contacts a residue makes
13 with other residues (96–99), and so on. Usually, the conformational property of
14 interest is discretized into a number of states, and a neural network algorithm
15 returns numerical values which correlate with the probabilities of occurrences
16 of those states.

17 One can combine the neural network algorithms for predicting conforma-
18 tional properties with the Bayesian probability formalism that has been used to
19 construct various knowledge-based functions. This leads to a class of scoring
20 functions that give log-odd scores, indicating whether a given structure is
21 native-like or not, and that have in their core a neural network component.
22 In the following subheadings, we review a standard formulation of the neural
23 network algorithm that is used to predict conformational properties of residues
24 in a protein sequence. We then describe how the neural network and the
25 Bayesian frameworks are combined to form several neural network-based
26 scoring functions.

27 *2.3.1. Neural Network Algorithms for Predicting Local Structures*

28
29 For concreteness, we consider the prediction of the degree of solvent
30 accessibility of individual residues along a given protein sequence, with the
31 degree discretized into three states: low, medium, and high. The now standard
32 approach, introduced in ref. 93 and improved upon in ref. 94, uses a feed-
33 forward neural network. The input to the network is a window of sequence
34

35 ←
36 Fig. 2. interpolation to the discrete distance bins of the RAPDF function. The score
37 $s_c(d_{ab})$ of a given pair of atom types at any distance within the 20 Å cutoff can be
38 uniquely defined by interpolating across the discrete bins of $s(d_{ab})$. The resulting
39 scoring function is termed as the interpolated RAPDF (IRAPDF).

01 profile corresponding to a consecutive sequence of residues. Such a windowed
 02 sequence profile can be obtained by following a procedure described in **ref.**
 03 **94**. The protein sequence of interest is employed as input to PSI-BLAST (**100**),
 04 which generates a position-specific scoring matrix (PSSM) associated with that
 05 sequence. The PSSM consists of $20 \times M$ entries, where M being the length of
 06 the sequence, and each entry in a column gives the log-likelihood for one of the
 07 twenty possible amino acid substitutions for the residue position of interest. The
 08 standard logistic transform is then applied to each entry of the PSSM, so that
 09 these values are rescaled to the 0–1 range, appropriate to serve as neural network
 10 inputs. The neural network itself can consist of one or more hidden layers, and
 11 its output layer comprises three output units, representing the low, medium, and
 12 high solvent accessibility states, respectively. Training of the network is done
 13 with back-propagation (**101**), using the database of experimentally determined
 14 protein structures we have already described in **Subheading 2.2.2**. Given a
 15 window of sequence profile of the residue of interest (i.e., the sequence profile
 16 of the residue as well as those of the neighboring residues), the resulting neural
 17 network returns a numerical value in each output unit correlating with the
 18 probability with which the residue assumes the corresponding state.

19 2.3.2. Combining the Neural Network Algorithms with the Bayesian 20 Probability Formalism

21 To describe how one combines the Bayesian and the neural network frame-
 22 works to construct new scoring functions, for concreteness, suppose once again
 23 that the conformational property of interest is the degree of solvent accessi-
 24 bility. Using the language of the preceding subheadings, we want to calculate
 25 the probability that a given structure belongs to the subset of correct structures,
 26 given the associated conformational string $\{q_a^i\}$. Here, $q_a^i \in \{l, m, h\}$, where
 27 l represents low solvent accessibility state, m medium, and h high, i is the
 28 residue number, and a is the amino acid type. A scoring function described in
 29 Eq. 10 now takes the following form:
 30

$$31 \quad S(\{q_a^i\}) = - \sum_i \log \left[\frac{P(q_a^i|C)}{P(q_a^i)} \right] \quad (13)$$

32
 33
 34 $P(q_a^i|C)$ is simply the (posterior) probability of residue i taking on a particular
 35 solvent accessibility state q_a^i in a native structure. With an additional processing
 36 step involving the nearest-neighbor approach of Yi and Lander (**102**) to be
 37 discussed in detail in the next subheading, this probability can be estimated
 38 by using the neural network algorithm previously described. $P(q_a^i)$, on the
 39 contrary, is the (prior) probability that the residue is observed to assume the

Scoring Functions for De Novo Protein Structure Prediction Revisited 257

01 solvent accessibility state q_a^i in any native-like or non-native structure. It can
 02 be estimated using the formula

$$03 \quad P(q_a) \equiv \frac{N(q_a)}{\sum_{q \in \{l, m, h\}} N(q_a)} \quad (14)$$

04 where $N(q_a)$ is the number of occurrences of the amino acid type a taking on the
 05 solvent accessibility state q in some database of structures, and $\sum_{q \in \{l, m, h\}} N(q_a)$
 06 is the total number of occurrences of the amino acid type a in that database.
 07 Again, the question is: What is an appropriate database from which to tabulate
 08 the counts? We can use the same approach adopted by Samudrala and Moult in
 09 **ref. 34**, arguing that the set of possible compact conformations is a good choice
 10 of prior distribution. Then, the database to use will simply be the database
 11 of the experimentally determined structures. Alternatively, we can employ a
 12 database of decoy structures. Such a database can be created by applying a de
 13 novo conformational space sampling protocol to generate n decoy structures
 14 (for example, $n = 10$) for each protein sequence that appears in the database
 15 of the experimentally determined structures and then gathering the resulting
 16 decoys.

17 We note that as $P(q_a^i|C)$ is estimated by the neural network algorithm with a
 18 window of sequence profile as its input, the influence of the neighbors of residue
 19 i on its conformation is automatically taken into account. Thus, the posterior
 20 probability that residue i assumes a particular conformation is calculated in the
 21 context of its surrounding environment. In contrast, the probability distribution
 22 $P(q_a)$ is compiled on a “single-residue” basis. Thus, $P(q_a)$ can be viewed as
 23 the tendency of the amino acid type a to adopt a certain conformation averaged
 24 over the various types of neighborhood environments.

25 For further illustration, we generate a neural network-based Bayesian scoring
 26 function for each of the following conformational properties: the virtual torsion
 27 angle, the virtual bending angle, and the degree of solvent accessibility. The
 28 virtual torsion angle and the virtual bending angle are calculated by the DSSP
 29 program (**103**). Specifically, given a residue i of interest, the virtual torsion
 30 angle for i is the dihedral angle defined by the C_α atoms of residues $i - 1$,
 31 i , $i + 1$, and $i + 2$. The virtual bending angle is the bending angle defined by
 32 the C_α atoms of residues $i - 2$, i , and $i + 2$. Solvent accessibility is the residue
 33 water exposed surface in \AA^2 . To implement the scoring functions, the virtual
 34 torsion angle are manually divided into two discrete states, whereas the virtual
 35 bending angle and the degree of solvent exposure are each manually divided
 36 into three discrete states.

2.3.3. Training and Post-Processing of the Neural Network

The Stuttgart Neural Network Simulator (**104**) is a versatile and convenient tool to configure and train the neural networks for predicting the various conformational properties. The network configurations follow the description given in **Subheading 2.3.1**. The input layer receives a window of sequence profile. The window size typically ranges from 1 to 17 consecutive residues. The network has a single hidden layer and an output layer of two or three units representing two or three discrete states. See **Fig. 3** for an illustration.

We divide the database of experimentally determined structures into two equal subsets *A* and *B*, which are alternately used as the training and the test sets. The neural network training is done in batch mode using standard back-propagation, and the cycle of batch-mode training is repeated until the test error reaches a minimum. We note that two neural networks are obtained at the conclusion of the training—one (denoted as NN_A) trained with subset *A* and tested with subset *B* and another one (denoted as NN_B) trained with subset *B* and tested with subset *A*.

Given a residue of interest together with its windowed sequence profile, it is desired to extract from NN_A and NN_B the posterior probabilities with which the residue assumes each of the three states, say in the case of solvent accessibility prediction (two states in the case of virtual torsion angle prediction and three states in the case of virtual bending angle prediction). To this end, the nearest-neighbor approach of Yi and Lander (**102**) is employed: The output layer of NN_A gives a 3-tuple vector (s_{lA}, s_{mA}, s_{hA}) . The closeness of this vector with respect to vectors corresponding to all instances in the test set can be calculated through the Euclidean measure

$$((s_{lA} - s_{lA}^g)^2 + (s_{mA} - s_{mA}^g)^2 + (s_{hA} - s_{hA}^g)^2)^{1/2} \quad (15)$$

where *g* stands for instance *g* in the test set. The *k*-nearest neighbors [e.g., the closest 5% of all instances in the test set with respect to (s_{lA}, s_{mA}, s_{hA})] are then determined, and the actual solvent accessibility states of those nearest neighbors are tabulated, yielding the counts (c_{lA}, c_{mA}, c_{hA}) . The same procedure is repeated with NN_B . The probability that the residue of interest takes on each of the three states is thus estimated by

$$P(s_q) = \frac{c_{qA} + c_{qB}}{\sum_{r \in \{l, m, h\}} c_{rA} + c_{rB}} \quad (16)$$

where *q* stands for low, medium, or high accessibility state. Equation 16 supplies the posterior probabilities required in Eq. 13 for score calculation.

