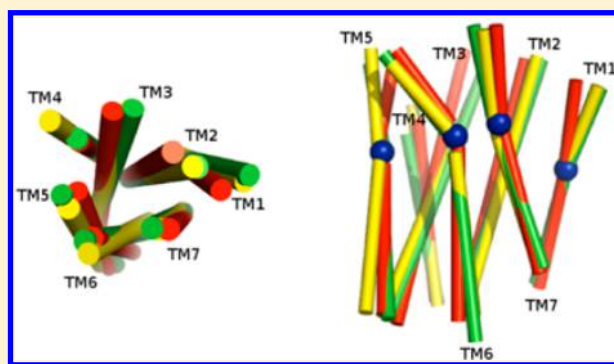# Kink Characterization and Modeling in Transmembrane Protein Structures

Tim Werner and W. Bret Church*

Group in Biomolecular Structure and Informatics, Faculty of Pharmacy, The University of Sydney, Sydney NSW 2006, Australia

Ⓢ Supporting Information

**ABSTRACT:** Kinks have been observed to provide important functional and structural features for membrane proteins. Despite their ubiquity in membrane proteins, and their perceived importance, no protein modeling methods explicitly considers kinks. In spite of the limited data for transmembrane proteins, we were able to develop a knowledge-based modeling method for introducing kinks, which we demonstrate can be exploited in modeling approaches to improve the quality of models. The work entailed a thorough analysis of the available high resolution membrane protein structures, concomitantly demonstrating the complexity of the structural considerations for kink prediction. Furthermore, our results indicate that there are systematic and significant differences in the sequence as well as the structural environment between kinked and nonkinked transmembrane helices. To the best of our knowledge, we are reporting a method for modeling kinks for the first time.

## INTRODUCTION

Proteome analysis of several organisms reveals that the proteomes are constituted by around 20−25% membrane proteins.[1,2] This class of protein is also a major target for the pharmaceutical industry. Integral membrane proteins are generally characterized by a hydrophobic domain which serves to retain the protein in the membrane, and although it can be relatively small, it can define the function and structure of the protein. Because of the lipophilic environment of the membrane surrounding integral membrane proteins, the sequences and structures differ significantly to soluble proteins. One of the most important subtypes of membrane proteins are those which mediate between two cellular compartments, performing a signaling role, and hence transmitting information across the membrane rather than physically transporting cargo.
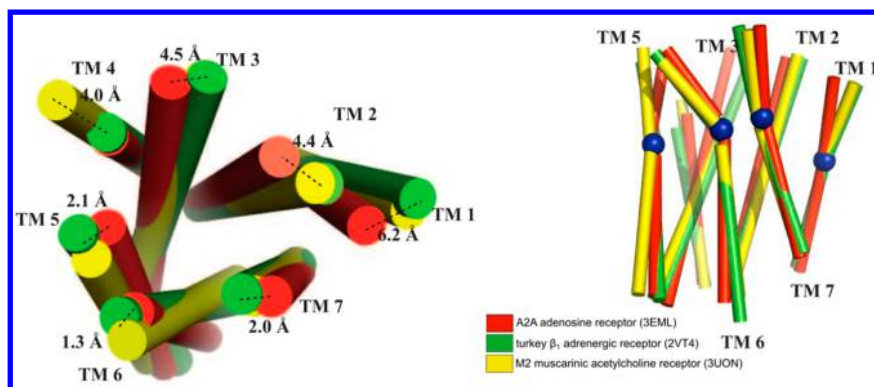
The most routinely used method for modeling three-dimensional structures due to its accuracy is homology modeling,[3] which uses known structures as a template to inform the unknown structure (target). As such, homology modeling relies on experimentally determined protein structures, and those with a sequence identity high enough to the target protein can achieve high accuracy models. The greater the similarity in target and experimental structure as measured by sequence comparisons the greater the structural similarity. Chothia and Lesk[4] provided the first quantification of the general extent to which structural changes were related to the relative magnitude of the sequence changes. These observations and the original development of the homology modeling method were prior to any appreciable progress in the structural determinations of integral membrane proteins.