Scoring Functions for De Novo Protein Structure Prediction Revisited 259

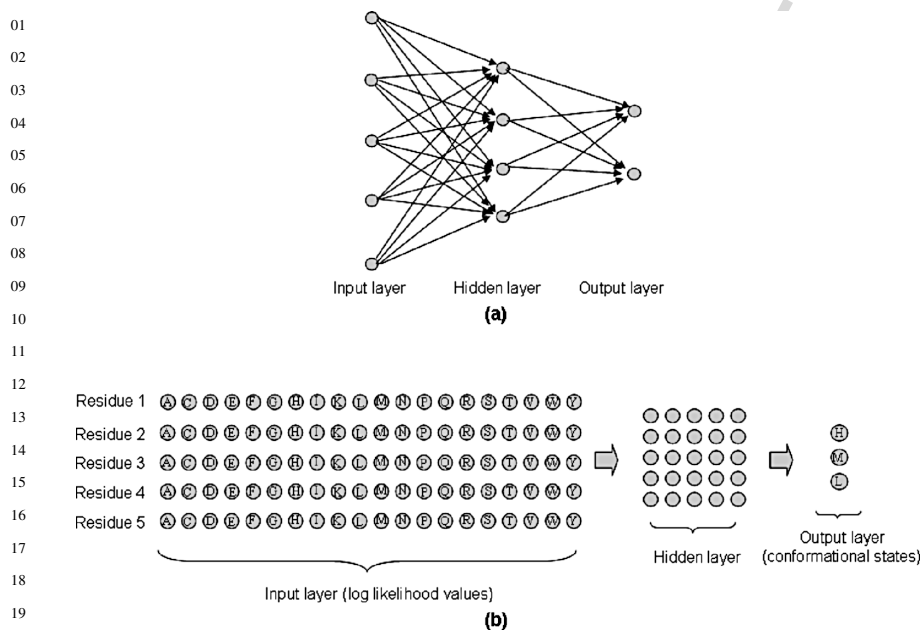


Fig. 3. Schematic diagrams of the neural networks used to predict conformational property given a sequence profile. **(A)** A fully connected neural network with input (5 units), hidden (4 units), and output (2 units) layers. Every unit in the input layers is connected with every unit in the hidden layers. The same holds true for the hidden and the output layers. **(B)** The typical size of a neural network we use for constructing the knowledge-based functions. In this example, the window size of the input sequence profile is five residues. Each residue provides twenty input units, representing the log-likelihood values for the twenty possible amino acid substitutions for that residue position. The hidden layer consists of 25 units. The output layer has three units. In the case of solvent accessibility prediction, these output units correspond to low, medium, and high solvent accessibility states, respectively. The input and the hidden layers, and the hidden and the output layers, are fully connected as in **(A)**, but for simplicity, the connections are not shown.

2.3.4. Decoy Sets and Evaluation of the Knowledge-Based Scoring Functions

One evaluates the usefulness of a scoring function by examining the ability of the scoring function to distinguish native-like conformations from non-native ones. This is achieved through generating test decoy sets and testing

01 the performance of the function on those sets. There are various approaches to
02 generate test decoys. For example, they can be created by sampling discrete-
03 state models starting from a native conformation (105), having amino acid
04 sequences with known folds mounted onto different folds (106,107), and using
05 crystal structures of various resolutions (85). Databases of test decoy sets
06 have been created to enable the evaluation of scoring functions on multiple
07 types of decoys (108–110). An approach most relevant to evaluating scoring
08 functions for de novo structure prediction is to create test decoys through de
09 novo conformational space sampling. A typical de novo conformational space
10 sampling protocol consists of an MCSA search procedure guided by a set of
11 energy functions, with move set based on lattice models (111,112), fragment
12 substitution (113,114), or continuous torsional distributions (81).

13 There are several commonly used measures for evaluating the usefulness of
14 scoring functions. The $\log P_{B1}$ measure is the log probability of selecting the
15 lowest C_α root mean square deviation (RMSD) conformation in a test decoy
16 set, calculated with the formula

$$\log P_{B1} = \log_{10} \left(\frac{R_i}{n} \right) \quad (17)$$

17
18
19
20 Here, R_i is the C_α RMSD rank of the best scoring conformation in the test
21 set of n decoys. The $\log P_{B10}$ measure is the log probability of selecting the
22 lowest C_α RMSD conformation among the top-10 best-scoring conformations,
23 that is, instead of using the RMSD rank of the best-scoring conformation, the
24 best RMSD rank achieved among the top-10 best-scoring conformations is used
25 as R_i in Eq. 17. The CC measure is the correlation coefficient between the C_α
26 RMSDs and the scores generated by the scoring function. The enrichment ratio
27 measure is the fraction enrichment of the top 10% lowest RMSD conformations
28 in the top 10% best scoring conformations. Specifically, after a scoring function
29 is applied to a test decoy set, we count the number of decoys (denoted as a),
30 which are in the top 10% in terms of both their scores and their C_α RMSDs
31 relative to the native structure. The expected number in a random distribution
32 is $10\% \times 10\% \times$ (number of decoys in the set) (denoted as b). The enrichment
33 ratio is a/b . A value above 1 indicates enrichment over the random distribution.
34 The four evaluation measures are illustrated in an example in Fig. 4.

35 To examine the utility of the knowledge-based scoring functions in decoy
36 discrimination, we apply both the RAPDF and the neural network-based
37 functions to 41 test decoy sets of varying quality generated with de novo
38 conformational space sampling. Each decoy set contains approximately 10,000
39 decoy conformations. Table 1 summarizes the PDB identifiers and the SCOP

Scoring Functions for De Novo Protein Structure Prediction Revisited 261

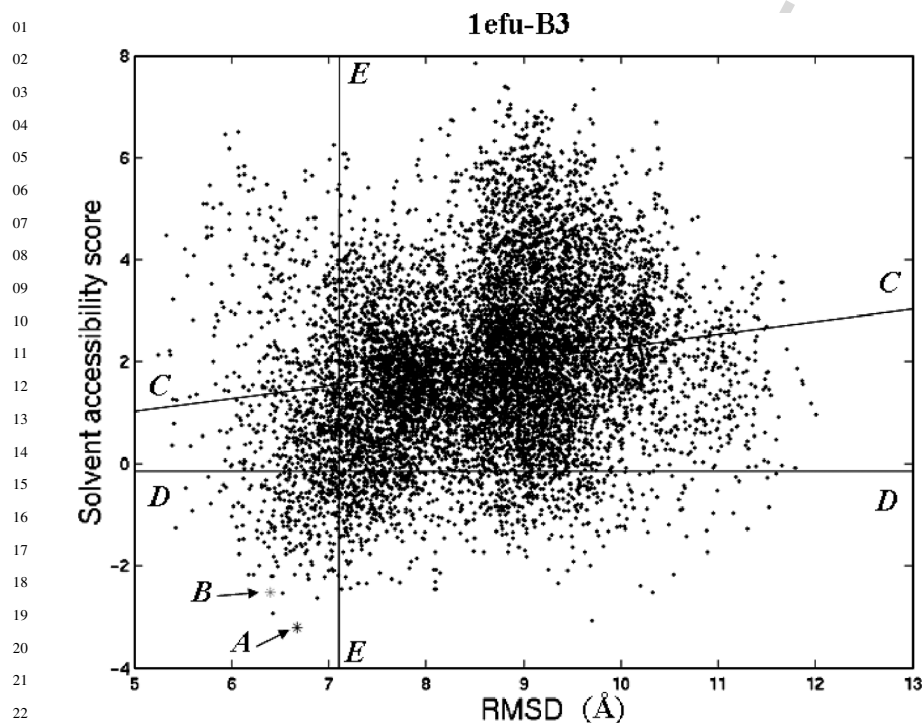


Fig. 4. Measures for evaluating scoring functions. $\log P_{B1}$ is the log probability of selecting the lowest C_α RMSD conformation in a test decoy set (point A), which is -1.42 in this example. $\log P_{B10}$ is the log probability of selecting the lowest C_α RMSD conformation among the top 10 best-scoring conformations in a test decoy set (point B), which is -1.76 in this example. The correlation coefficient between the C_α RMSDs and the scores is equal to the slope of line C-C and has the value of 0.25 in the present case. Line D-D represents the top 10% score cutoff for the decoy set. By counting the number of decoys below this line, which are also within the top 10% RMSD cutoff (left of line E-E), and dividing this number by the expected value for a random distribution, an enrichment ratio of 2.7 is obtained. Different measures are needed dependent on the specific purposes and roles of the scoring functions.

classifications of the 41 protein sequences used in generating the test decoy sets. Also included is the C_α RMSD of the best decoy relative to the corresponding native structure in each test set. Among them, fifteen test decoy sets have their best structures below 6 \AA C_α RMSD relative to their native conformations.

01 **Table 1**
 02 **List of the Protein Sequences Used in Generating the Test Decoy Sets**

03 Protein	04 SCOP classifications	Length	Minimum RMSD
05 1b0n-A2	a.35.1.3 (A:1–68)	68	2.729
06 1b33-N	d.30.1.1 (N:)	67	7.349
07 1b34-A	b.38.1.1 (A:)	80	7.943
08 1b4b-A	d.74.2.1 (A:)	71	5.506
09 1b79-A	a.81.1.1 (A:)	102	5.29
10 1ck9-A	d.79.3.1 (A:)	104	7.661
11 1ctf	d.45.1.1 (–)	68	4.37
12 1dgn-A	a.77.1.1 (A:)	89	4.482
13 1dj8-A	a.57.1.1 (A:)	79	5.092
14 1dtj-A	d.51.1.1 (A:)	74	4.902
15 1e68-A	a.64.2.1 (A:)	70	3.794
16 1eai-C	g.22.1.1 (C:)	61	6.914
17 1edz-A2	c.58.1.2 (A:3–148)	146	9.277
18 1efu-B3	a.5.2.2 (B:1–54)	54	5.247
19 1ev0-A	d.71.1.1 (A:)	58	6.641
20 1f53-A	b.11.1.4 (A:)	84	9.123
21 1fc3-A	a.4.6.3 (A:)	119	8.184
22 1fmt-A1	b.46.1.1 (A:207–314)	108	7.385
23 1g6e-A	b.11.1.6 (A:)	87	7.891
24 1g7d-A	a.71.1.1 (A:)	106	5.867
25 1goi-A1	b.72.2.1 (A:447–498)	52	6.111
26 1gut-A	b.40.6.1 (A:)	67	6.459
27 1h5p-A	b.99.1.1 (A:)	95	8.223
28 1h8a-C1	a.4.1.3 (C:87–143)	57	2.941
29 1ijy-A	a.141.1.1 (A:)	122	7.916
30 1ira-Y1	b.1.1.4 (Y:1–101)	101	8.317
31 1iwg-A1	d.58.44.1 (A:38–134)	97	5.7
32 1jju-A3	b.1.18.14 (A:274–351)	78	6.614
33 1jos-A	d.52.7.1 (A:)	100	5.302
34 1jyg-A	a.60.11.1 (A:)	69	3.471
35 1k2y-X2	c.84.1.1 (X:155–258)	104	6.889
36 1ktz-B	g.7.1.3 (B:)	106	8.586
37 1l9l-A	a.64.1.1 (A:)	74	4.041
38 1msp-A	b.1.11.2 (A:)	124	9.932
39 1n69-A	a.64.1.3 (A:)	78	6.753
40 1qu6-A1	d.50.1.1 (A:1–90)	90	8.597
41 1rie	b.33.1.1 (–)	127	9.548
42 1sra	a.39.1.3 (–)	151	8.781

Scoring Functions for De Novo Protein Structure Prediction Revisited 263

01	1sro	b.40.4.5 (-)	76	6.031
	2igd	d.15.7.1 (-)	61	6.508
02	7gat-A	g.39.1.1 (A:)	66	7.248

03
04 Each row lists the Protein Data Bank (PDB) identifier of the sequence, the SCOP classification,
05 the length of the protein sequence, and the C_{α} RMSD of the best decoy structure relative to the
06 native conformation in the test decoy set. Each test decoy set contains $\sim 10,000$ decoys. Fifteen
07 test decoy sets have their best structures below 6 \AA C_{α} RMSD relative to their corresponding
08 native conformations. Twenty-four test decoy sets have their best structures below 7 \AA C_{α} RMSD
09 relative to their corresponding native conformations.

10
11 Twenty-four decoy sets have their best structures below 7 \AA C_{α} RMSD relative
12 to their native conformations, and so on. For illustration purpose, we employ
13 the enrichment ratio measure to evaluate the scoring functions. The results are
14 displayed in **Fig. 5**. From the figure, we observe that the RAPDF function gives
15 uniform performance for decoy discrimination across decoy sets of different
16 quality, whereas the neural network-based scoring functions tend to perform
17 better for decoy sets with better quality.

18 **2.4. Some Other Knowledge-Based Scoring Functions in the Recent** 19 **Literature**

20
21 In the formulation of the RAPDF scoring function as well as of the other
22 pairwise distance preference functions described in **refs. 11,77,87** and **(88)**,
23 the solvation effect is not explicitly modeled. However, as we have previously
24 discussed, as protein folding occurs in the aqueous environment, a careful
25 accounting of the solvent effect is important in determining the native conform-
26 ation. In this regard, McConkey et al. **(115)** quantify contact surfaces of atoms
27 by integrating the solvent accessible surface and the inter-atomic contacts into
28 one quantity and construct an all-atom contact potential based on the contact
29 preferences of 167 residue-specific atom types with 168 possible contact types
30 (167 possible atom contact types and one solvent contact). They demonstrate
31 that this all-atom contact potential delivers satisfactory performance for distin-
32 guishing native conformations from decoy structures.

33 Another possible approach to augment the pairwise distance preference
34 scoring functions is by considering various multi-body geometric properties.
35 In **ref. 116**, a four-body SNAPP potential involving the tiling of protein struc-
36 tures with tetrahedra having the center of mass of each amino acid side-chain
37 at each vertex is introduced. This formulation results in 8855 possible tetra-
38 hedron types with the corresponding log-likelihoods computed from structural
39 databases. It is found that the SNAPP potential is accurate in predicting the

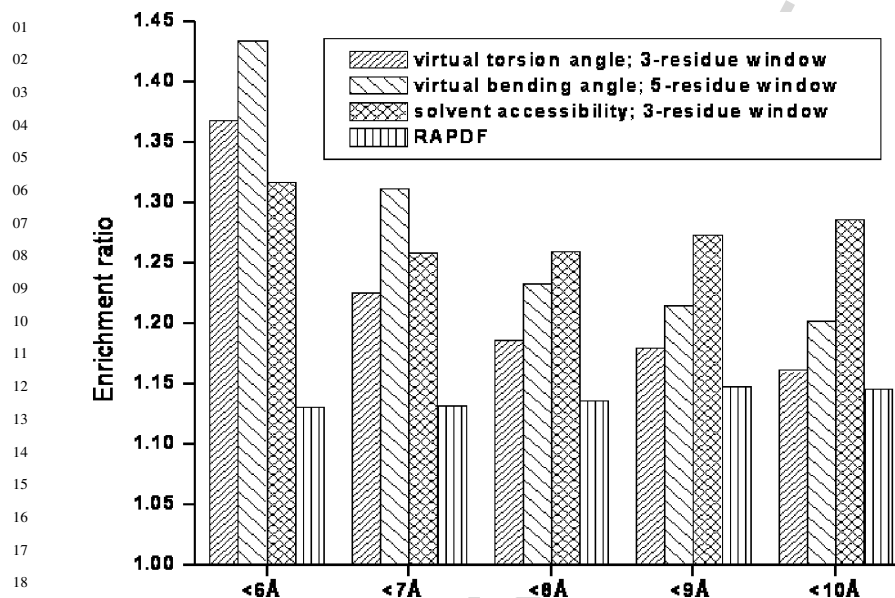


Fig. 5. Performances of the various knowledge-based scoring functions. The functions are evaluated using the average enrichment ratios on test decoy sets of varying quality. For example, the first four bars indicates the average enrichment ratios attained by the individual functions for the test decoy sets that contain structures of less than 6 Å C_{α} RMSD relative to the native conformations. The following scoring functions are examined in the figure: a neural network-based virtual torsion angle scoring function with a three-residue window; a neural network-based virtual bending angle scoring function with a five-residue window; a neural network-based solvent accessibility scoring function with a three-residue window; and the all-atom distance-dependent conditional probability function.

effects of hydrophobic core mutations. A similar four-body scoring function derived through the Delauney tessellation of side-chain centroids of amino acids is shown to be able to distinguish native conformation from partially unfolded and deliberately misfolded structures (117). On the basis of the work of Professor Banavar and his colleagues, Ngan et al. (118) construct a three-body knowledge-based potential involving the radii of curvature formed among triplets of residues in protein conformations. The resulting residue-triplet function is shown to be of utility in discriminating native-like conformations from non-native structures. Finally, Li et al. (119) introduce a knowledge-based

Scoring Functions for De Novo Protein Structure Prediction Revisited 265

01 scoring function based on the edge simplices from the alpha shape of the
02 protein structure. Formally, their statistical alpha contact potential is a two-body
03 scoring function, and their definition of contact is when atoms from non-bonded
04 residues share a Voronoi edge, with the edge at least partially contained in
05 the body of the protein. This formulation has the benefit of avoiding spurious
06 contact between two residues when a third residue is between them. The authors
07 have shown that the alpha contact potential performs comparably with other
08 atom-based potentials, while requiring fewer parameters.