However, homology modeling can be applied to membrane proteins despite the systematic structural differences between membrane proteins and soluble proteins.[5] Also, many homology modeling approaches generate more than one model, and therefore model quality assessment algorithms are employed to identify the high quality models. To identify the high quality models the algorithms for the assessment often rely on knowledge-based methods, which are themselves derived from the databases almost exclusively constituted by water-soluble proteins and hence such algorithms need to be adapted for membrane protein use.[6,7] To our best knowledge, despite the need of model quality assessment algorithms optimized for membrane proteins, there are only two methods developed for membrane proteins.[8,9] Apart from the preponderance of hydrophobic residues in integral membrane protein domains, which might be anticipated and which reduces the threshold of sequence similarities for structural similarity seen in soluble proteins, the two-dimensionality of the membrane, although fluid, imposes different protein architectural constraints.

One of the widely observed features in transmembrane proteins related to the overall architecture are kinks, which are not in such evidence in soluble proteins, but that play an important role for the structure and the function of transmembrane proteins. Kinks have been observed to provide sites of flexibility,[10] form cavities essential for ligand binding[11] or the active site of catalysis,[12] and can be wedged by water which provides hydrophilic contacts within the hydrophobic

**Figure 1.** Comparison of human A2A adenosine receptor (red) with turkey $\beta_2$-adrenergic receptor (2VT4) (green) and M2 muscarinic acetylcholine receptor (3UON) (yellow). The structural differences within the transmembrane regions are caused by kinks (blue balls).

**Table 1. Average Kink Size of the 7 TM Helices in All 14 GPCR Structures**

|         | TM1    | TM2   | TM3    | TM4    | TM5    | TM6    | TM7   |
|---------|--------|-------|--------|--------|--------|--------|-------|
| average | 11.5°  | 34°   | 14.07° | 31.36° | 41.5°  | 28.43° | 23.5° |
| std dev | 3.88°  | 9.21° | 6.49°  | 9.60°  | 16.80° | 3.34°  | 6.57° |

environment of the membrane.[13] Despite their ubiquity in membrane proteins, and their perceived importance, none of the protein modeling methods explicitly considers kinks. In the case of homology modeling, kinks are modeled by the implicit adoption of the kink position and the kink size dictated by the templates. Studies of primary sequence at kinks have been performed as a way of understanding their structural origins, but none of them take the explicit structural environment of kinks into account.[14] Methods have been developed to predict the kink position from sequence,[15,16] but they do not predict the kink size. Methods utilizing ligand binding when searching for plausible conformations have been reported and have specifically used standard helix descriptions such as tilt, rotation and translation to seek models (for a review, see ref 17). Such methods have demonstrated particular success in positioning ligands in two GPCRs but are dependent on ligand knowledge as are the other methods that have been reported to have some success in modeling helices in the vicinity of kinks.[18]

G-protein coupled receptors (GPCRs) are one of the most important members of transmembrane proteins accounting for around 30% of all drug targets[19] and hence command much attention. Today, there are 14 distinct GPCR crystal structures available, which are not only valuable in their own right but could be employed to provide models of many more GPCRs by homology modeling. With the membrane domain consisting of only seven transmembrane helices, they are relatively small. However, for the human GPCRs, the sequence identity between target and template is not sufficiently high enough to correctly select the best available template. To illustrate, the human $A_{2A}$ adenosine receptor (PDB ID: 3EML)[20] is an example in which a template with lower sequence identity is more appropriate because it has sequence identities to turkey $\beta_1$ adrenergic receptor (PDB ID: 2VT4)[21] and human $M_2$ muscarinic acetylcholine receptor (PDB ID: 3UON)[22] of 36.5% and 26.8%, respectively, in the transmembrane region. However, the RMSD in the transmembrane region between human A2A adenosine receptor for the comparison with turkey $\beta_1$ adrenergic receptor is higher (2.7 Å) than the RMSD for that with the human $M_2$ muscarinic acetylcholine receptor (2.3 Å). Significantly though too, neither GPCR structure is able to satisfactorily describe the kinks in the human A2A adenosine

receptor adequately (Figure 1). Kinks are seen in most of the transmembrane helices of all the GPCR structures, though the kink size is seen to vary significantly with the highest diversity observed in transmembrane helix 5 (TM5). Whereas the kink size deviation in transmembrane helix 1 (TM1) and transmembrane helix 6 (TM6) is lowest in all structures (Table 1 and Supporting Information Figure S1). Because of the diversity of the kinks even in proteins of relatively high sequence identity, and the difficulty in modeling the kinks based on a fixed template, new modeling methods focusing on kink regions are necessary to achieve high quality models.