09 In summary, the construction of a knowledge-based scoring function involves
10 the following steps: (1) selection of a conformational property that differs
11 between native-like and non-native structures; (2) compilation of the posterior
12 probability distributions of this conformational property by direct counting or
13 through statistical techniques such as neural network, based on a database of
14 experimentally determined structures; (3) derivation of the prior probability
15 distributions based on a database of decoy structures or through simplifying
16 assumptions such as the averaging-over-atom-types argument of Samudrala and
17 Moulton (34), the quasi-chemical approximation of Lu and Skolnick (87), or the
18 uniform distribution argument of Zhou and Zhou (88); and (4) formation of the
19 log-odd scores from the prior and posterior probabilities. Step 1 is perhaps the
20 most critical step and is largely dependent on one's insights into the physical
21 and chemical processes involved in protein folding and by trial and error. In
22 step 2, the selection of appropriate statistical techniques is heavily influenced
23 by the size and quality of the available data set, because these factors have a
24 direct impact on determining whether certain statistical assumptions (e.g., the
25 conditional independence assumption in Eq. 7) are needed.

27 **2.5. The Design of Decoy Filters**

28
29 As we have discussed, conformational search algorithms produce a multitude
30 of candidate conformations. Various scoring functions can be combined into a
31 filter to distill this vast collection of decoys, to retain those that are native-like.
32 An approach to constructing such a filter is to assign weights to the different
33 scoring functions, such that the resulting linear combination of the scores gives
34 the overall quantitative assessment of a decoy structure of interest. The weights
35 used in the linear combination can be derived by performing logistic regression
36 on test decoy sets. Specifically, native-like decoys (determined by a suitably
37 chosen C_{α} RMSD cutoff) in each test set are labeled as belonging to class 1,
38 and the rest labeled as class 0. The normalized scores for an individual decoy
39 become the independent variables (x_j ; $j = 1 \dots k$; $k =$ the total number of score

01 types), whereas its associated class label forms the dependent variable (p),
02 which are then used to fit the following equation to obtain the weights w_j s:

$$03 \quad \log\left(\frac{p}{1-p}\right) = \alpha + w_1x_{1,i} + \dots + w_kx_{k,i} \quad (18)$$

04
05
06 Here, α is a constant representing the intercept. i ranges from 1 to N , and
07 N is the total number of decoys. Normalization of a scoring function can be
08 achieved by subtracting its mean and dividing by its standard deviation, where
09 the mean and the standard deviation are computed over all decoys within a test
10 set, or by replacing the raw score of a decoy with its rank and then dividing
11 by the total number of decoys in the test set. Techniques such as leave-one-
12 out cross-validation and forward and backward stepwise regression can be
13 applied to determine which independent variables are helpful in assessing the
14 accuracy of a given decoy structure and which can be discarded. Essentially,
15 functions describing useful orthogonal characteristics of protein native conforma-
16 tions will receive large weights, whereas those that are less useful or containing
17 overlapping information will have smaller or zero weights. Finally, alternative
18 approach to performing logistic regression is also possible, for example, by
19 replacing it with machine-learning techniques such as the neural network or SVM.
20 The decision is again influenced by the size and quality of the available test data.

21 **2.6. Further Enhancement of Decoy Selection Through Conformer** 22 **Clustering and High-Resolution Refinement**

23
24 Conformer clustering and high-resolution refinement are often used as
25 additional steps in the decoy selection process to further refine the set of
26 native-like conformations retained by the decoy filter. The idea of conformer
27 clustering is based on the following observation: Conformers with correct folds
28 are in general similar to other conformers with correct folds. On the contrary,
29 it is unlikely that multiple conformers share the same mistake, and therefore,
30 conformers with incorrect folds are in general dissimilar to each other as well
31 as to conformers with correct folds. Hence, the conformers that are most similar
32 to the others, that is, those at the cluster centers of the conformational distri-
33 bution, will tend to be the correct ones. Various metrics are used to describe the
34 conformational distribution, including pairwise RMSD, pairwise RMSD with
35 cutoffs, and number of neighbors (16,120). Heuristic schemes such as k -mean
36 clustering, visual inspection following dimensionality reduction, and iterative
37 sampling (121) can be used to locate these cluster centers.

38 **Figure 6** illustrates the performance of a conformer-clustering algorithm [the
39 density score function available in the RAMP package (122)] in distinguishing

Scoring Functions for De Novo Protein Structure Prediction Revisited 267

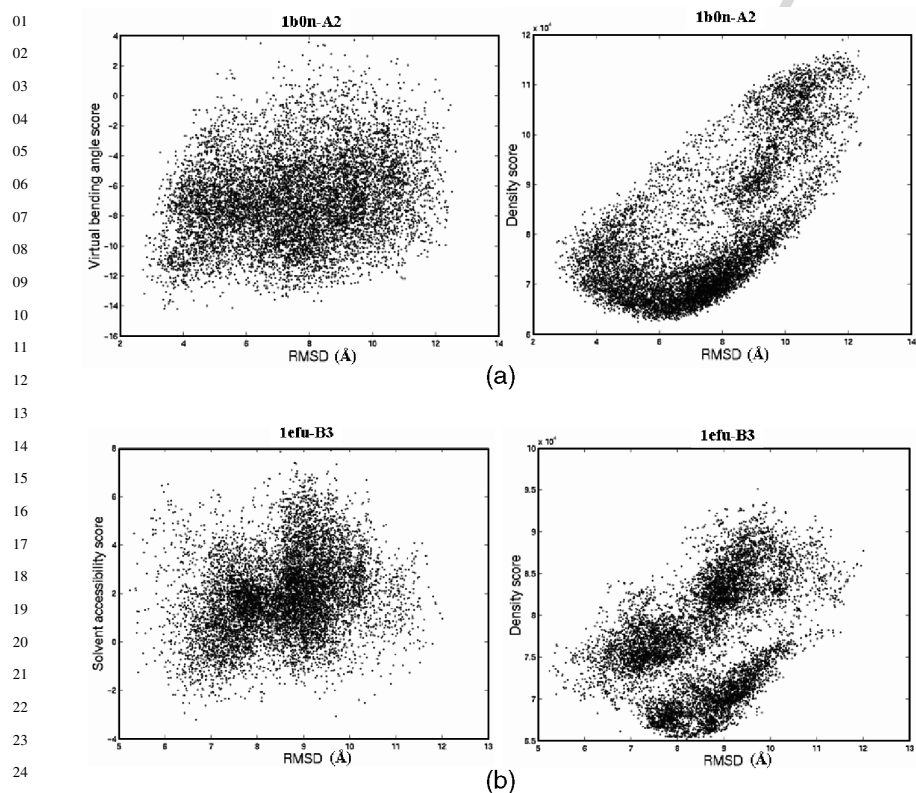


Fig. 6. The comparison of some knowledge-based scoring functions and the density score function in discriminating decoys. In **(A)**, the virtual bending angle scoring function is compared to the density score function, whereas in **(B)**, the solvent accessibility scoring function is compared to the density score function. The diagrams show that the density score function produces improved correlation between the C_{α} RMSDs and the scores in both cases, suggesting that conformer clustering is useful as a complementary step in decoy selection.

native-like structures from non-native conformations. Compared with the neural network-based virtual bending angle and solvent accessibility scoring functions, the density score function produces results that show improved correlation between the C_{α} RMSDs and the generated scores. This observation suggests that applying conformer clustering in addition to using scoring functions as filter can enhance the overall ability to select native-like structures from decoys.

01 The goal of high-resolution refinement is to further optimize the remaining
02 candidate structures that have passed through the decoy filtering and conformer
03 clustering stages. The optimization is carried out by making small perturbations
04 to a candidate structure guided by a highly detailed energy potential. One of
05 the most notable methods is that of Misura et al., which has been shown to be
06 effective in the Sixth Critical Assessment of Techniques for Protein Structure
07 Prediction (CASP-6) (*123,124*). It involves applying perturbations to backbone
08 and side-chain torsion angles using an all-atom force field. The force field
09 consists of a standard 6–12 Lennard–Jones potential for Van der Waals packing,
10 the implicit solvation model of Lazaridis and Karplus describing dielectric
11 screening (*73*), and a new orientation-dependent hydrogen bonding term (*125*).
12 The hydrogen-bonding term is derived based on observed geometrical param-
13 eters of hydrogen bonds in high-resolution crystal structures of proteins. Using
14 this combined physics-based and knowledge-based function as part of their
15 prediction protocol, Bradley et al. have reported success in high-resolution
16 structure prediction of less than 1.5 Å for protein domain of less than 85
17 residues (*124*).

18 A summary of the scoring functions discussed in this chapter can be found
19 in Table 2. We should note that there are other means to guide conforma-
20 tional search and decoy filtering besides using scoring functions. For example,
21 filtering schemes based on contact order (*126*) and beta sheet topology (*127*)
22 have been found to be beneficial in enriching the ensemble quality of decoy
23 structures.

24 25 **3. Discussion and Conclusion** 26

27 A main objective of the structural genomic initiatives, spurred by large-scale
28 genome sequencing efforts, is to determine as many protein folds as possible.
29 The need to determine protein structures rapidly and inexpensively in turn leads
30 to an increased interest in computational protein structure prediction, the two
31 main approaches of which being homology modeling and de novo structure
32 prediction.

33 The key components in de novo protein structure prediction are conforma-
34 tional space sampling and decoy selection. Scoring functions are employed in
35 both the conformational sampling stage and the decoy selection stage. In the
36 first stage, a selected combination of scoring functions approximates the energy
37 landscape of the conformational space, and conformational search algorithms
38 generate trajectories leading to the landscape minima, whereas in the second
39 stage, another set of possibly different scoring functions are used as filter to

*Scoring Functions for De Novo Protein Structure Prediction Revisited 269*01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39**Table 2**
A list of the scoring functions discussed in Section 2

Scoring function	Subheading	Usage	Description
Class I force field	2.1.1.	Conformational space search	Physics-based force field modeling bonded and non-bonded interactions among atoms
RAPDF	2.2.1.	Conformational space search/decoy filtering	Knowledge-based potential describing atom–atom distance preferences
IRAPDF	2.2.3.	Conformational space search/decoy filtering	Continuous version of the RAPDF function
Neural network knowledge-based functions	2.3.	Conformational space search/decoy filtering	Incorporation of neural network into the Bayesian probability framework to describe various conformational properties
Atom–atom contact scoring function	2.4.	Conformational space search/decoy filtering	Knowledge-based atom–atom contact preference function taking solvent accessibility into account
SNAPP potential	2.4.	Conformational space search/decoy filtering	A four-body knowledge-based function describing tiling of protein structures with tetrahedra
Four-body contact scoring function	2.4.	Conformational space search/decoy filtering	A four-body knowledge-based function based on Delauney tessellation of side chain
Residue triplet scoring function	2.4.	Conformational space search/decoy filtering	A three-body knowledge-based function based on the radii of curvature formed among triplets of residues

(Continued)

01 **Table 2**
02 **(Continued)**

03 Scoring function	04 Subheading	05 Usage	06 Description
07 Alpha contact potential	08 2.4.	09 Conformational space search/decoy filtering	10 A two-body knowledge-based function based on edge simplices from the alpha shape of the protein structure
11 Structure refinement potential of Misura et al.	12 2.6.	13 High-resolution refinement	14 A combined physics- and knowledge-based function modeling Van der Waals interaction, solvent effects, and hydrogen bonding

15 Each row gives the name of the scoring function, the subheading in which it is discussed, its
16 usage, and a brief description of its components.
17
18
19

20 retain a collection of the native-like structures. Conformer clustering and high-
21 resolution refinement can also be used as additional steps to further refine this
22 collection. In this chapter, we have studied some examples of the physics-
23 based and knowledge-based scoring functions. For the physics-based approach,
24 the Class I force field and its extensions as well as solvation modeling were
25 discussed. For the knowledge-based approach, we studied the Bayesian proba-
26 bility formalism and used it to derive the RAPDF (34). In addition, we detailed
27 the construction of the neural network-based Bayesian scoring functions. The
28 Bayesian probability formalism was combined with the neural network method-
29 ology to construct various types of log-likelihood scoring functions. Then,
30 we described some of the new knowledge-based scoring functions from in
31 the recent literature. These functions extend the pairwise distance preference
32 scoring functions in various ways, for example, by explicitly modeling the
33 solvent effects and by considering multi-body geometric arrangements and
34 interactions. Finally, we briefly discussed conformer clustering and described
35 a detailed energy potential used for high-resolution refinement. In general,
36 because of the weaknesses of solvent and electrostatic modeling, simulations
37 attempting to fold proteins de novo from physics-based scoring functions alone
38 do not perform satisfactorily. The statistical models that are used to construct
39 knowledge-based functions provide added flexibilities over direct physical

Scoring Functions for De Novo Protein Structure Prediction Revisited 271

01 modeling, and as a result, most of the successful de novo structure prediction
02 protocols have both physics-based and knowledge-based components.

03 Scoring function design remains a very difficult problem. None of the
04 existing physics-based and knowledge-based functions can faithfully reproduce
05 the true energy landscape of the protein conformational space, and none of
06 them can consistently and reliably select native-like conformations from non-
07 native structures for a broad spectrum of proteins. The difficulty is mainly
08 because the physical and statistical models considered so far in the literature
09 cannot well approximate the quantum mechanical character of intra-molecular
10 and solvent-protein interactions. Furthermore, scoring functions describing
11 truly orthogonal characteristics of protein native conformations are difficult
12 to discover, especially for the knowledge-based functions that are the sum of
13 many constituent effects. Thus, it is of practical interest to continue devel-
14 oping various types of new scoring functions, to exploit their differences, and
15 to capture the cumulative effect of incremental enrichments. Fortunately, the
16 increase in the size of the PDB together with increased computational power
17 means that the construction of more sophisticated knowledge-based scoring
18 functions are now possible. More realistic electrostatics and solvation models
19 are also being developed, increasing the capabilities of the physics-based force
20 fields. These advances will play important roles to improving the state of the
21 art of protein folding simulation and de novo structure prediction.

22 Acknowledgments

23 We thank Drs. Enoch Huang and Britt Park for their earlier edition on
24 scoring functions for de novo protein structure prediction and the anonymous
25 reviewer for the many helpful suggestions. This work is supported in part by a
26 Searle Scholar Award, NSF Grant DBI-0217241, an NSF CAREER award, and
27 NIH Grant GM068152 to R.S. and the University of Washington's Advanced
28 Technology Initiative in Infectious Diseases.

30 References

- 31 1. Brenner, S., Levitt, M. (2000) Expectations from structural genomics. *Protein*
32 *Sci.*, **9**, 197–200.
- 33 2. Brenner, S.E. (2001) A tour of structural genomics. *Nat. Genet.*, **210**, 801–809.
- 34 3. Burley, S.K. (2000) An overview of structural genomics. *Nat. Struct. Biol.*, **7**
35 **(Suppl)**, 932–934.
- 36 4. Heinemann, U., Illing, G., Oschkinat, H. (2001) High-throughput three-
37 dimensional protein structure determination. *Curr. Opin. Biotech.*, **12**, 348–354.
- 38 5. Bonneau, R., Baker, D. (2001) Ab initio protein structure prediction: progress
39 and prospects. *Annu. Rev. Biophys. Biomol. Struct.*, **30**, 173–189.