The aim of this study was to develop a membrane-protein specific knowledge-based scoring function that predicts kink sizes in structural models. The main intention was to use information from known structures to improve models of unknown structures with the critical use of improved kink modeling. A database was constructed from which the scoring function was derived and was also used to analyze the sequence near the kinks as well as the structural environment of the kinks. We anticipate that our method will assist in the overall structure predictions of transmembrane helical proteins without depending on ligand binding knowledge and additionally can be combined with other methods.

## ■ METHODS

The knowledge-based scoring approach benefits from a comprehensive data set from which to derive the knowledge, to evaluate the scores and to optimize the scoring parameters. There are no standard publicly available nonredundant sets for membrane proteins and the importance to have up to date data sets in the fast growing field of membrane proteins makes it worthwhile to generate new data sets to derive the score, optimize the parameters and testing. In this work, a data set for training and testing as well as a data set based on GPCR structures for an independent test set to quantify the scoring performance was generated. The method for generating the training and test data set for membrane proteins adapted from that of Ray et al.[9] and is described in the following section.

**Training and Test Data.** Here, 358 PDB files containing $\alpha$ helical transmembrane proteins were downloaded from OPM (Orientations of Proteins in Membranes)[23] in 2012. The 358

PDB files contain 538 individual polytopic transmembrane polypeptide structures. Structures with resolution lower than 3.5 Å were removed from the data set and then the remaining structures (304) were clustered by cdhit[24] using a 40% sequence identity threshold. Afterward, a pairwise sequence alignment of all versus all centroids from each cluster was performed and only sequences with 35% identity or less to all other centroids were used for further analyses (102 structures). The 102 crystal structures were used as the training set for deriving the score as well as performing the kink analysis.
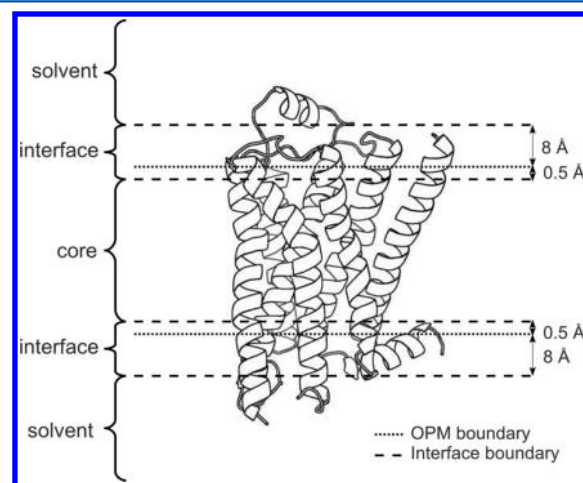
The 102 structures were structurally pairwise aligned to all OPM database structures with the same number of transmembrane helices by TMalign[25] in order to generate a data set for 5-fold cross validation containing the crystal structures as well as decoys. Structures with more than 100 residues in the alignment and with sequence identity between 20% and 90% were used as templates for the homology modeling approach in the following step. For 45 of the 102 structures, templates which fulfilled these criteria were found, and these 45 structures, now serving as targets, were realigned to their templates in multiple alignments using CLUSTALW2.[26] The resulting multiple alignments of the 45 structures were then used to generate 10 homology models for each target-template pair in the alignments using MODELLER9 V8[27] and the resultant models in turn were clustered with MaxCluster[28] based on RMSD so that for every target 10 clusters were created. The representative centroid structure for each cluster was determined, as the one whose sum of RMSD to all other structures of the cluster was the smallest. These centroids were then chosen as decoys which resulted overall in a total of 450 decoys. Decoy subsets are then chosen by selecting randomly the target and adding the corresponding decoys to the subset so that 5 different subsets, each contains 90 decoys, were created. The corresponding crystal structures (targets) were added to each of the decoy subsets. The sets were used for the 5-fold cross validation procedure which was used in the parameter optimization and the weighting of each scoring component. The 57 structures for which there was no suitable template were used for scoring in all 5-fold cross validation data sets. The average RMSD in the test sets is 3.89 Å [std dev 2.31 Å], the highest RMSD is 10.32 Å, and the lowest is 0.3 Å.