- 01 6. Anfinsen, C.B., Haber, E., Sela, M., White, F.H., Jr. (1961) The kinetics of
02 formation of active ribonuclease during oxidation of the reduced polypeptide
03 chain. *Proc. Natl. Acad. Sci. U. S. A.*, **47**, 1309–1314.
- 04 7. Doolittle, R. (1981) Similar amino acid sequences: chance or common ancestry?
05 *Science*, **214**, 149–159.
- 06 8. Sander, C., Schneider, R. (1991) Database of homology-derived protein structures
07 and the structural meaning of sequence alignment. *Proteins*, **9**, 56–68.
- 08 9. Murzin, A., Bateman, A. (1997) Distance homology recognition using structural
09 classification of proteins. *Proteins*, **29S**, 105–112.
- 10 10. Bowie, J., Luthy, R., Eisenberg, D. (1991) Method to identify protein sequences
11 that fold into a known three-dimensional structure. *Science*, **253**, 164–170.
- 12 11. Jones, D., Taylor, W., Thornton, J. (1992) A new approach to protein fold
13 recognition. *Nature*, **258**, 86–89.
- 14 12. Moult, J., Fidelis, K., Zemla, A., Hubbard, T. (2003) Critical assessment of
15 methods of protein structure prediction (CASP): round V. *Proteins*, **53**, 334–339.
- 16 13. Moult, J., Fidelis, K., Rost, B., Hubbard, T., Tramontano, A. (2005) Critical
17 assessment of methods of protein structure prediction (CASP) – round 6. *Proteins*,
18 **61**, 3–7.
- 19 14. Lee, J., Liwo, A., Ripoll, D., Pillardy, J., Scheraga, J. (1999) Calculation of protein
20 conformation by global optimization of a potential energy function. *Proteins*, **S3**,
21 204–208.
- 22 15. Samudrala, R., Xia, Y., Huang, E., Levitt, M. (1999) Ab initio protein structure
23 prediction using a combined hierarchical approach. *Proteins*, **S3**, 194–198.
- 24 16. Simons, K., Bonneau, R., Ruczinski, I., Baker, D. (1999) Ab initio structure
25 prediction of CASP3 targets using ROSETTA. *Proteins*, **S3**, 171–176.
- 26 17. Samudrala, R., Xia, Y., Levitt, M., Huang E.S. (1999) A combined approach for
27 ab initio construction of low resolution protein tertiary structures from sequence,
28 in *Proceedings of the Pacific Symposium on Biocomputing* (Altman, R. B.,
29 Dunker, A.K., Hunter, L., Klein, T.E., Lauderdale, K., eds.), World Scientific
30 Press, Singapore, pp. 505–516.
- 31 18. Samudrala, R., Levitt, M. (2002) A comprehensive analysis of 40 blind protein
32 structure predictions. *BMC Struct Biol*, **2**, 3–18.
- 33 19. Moult, J., Hubbard, T., Bryant, S.H., Fidelis, K., Pedersen, J.T. (1997) Critical
34 assessment of methods of protein structure prediction (CASP): round II. *Proteins*,
35 **29**, 2–6.
- 36 20. Moult, J., Hubbard, T., Fidelis, K., Pedersen, J.T. (1999) Critical assessment of
37 methods of protein structure prediction (CASP): round III. *Proteins*, **37**, 2–6.
- 38 21. Moult, J., Fidelis, K., Zemla, A., Hubbard, T. (2001) Critical assessment of
39 methods of protein structure prediction (CASP): round IV. *Proteins*, **45**, 2–7.
22. Brooks, B., Bruccoleri, R., Olafson, B., States, D., Swaminathan, S., Karplus, M.
(1983) CHARMM: a program for macromolecular energy, minimization, and
dynamics calculations. *J. Comp. Chem.*, **4**, 187–217.

Scoring Functions for De Novo Protein Structure Prediction Revisited 273

- 01 23. Weiner, S., Kollman P., Nguyen, D., Case, D. (1986) An all atom force field for
02 simulations of proteins and nucleic acids. *J. Comp. Chem.*, **7**, 230–252.
- 03 24. Jorgensen, W., Tirado-Rives, J. (1988) The OPLS potential function for proteins.
04 Energy minimisations for crystals of cyclic peptides and crambin. *J. Amer. Chem.*
05 *Soc.*, **110**, 1657–1666.
- 06 25. MacKerell, A.D., Jr., Bashford, D., Bellott, M., Dunbrack, R.L., Jr.,
07 Evanseck, J.D., et al. (1998) All-atom empirical potential for molecular modeling
08 and dynamics studies of proteins. *J. Phys. Chem. B*, **102**, 3586–3616.
- 09 26. Cornell, W.D., Cieplak, P., Bayly, C.I., Gould, I.R., Merz, K.M., Jr.,
10 Fergusson, D.M., Spellmeyer, D.C., Fox, D.C., Caldwell, J.W., Kollman, P.A.
11 (1995) A second generation force field for the simulation of proteins and nucleic
12 acids. *J. Amer. Chem. Soc.*, **117**, 5179–5197.
- 13 27. Nemethy, G., Gibson, K.D., Palmer, K.A., Yoon, C.N., Paterlini, G., Zagari, A.,
14 Rumsey, S., Scheraga, H.A. (1992) Energy parameters in peptides: improved
15 geometrical parameters and non-bonded interactions for use in the ECEPP/3
16 algorithm, with application to proline-containing peptides. *J. Phys. Chem.*, **96**,
17 6472–6484.
- 18 28. Wodak, S., Rooman, M. (1993) Generating and testing protein folds. *Curr. Opin.*
19 *Struct. Biol.*, **3**, 247–259.
- 20 29. Sippl, M. (1995) Knowledge based potentials for proteins. *Curr. Opin. Struct.*
21 *Biol.*, **5**, 229–235.
- 22 30. Gilis, D., Rooman, M. (1996) Stability changes upon mutation of solvent-
23 accessible residues in proteins evaluated by database-derived potentials. *J. Mol.*
24 *Biol.*, **257**, 1112–1126.
- 25 31. Jernigan, R.L., Bahar I. (1996) Structure-derived potentials and protein simula-
26 tions. *Curr. Opin. Struct. Biol.*, **6**, 195–209.
- 27 32. DeBolt, S.E., Skolnick, J. (1996) Evaluation of atomic level mean force potentials
28 via inverse refinement of protein structures: atomic burial position and pairwise
29 non-bonded interactions. *Protein Eng.*, **8**, 637–655.
- 30 33. Zhang, C., Vasmatazis, G., Cornette, J.L., DeLisi, C. (1997) Determination of
31 atomic desolvation energies from the structures of crystallised proteins. *J. Mol.*
32 *Biol.*, **267**, 707–726.
- 33 34. Samudrala, R., Moult, J. (1998) An all-atom distance-dependent conditional
34 probability discriminatory function for protein structure prediction. *J. Mol. Biol.*,
35 **275**, 895–916.
- 36 35. Huang, E.S., Samudrala, R., Park, B.H. (2000) Scoring functions for ab initio
37 protein structure prediction. *Methods Mol. Biol.*, **143**, 223–245.
- 38 36. Hartree, D.R. (1957) *The Calculation of Atomic Structure*. John Wiley & Sons,
39 New York.
- 37 37. Hohenberg, P., Kohn, W. (1964) Inhomogeneous electron gas. *Phys. Rev.*,
38 **136**, 864.

- 01 38. Kauzmann, W. (1959) Some factors in the interpretation of protein denaturation.
02 *Adv. Protein Chem.*, **14**, 1–64.
- 03 39. Dill, K.A. (1990) Dominant forces in protein folding. *Biochemistry*, **29**,
04 7133–7155.
- 05 40. Morozov, A.V., Kortemme, T., Tsemekhman, K., Baker, D. (2004) Close
06 agreement between the orientation dependence of hydrogen bonds observed in
07 protein structures and quantum mechanical calculations. *Proc. Natl. Acad. Sci.*
08 *U. S. A.*, **101**, 6946–6951.
- 09 41. Weiner, P.K., Kollman P.A. (1981) AMBER: Assisted model building with
10 energy refinement. A general program for modeling molecules and their interac-
11 tions. *J. Comp. Chem.*, **2**, 287–303.
- 12 42. Brooks, B.R., Bruccoleri, R.E., Olafson, B.D., States, D.J., Swaminathan, S.,
13 Karplus, M. (1983) CHARMM: a program for macromolecular energy,
14 minimization, and dynamics calculations. *J. Comp. Chem.*, **4**, 187–217.
- 15 43. Levitt, M., Hirshberg, M., Sharon, R., Daggett, V. (1995) Potential energy
16 function and parameters for simulations of the molecular dynamics of proteins
17 and nucleic acids in solution. *Comp. Phys. Comm.*, **91**, 215–231.
- 18 44. Levitt, M. (1983) Molecular dynamics of native protein. I. Computer simulation
19 of trajectories. *J. Mol. Biol.*, **168**, 595–617.
- 20 45. Daggett, L.P., Saccaan, A.I., Akong, M., Rao, S.P., Hess, S.D., Liaw, C.,
21 Urrutia, A., Jachec, C., Ellis, S.B., Dreessen J, et al. (1995) Molecular
22 and functional characterization of recombinant human metabotropic glutamate
23 receptor subtype 5. *Neuropharmacology*, **34**, 7133–7155.
- 24 46. Levitt, M. (1983) Protein folding by restrained energy minimization and
25 molecular dynamics. *J. Mol. Biol.*, **170**, 723–764.
- 26 47. Brunger, A.T., Clore, G.M., Gronenborn, A.M., Karplus, M. (1986) Three-
27 dimensional structure of proteins determined by molecular dynamics with inter-
28 proton distance restraints: application to crambin. *Proc. Natl. Acad. Sci. U. S. A.*,
29 **83**, 3801–3805.
- 30 48. Ferguson, D.M., Kollman, P.A. (1991) Can the Lennard-Jones 6-12 function
31 replace the 10–12 form in molecular mechanics calculations? *J. Comput. Chem.*,
32 **12**, 620–626.
- 33 49. Halgren, T.A. (1992) Representation of van der Waals (vdW) interactions in
34 molecular mechanics force fields: potential form, combination rules, and vdW
35 parameters. *J. Am. Chem. Soc.*, **114**, 7827–7843.
- 36 50. Halgren, T.A. (1996) Merck molecular force field. I. Basis, form, scope, param-
37 eterization, and performance of MMFF94. *J. Comput. Chem.*, **17**, 490–519.
- 38 51. Hart, J.R., Rappe, A.K. (1992) van der Waals functional forms for molecular
39 simulations. *J. Chem. Phys.*, **97**, 1109–1115.
52. Buckingham, A.D., Fowler, P.W. (1985) A model for the geometries of van der
Waals complexes. *Can. J. Chem.*, **63**, 2018.