**Target Function.** We have used the z-score to measure the validity of the scoring function. The z-score[29,30] is used to describe how well the resulting scores differentiate the native fold of a protein from an ensemble of misfolded structures: $Z_{score} = (\langle E_{misfolds} \rangle - E_{native})/\sigma_{misfolds}$ where $E_{native}$ is the score of the native structure of a protein, $\langle E_{misfolds} \rangle$ is the average score of the test subset of misfolded structures and $\sigma_{misfolds}$ is the standard deviation of the score in the test subset. The average $\langle z\text{-score} \rangle$ obtained from the training and test set over all 5 subsets was used to measure the overall performance of the scoring function and to optimize the parameters for the scoring components. In addition to the z-score, the rank and the Pearson's correlation coefficient were used to quantify the performance of the scoring function within the independent test set. However, these methods were not used for selecting the parameters and weighting the scores.

**Scoring Function.** The scoring function for the kink modeling consists of four scoring components: residue–residue contact score, atom–atom contact score, membrane burial propensity score, and SASA score. The parameters for each scoring component were individually optimized, whereas the

weights for the individual component scores were optimized together.

**Residue–Residue Contacts.** This describes the distribution of residue–residue contacts. The residues of each structure were classified as originating from one of either a core, interface or solvent region. Beginning from the OPM boundary definition, the interface region was defined by extending both OPM boundaries toward the extracellular and cytoplasmic side. All possible combinations of the two between the maximum +10 Å and minimum −10 Å with a step size of 0.5 Å from the OPM boundaries were tested using the 5-fold cross validation data set. The interface region with the highest $\langle z\text{-score} \rangle$ (3.48; std dev 0.4) was obtained by extending the OPM boundaries 8 Å toward the solvent and 0.5 Å toward the center of the membrane. The size of the interface defines also the size of the core and solvent region (Figure 2), which vary for each protein, depending on the OPM boundary definition.



**Figure 2.** Schematic overview of the three regions (solvent, interface, core) each structure was divided in for the residue–residue and atom–atom contact score. The OPM boundary definition was extended toward the extracellular (8 Å) and cytoplasmic (0.5 Å) side forming the interface region. On the basis of the interface region, the solvent and membrane were defined as solvent region and core region, respectively. The thickness of the core varies for each protein, depending on the OPM boundary definition.

In each region, the observed number of contacts between two residues was counted to derive the probability (observed probability) that two residue types are in contact. Two residues are considered to be in contact if at least one atom of each residue is in contact. Two atoms are in contact if the distance between the two atoms is smaller than or equal to the sum of both van der Waals (vdW) radii[31] plus a threshold of 0.6 Å. All hydrophobic and acceptor–donor contacts were considered and the remaining contacts were ignored. The definition of acceptor and donor contacts is described in Supporting Information Table S1. Contacts between atoms which do not appear in the table were considered to be hydrophobic contacts.

For each region the fraction of residue type $r$ and of type $r'$ is derived from the database and therefore the expected probability to observe randomly a contact of type $rr'$ is

$$P(rr'|region_1 region_2)_{exp} = P(r|region_1)P(r'|region_2)$$

where region$_1$ and region$_2$ are core, interface, or solvent, $P(rr'|$ region$_1$region$_2$)$_\text{exp}$ is the probability of a random contact of type $rr'$ of residue type $r$ in region$_1$ and of residue type $r'$ in region$_2$. $P(r|$region$_1$) and $P(r'|$region$_2$) are the fractions of residue type $r$ and type $r'$ in the corresponding regions. The observed and the expected probability of a contact of type $rr'$ can be combined into a log odds score over all residue−residue contact pairs $ij$: $S_\text{residue\_residue} = \sum_i\sum_{j>i} - \log(P(r_ir'_j|\text{region})_\text{obs}/P(r_ir'_j|\text{region})_\text{exp})$ where $r_i$ is the residue type of residue $i$ and $r'_j$ is the residue type of residue $j$. The magnitude of the log odd score $-\log(P(r_ir'_j|\text{region})_\text{obs}/P(r_ir'_j|\text{region})_\text{exp})$ gives a measure of how "nonrandomly" the pair $r_ir'_j$ occurs and it is positive if $r_ir'_j$ is observed less often than expected and negative if $r_ir'_j$ is observed more often. Two scores were created: one for long-range interactions (contacts between residues five or more residues away in sequence) and one for short-range interactions (contacts between residues less than five residues away in the sequence). To overcome the problem of undersampling, the score is set to 0 for a specific residue−residue contact if the number of observed residue−residue contacts is less than 2.

**Atom−Atom Contacts.** The score was calculated similarly to the residue−residue contact score. Each residue was categorized into one of core, interface, or solvent region according to the residue−residue contact score (⟨$z$-score⟩ 1.81; std dev 0.2). The highest ⟨$z$-score⟩ was observed by extending the OPM boundaries 8.5 Å toward the solvent and 0.5 Å toward the center of the membrane (⟨$z$-score⟩ 1.93; std dev 0.3) but because the ⟨$z$-score⟩ decreased by only 0.12 when using the same definition as for the residue−residue score, we decided to use these definitions for both the atom−atom and residue−residue scores to achieve a consistent region size. Atoms were grouped into 13 different atom types based on chemical properties.[32] The definition used for two atoms to be into contact was the same as used for the residue−residue contacts (see residue−residue contacts section). Only contacts between atoms more than five residues away in sequence were considered. The log-odds score ($S_\text{atom\_atom}$) over all atom−atom contact pairs was calculated analogously to the residue−residue contact score (see section above) except using the 13 atom types instead of the residues.

**Membrane Burial Propensity.** The membrane burial propensity component describes the propensity of a residue type to be buried at a particular depth within the membrane. The value of the score changes along the membrane normal, $z$, and the score is zero outside the membrane domain. Several sizes for the membrane domain thickness were tested but the highest ⟨$z$-score⟩ (4.6; std dev 1.7) was achieved defining membrane domain as a region spanning from 7 Å ($z = 0$) in the direction of the solvent starting from the OPM boundary definition, whereas the distance to the membrane center was set fixed to 15 Å ($z = 22$) (see Supporting Information Figure S2). The score as implemented does not distinguish between the sides of the membrane, but rather considers the observation that certain residues have preferences to be buried in the membrane in a certain depth, for example tryptophan can be often found in the interface region.[33] The scoring function is based on Ulmschneider et al.[34] with a difference that the score is normalized by the residue distribution in the solvent region: $S(r, z) = -\log(P(r|z)/P(r|z = 0))$, where $z$ can obtain a value between 0 and 22; 0 indicates that the residue is outside the membrane domain which will give this residue a membrane burial score of 0 and will therefore not be considered for the scoring. $P(r|z)$ is the probability of observing the residue type $r$

in depth $z$, and $P(r|z = 0)$ is the probability of observing the same residue type in the solvent. The score was fitted with a smoothed function due to the variations anticipated in a small number of available structures. The density function used is the Gaussian, $D_r(z) = \sum_i^n G_{z_i}^\sigma(z)$, where $n$ is the frequency of residue type $r$ observed in the particular depth $z$. The index $i$ indicates the $i$th observation, and $z_i$ is the burial depth of the $i$th observation; $\sigma$ is the Gaussian smoothing parameter, which was optimized using the 5-fold cross validation data set and is set to 2. Using the density function, the score over all residues $i$ is the following: $S_\text{membrane\_burial} = \sum_i - \log\{[D_{r_i}(z_i)/\sum_a D_a(z_i)]/[D_{r_i}(z = 0)/\sum_a D_a(z = 0)]\}$ where the summations are over all residues $i$ and the 20 amino acid types denoted by $a$, $r_i$ is the residue type of residue $i$, $D_{r_i}(z = 0)$ is the density at $z = 0$, and $D_{r_i}(z_i)$ as well as $D_a(z_i)$ at the given $z$ value.