Scoring Functions for De Novo Protein Structure Prediction Revisited 275

- 01 53. Sokalski, W.A., Shibata, M., Ornstein, R.L., Rein, R. (1993) Point charge
02 representation of multicenter multipole moments in calculation of electrostatic
03 properties. *Theor. Chim. Acta*, **85**, 209–216.
- 04 54. Stone, A.J. (1981) Distributed multipole analysis, or how to describe a molecular
05 charge distribution. *Chem. Phys. Lett.*, **83**, 233–239.
- 06 55. Kosov, D., Popelier, P.L.A. (2000) Atomic partitioning of molecular electrostatic
07 potentials. *J. Phys. Chem. A*, **104**, 7339–7345.
- 08 56. Cieplak, P., Caldwell, J., Kollman, P. (2001) Molecular mechanical models
09 for organic and biological systems going beyond the atom centered two body
10 additive approximation: aqueous solution free energies of methanol and N-methyl
11 acetamide, nucleic acid base, and amide hydrogen bonding and chloroform/water
12 partition coefficients of the nucleic acid bases. *J. Comput. Chem.*, **22**, 1048–1057.
- 13 57. Kaminski, G.A., Stern, H.A., Berne, B.J., Friesner, R.A., Cao, Y.X.,
14 Murphy, R.B., Zhou, R., Halgren, T.A. (2002) Development of a polarizable
15 force field for proteins via ab initio quantum chemistry: first generation model
16 and gas phase tests. *J. Comput. Chem.*, **23**, 1515–1531.
- 17 58. Ren, P., Ponder, J.W. (2003) Polarizable atomic multipole water model for
18 molecular mechanics simulation. *J. Phys. Chem. B*, **107**, 5933–5947.
- 19 59. Jorgensen, W.L. (1981) Transferable intermolecular potential functions for water,
20 alcohols, and ethers. Application to liquid water. *J. Am. Chem. Soc.*, **103**,
21 335–340.
- 22 60. Jorgensen, W.L., Chandrasekhar, J., Madura, J.D., Impey, R.W., Klein, M.L.
23 (1983) Comparison of simple potential functions for simulating liquid water. *J.*
24 *Chem. Phys.*, **79**, 926–935.
- 25 61. Berendsen, H.J.C., Grigera, J.R., Straatsma, T.P. (1987) The missing term in
26 effective pair potentials. *J. Phys. Chem.*, **91**, 6269–6271.
- 27 62. Levitt, M., Hirshberg, M., Sharon, R., Laidig, K.E., Daggett, V. (1997) Calibration
28 and testing of a water model for simulation of the molecular dynamics of proteins
29 and nucleic acids in solution. *J. Phys. Chem. B*, **101**, 5051–5061.
- 30 63. York, D.M., Darden, T., Pedersen, L.G. (1993) The effect of long-range electro-
31 static interactions in simulations of macromolecular crystals: a comparison of the
32 Ewald and truncated list methods. *J. Chem. Phys.*, **99**, 8345–8348.
- 33 64. Darden, T., York, D., Pedersen, L. (1993) Particle mesh Ewald: an $N^2 \log(N)$
34 method for Ewald sums in large systems. *J. Chem. Phys.*, **98**, 10089–10092.
- 35 65. Gouy, M. (1910) Sur la constitution de la charge électrique a la surface d'un
36 électrolyte. *Journ. Phys.*, **9**, 457–468.
- 37 66. Gilson, M.K., Honig, B. (1988) Calculation of the total electrostatic energy of a
38 macromolecular system: solvation energies, binding energies, and conformational
39 analysis. *Proteins*, **4**, 7–18.
67. Nicholls, A., Honig, B. (1991) A rapid finite difference algorithm, utilizing
successive over-relaxation to solve the Poisson-Boltzmann equation. *J. Comp.*
Chem., **12**, 435–445.

- 01 68. Bashford, D., Case, D.A. (2000) Generalized Born models of macromolecular
02 solvation effects. *Annu. Rev. Phys. Chem.*, **51**, 129–152.
- 03 69. de Bakker, P.I.W., DePristo, M.A., Burke, D.F., Blundell, T.L. (2003) Ab initio
04 construction of polypeptide fragments: accuracy of loop decoy discrimination by
05 an all-atom statistical potential and the AMBER force field with the generalized
06 born solvation model. *Proteins*, **51**, 21–40.
- 07 70. Fogolari, F., Brigo, A., Molinari, H. (2003) Protocol for MM/PBSA molecular
08 dynamics simulations of proteins. *Biophys. J.*, **85**, 159–166.
- 09 71. Warshel, A., Levitt, M. (1976) Theoretical studies of enzymic reactions –
10 dielectric, electrostatic and steric stabilization of carbonium-ion in reaction of
11 lysozyme. *J. Mol. Biol.*, **103**, 227–249.
- 12 72. Gelin, B.R., Karplus, M. (1979) Side-chain torsional potentials: effect of
13 dipeptide, protein, and solvent environment. *Biochemistry*, **18**, 1256–1268.
- 14 73. Lazaridis, T., Karplus, M. (1999) Effective energy function for proteins in
15 solution. *Proteins*, **35**, 133–152.
- 16 74. Mallik, B., Masunov, A., Lazaridis, T. (2002) Distance and exposure dependent
17 effective dielectric function. *J. Comp. Chem.*, **23**, 1090–1099.
- 18 75. Moulton, J. (1997) Comparison of database potentials and molecular mechanics
19 force fields. *Curr. Opin. Struct. Biol.*, **7**, 194–199.
- 20 76. Eisenberg, D., Weiss, R.M., Terwillinger, T.C. (1982) The helical hydrophobic
21 moment: a measure of the amphiphilicity of a helix. *Nature*, **299**, 371–374.
- 22 77. Sippl, M.W., S. (1992) Detection of native-like models for amino acid sequences
23 of unknown three-dimensional structure in a database of known protein conform-
24 ations. *Proteins*, **13**, 258–271.
- 25 78. Jones, D.T. (2001) Predicting novel protein folds by using FRAGFOLD. *Proteins*,
26 **45**, 127–132.
- 27 79. Zhang, Y., Skolnick, J. (2004) Tertiary structure predictions on a comprehensive
28 benchmark of medium to large size proteins. *Biophys. J.*, **87**, 2647–2655.
- 29 80. Boniecki, M., Rotkiewicz, P., Skolnick, J., Kolinski, A. (2003) Protein fragment
30 reconstruction using various modeling techniques. *J. Comput. Aided Mol. Des.*,
31 **17**, 725–738.
- 32 81. Hung, L.H., Ngan, S.C., Liu, T., Samudrala, R. (2005) PROTINFO: new
33 algorithms for enhanced protein structure predictions. *Nucleic Acids Res.*, **33**,
34 W77–W80.
- 35 82. Westbrook, J., Feng, Z., Chen, L., Yang, H., Berman, H.M. (2003) The Protein
36 Data Bank and structural genomics. *Nucleic Acids Res.*, **31**, 489–491.
- 37 83. Bourne, P.E., Address, K.J., Bluhm, W.F., Chen, L., Deshpande, N., Feng, Z.,
38 Fleri, W., Green, R., Merino-Ott, J.C., Townsend-Merino, W., Weissig, H.,
39 Westbrook, J., Berman, H.M. (2004) The distribution and query systems of the
RCSB Protein Data Bank. *Nucleic Acids Res.*, **32**, D223–D225.
84. Chandonia, J.M., Hon, G., Walker, N.S., LoConte, L., Koehl, P., Levitt, M.,
Brenner, S.E. (2004) The ASTRAL compendium in 2004. *Nucleic Acids Res.*,
32, D189–D192.