**Solvent Accessible Surface Area.** The solvent accessible surface area (SASA) distribution is similar to the membrane burial propensity score with the exception that each residue was grouped as either buried or exposed. The solvent accessible surface was calculated using the relative exposure of the side chains calculated by NACCESS.[35] Different values for grouping residues into buried and exposed were tried but a threshold of 35% was found to achieve the highest ⟨$z$-score⟩ (2.9; std dev 0.6), which means that residues with a side chain exposure of 35% or more were categorized as exposed. There are 20 different residues types and each of these residue types can be categorized as buried or exposed resulting into 40 different residue grouping pair types (see Supporting Information Figure S3). The score was summarized over all residues and their corresponding residue grouping pair $i$: $S_\text{SASA} = \sum_i - \log\{[D_{r_i}(z_i)/\sum_a D_a(z_i)]/[D_{r_i}(z = 0)/\sum_a D_i(z = 0)]\}$ where the summation in the log odd part is over all 40 residue grouping pair types denoted by $a$, $z$ is the burial depth, and $r_i$ is the residue type of the residue grouping pair $i$.

**Weight Optimization.** The weight of each scoring component was optimized by using the method of steepest descent and the 5-fold cross validation data set. The ⟨$z$-score⟩ was used as the target function. In 10 000 steps, the weights were randomly chosen between 0 and 10 and each weight was then optimized by the steepest descent method changing the weights in each step by 0.01 until a maximum number of steps (1000) are reached or the ⟨$z$-score⟩ difference between two steps is smaller than 0.001. The optimal weights for the individual components are: 0.12 for the residue−residue contacts, 0.31 for the atom−atom contacts, 1.79 for the membrane burial propensity and 0.47 for the SASA score.

**Overall Score.** All four scoring components were weighted and combined into one score and divided by the number of residues $N$. Thus, the complete formula to calculate the score is as follows:

$$\text{score} = (0.12 S_\text{residue\_residue} + 0.31 S_\text{atom\_atom} + 1.79 S_\text{membrane\_burial} + 0.47 S_\text{SASA})/N$$

The overall score is divided by the number of residues to be normalize for comparing structures of different size, but still anticipate to account for the observation that well modeled structures have more inter-residue interactions.[36] The four scoring components were chosen to reflect the anticipated change in conditions a kink will introduce to a protein structure. The residue−residue and atom−atom scoring components contribute to changes of the helix packing,

whereas the membrane-burial component reflects the positioning of residues along the membrane normal $z$ when a kink is introduced. Finally, the SASA scoring component rewards or penalizes the change of exposure of residues to the environment, which can occur when kinks of varying sizes are introduced.

**Kink Definition and Determination.** Kinked and straight helices were identified within the 102 crystal structures, compiled as the data set from which the scores were to be derived, by measuring the axis angle between two axes derived from a window of nine residues. For each protein chain, the transmembrane regions defined in the Orientation of Proteins in Membranes (OPM) were used as starting points. The N-terminal and C-terminal ends of the helices were then extended if the region was defined as helical in STRIDE[37] ($\alpha$, $3_{10}$, or $\pi$). A window of nine residues along the transmembrane helices was considered and for each window; the helical axis defined by the first four residues and the helical axis defined by the last four residues were calculated. The axis was calculated by finding points lying on the axis using the method of Kahn,[38] and then, a line was fit into these points using singular value decomposition to obtain a least-squares fit. The angle between the two axes was then calculated and assigned to the center residue in the windows. Only windows with the center residue in the transmembrane region were considered. To extract kinked and nonkinked segments, the nine residue windows were filtered using several steps. For each transmembrane helix, a list containing all windows was constructed. Beginning with the window with the highest axis angle, windows were removed from consideration by excluding all those whose center residue overlapped with the selected window. The procedure was then repeated using the remaining windows until all windows representing kinks were extracted from the list. This algorithm then supports the identification of more than one kink per helix. A window was defined as kinked if the angle between the two axes was 13° or larger. This value was chosen for consistency with previous methods.[15,39] The kink size is slightly lower than the average axis angle size in the database (15.2°). 1847 windows were defined within the database—712 were defined as kinked and 1135 as nonkinked. Around 60% of all transmembrane helices are kinked.