Scoring Functions for De Novo Protein Structure Prediction Revisited 277

- 01 85. Subramaniam, S., Tchong, D.K., Fenton, J. (1996) Knowledge-based methods
02 for protein structure refinement and prediction, in *Proceedings of the Fourth*
03 *International Conference on Intelligent Systems in Molecular Biology* (States, D.,
04 Agarwal, P., Gaasterland, T., Hunter, L. & Simth, R., eds.), AAAI Press, Menlo
05 Park, CA, pp. 218–229.
- 06 86. Avbelj, F., Moult, J. (1995) Role of electrostatic screening in determining protein
07 main chain conformational preferences. *Biochemistry*, **34**, 755–764.
- 08 87. Lu, H., Skolnick, J. (2001) A distance-dependent atomic knowledge-based
09 potential for improved protein structure selection. *Proteins*, **44**, 223–232.
- 10 88. Zhou, H., Zhou, Y. (2002) Distance-scaled, finite ideal-gas reference state
11 improves structure-derived potentials of mean force for structure selection and
12 stability prediction. *Protein Sci.*, **11**, 2714–2726.
- 13 89. Oppenheim, A.V., Schafer, R.W., Buck, J.R. (1999) *Discrete-Time Signal*
14 *Processing*, 2nd ed. Prentice Hall, Upper Saddle River, NJ.
- 15 90. Rost, B., Sander, C. (1994) Conservation and prediction of solvent accessibility
16 in protein families. *Proteins*, **20**, 216–226.
- 17 91. Ahmad, S., Gromiha, M.M. (2002) NETASA: neural network based prediction
18 of solvent accessibility. *Bioinformatics*, **18**, 819–824.
- 19 92. Kim, H., Park, H. (2004) Prediction of protein relative solvent accessibility with
20 support vector machines and long-range interaction 3D local descriptor. *Proteins*,
21 **54**, 557–562.
- 22 93. Rost, B., Sander, C. (1993) Prediction of protein secondary structure at better
23 than 70% accuracy. *J. Mol. Biol.*, **232**, 584–599.
- 24 94. Jones, D.T. (1999) Protein secondary structure prediction based on position-
25 specific scoring matrices. *J. Mol. Biol.*, **292**, 195–202.
- 26 95. Cuff, J.A., Barton, G.J. (1999) Application of enhanced multiple sequence
27 alignment profiles to improve protein secondary structure prediction. *Proteins*,
28 **40**, 502–511.
- 29 96. Lund, O., Frimand, K., Gorodkin, J., Bohr, H., Bohr, J., Hansen, J., Brunak, S.
30 (1997) Protein distance constraints predicted by neural networks and probability
31 density functions. *Protein Eng.*, **10**, 1241–1248.
- 32 97. Pollastri, G., Baldi, P., Fariselli, P., Casadio, R. (2002) Prediction of coordination
33 number and relative solvent accessibility in proteins. *Proteins*, **47**, 142–153.
- 34 98. Olmea, O., Valencia, A. (1997) Improving contact predictions by the combination
35 of correlated mutations and other sources of sequence information. *Fold Des.*, **2**,
36 S25–32.
- 37 99. Fariselli, P., Casadio, R. (1999) Neural network based predictor of residue
38 contacts in proteins. *Protein Eng.*, **12**, 15–21.
- 39 100. Altschul, S.F., Madden, T.L., Schaffer, A.A. (1997) Gapped BLAST and PSI-
BLAST: a new generation of protein database search programs. *Nucleic Acids*
Res., **25**, 3389–3402.

- 01 101. Rumelhart, D.E., Hinton, G.E., Williams, R.J. (1986) Learning representations
02 by back-propagating errors. *Nature*, **323**, 533–536.
- 03 102. Yi, T.-M., Lander, E.S. (1993) Protein secondary structure prediction using
04 nearest-neighbor methods. *J. Mol. Biol.*, **232**, 1117–1129.
- 05 103. Kabsch, W., Sander, C. (1983) Dictionary of protein secondary structure: pattern
06 recognition of hydrogen-bonded and geometrical features. *Biopolymers*, **22**,
07 2577–2637.
- 08 104. Zell, A., Mamier, G., Vogt, M., et al. (2005) The SNNS users manual version
09 4.1. Available at <http://www-ra.informatik.uni-tuebingen.de/snns>.
- 10 105. Park, B., Levitt, M. (1996) Energy functions that discriminate x-ray and near
11 native folds from well-constructed decoys. *J. Mol. Biol.*, **266**, 831–846.
- 12 106. Novotny, J., Brucoleri, R., Karplus, M. (1984) An analysis of incorrectly
13 folded protein models. Implications for structure predictions. *J. Mol. Biol.*, **177**,
14 787–818.
- 15 107. Holm, L., Sander, C. (1992) Evaluation of protein models by atomic solvation
16 preference. *J. Mol. Biol.*, **225**, 93–105.
- 17 108. Samudrala, R., Levitt, M. (2000) Decoys ‘R’ Us: a database of incorrect conformations
18 to improve protein structure prediction. *Protein Sci.*, **9**, 1399–1401.
- 19 109. Tsai J., B., R., Morozov, A.V., Kuhlman, B., Rohl, C.A., Baker, D. (2003) An
20 improved protein decoy set for testing energy functions for protein structure
21 prediction. *Proteins*, **53**, 76–87.
- 22 110. Park, B.H., Huang, E.S., Levitt, M. (1997) Factors affecting the ability of energy
23 functions to discriminate correct from incorrect folds. *J. Mol. Biol.*, **266**, 831–846.
- 24 111. Hinds, D.A., Levitt, M. (1992) A lattice model for protein structure prediction at
25 low resolution. *Proc. Natl. Acad. Sci. U. S. A.*, **89**, 2536–2540.
- 26 112. Park, B., Levitt, M. (1995) The complexity and accuracy of discrete state models
27 of protein structure. *J. Mol. Biol.*, **249**, 493–507.
- 28 113. Simons, K.T., Kooperberg, C., Huang, E., Baker, D. (1997) Assembly of protein
29 tertiary structures from fragments with similar local sequences using simulated
30 annealing and Bayesian scoring functions. *J. Mol. Biol.*, **268**, 209–225.
- 31 114. Hung, L.H., Samudrala, R. (2003) PROTINFO: secondary and tertiary protein
32 structure prediction. *Nucleic Acids Res.*, **31**, 3296–3299.
- 33 115. McConkey, B.J., Sobolev, V., Edelman, M. (2003) Discrimination of native
34 protein structures using atom-atom contact scoring. *Proc. Natl. Acad. Sci. U. S. A.*,
35 **100**, 3215–3220.
- 36 116. Carter, C.W., Jr., LeFebvre, B.C., Cammer, S.A., Tropsha, A., Edgell, M.H.
37 (2001) Four-body potentials reveal protein-specific correlations to stability
38 changes caused by hydrophobic core mutations. *J. Mol. Biol.*, **311**, 625–638.
- 39 117. Krishnamoorthy, B., Tropsha, A. (2003) Development of a four-body statistical
pseudo-potential to discriminate native from non-native protein conformations.
Bioinformatics, **19**, 1540–1548.

Scoring Functions for De Novo Protein Structure Prediction Revisited 279

- 01 118. Ngan, S.-C., Inonye, M.T, Samudrala, R. (2006) A knowledge-based scoring
02 function based on residue triplets for protein structure prediction. *Protein Eng.*,
03 **19**, 187–193.
- 04 119. Li, X., Hu, C., Liang, J. (2003) Simplicial edge representation of protein structures
05 and alpha contact potential with confidence measure. *Proteins*, **53**, 792–805.
- 06 120. Wang, K., Fain, B., Levitt, M., Samudrala, R. (2004) Improved protein structure
07 selection using decoy-dependent discriminatory functions. *BMC Struct. Biol.*, **4**, 8.
- 08 121. Zhang, Y., Skolnick, J. (2004) SPICKER: a clustering approach to identify near-
09 native protein folds. *J. Comput. Chem.*, **25**, 865–871.
- 10 122. Samudrala, R. (2006). RAMP Howto. Available at <http://software.compbio.washington.edu/ramp/ramp.html>
- 11 123. Misura, K.M.S., Baker, D. (2005) Progress and challenges in high-resolution
12 refinement of protein structure models. *Proteins*, **59**, 15–29.
- 13 124. Bradley, P., Misura, K.M.S., Baker, D. (2005) Toward high-resolution de novo
14 structure prediction for small proteins. *Science*, **309**, 1868–1871.
- 15 125. Kortemme, T., Morozov, A.V., Baker, D. (2003) An orientation-dependent
16 hydrogen bonding potential improves prediction of specificity and structure for
17 proteins and protein-protein complexes. *J. Mol. Biol.*, **326**, 1239–1259.
- 18 126. Bonneau, R., Ruczinski, I., Tsai, J., Baker, D. (2002) Contact order and ab initio
19 protein structure prediction. *Protein Sci.*, **11**, 1937–1944.
- 20 127. Bradley, P., Malmstrom, L., Qian, B., Schonburn, J., Chivian, D., Kim, D.E.,
21 Meiler, J., Misura, K.M., Baker D. (2005) Free modeling with Rosetta in CASP6.
22 *Proteins*, **61**, 128–134.
- 23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

UNCORRECTED PROOF

01 **QUERIES TO BE ANSWERED (SEE MARGINAL MARKS)**
02 **IMPORTANT NOTE: Please mark your corrections and answers to these**
03 **queries directly onto the proof at the relevant place. Do NOT mark your**
04 **corrections on this query sheet.**
05

06 **Chapter-10**

07 Query No.	Page No.	Line No.	Query
09 AQ1	262	34	The text 'Based on the work of Professor Banavar and his colleagues' has been changed to 'On the basis of the work of Professor Banavar and his colleagues'. Please check if this is OK.

13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39

UNCORRECTED