**Kink Sampling (Independent Test Set).** To validate the performance of the scoring function, 14 000 structures derived from the 14 GPCR crystal structures and 999 decoys for each GPCR structure were used to determine the ability of the scoring function to correctly find the crystal structure among a set of decoys with different kinks. The following describes the generation of the decoys for the independent test set. In each of the 14 GPCR crystal structures, the 7 transmembrane helices were defined by STRIDE accepting each residue in $\alpha$, $3_{10}$, or $\pi$ transmembrane helix. In each helix, the kink positions were defined by the kink determination algorithm and then one kink position for the sampling procedure was selected in each helix (see Supporting Information Table S2). To sample the kinks, the kink windows derived from the scoring data set (see kink definition and determination section) were used to replace the kink windows in the GPCR crystal structures. A Monte Carlo sampling was performed in 100 steps. In each step a kink windows from the crystal structure was selected randomly and replaced by a kink windows from the database following a Gaussian distribution, so that windows with a kink angle similar to the crystal structure kink angle are more likely to be selected than windows which will have a larger impact, measured by

RMSD, on the overall structure. After replacing the kink window with a new window, the N- and C-terminal region of the original GPCR helix without the side chains were translated and rotated so that the backbone atoms will fit to the new kink window. The new modeled helix was accepted if the new helix does not involve a stereochemical clash with any helix in the GPCR structure and the loops between the helices are long enough to connect the helices under the assumption that two residues are able to span 3.8 Å. Such a clash occurs if the distance between two atoms in the backbone is smaller than the sum of their vdW radii. After running the 100 steps, the new decoy structure was accepted and stored as a PDB file format. The procedure was repeated 999 times to get 999 decoys.

After sampling the structures, the side chains were added to the structure using SCRWL4.[40] The structures were minimized in 1000 steps using the steepest descent minimizer provided by the MMTK toolkit[41] in order to remove all clashes between the side chains and to refine the backbone in the interface between kink windows and N- and C-terminal region.

The crystal structures were included in the decoy set and the structures were scored using the scoring function. For the purpose of scoring, all the GPCR structures which were included in the 102 structures used for generating the score (PDB IDs: 3ODU, 2YDV, 2Y02, 3OAX, and 4EA3) were removed from the scoring database to avoid possible bias toward GPCR structures.
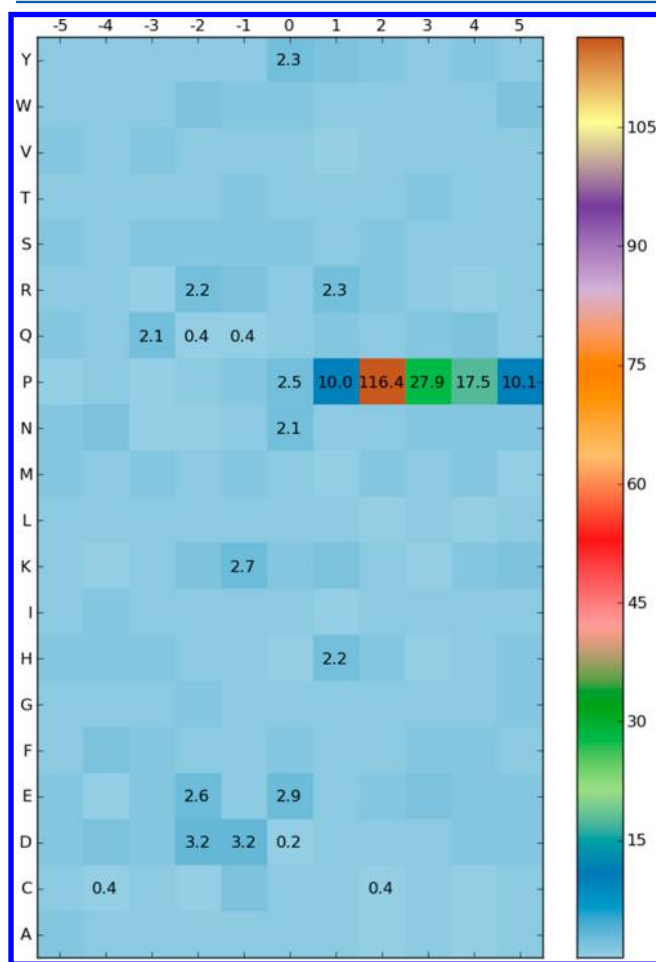
**GPCR Dock 2008 and GPCR Dock 2010 Data Sets.** In order to validate the performance of the proposed scoring function in comparison to potential functions developed specifically for membrane proteins, the scoring function and ProQM[9] were applied to the GPCR Dock 2008[42] and GPCR Dock 2010[43] data sets. The data sets contain models submitted to the community-wide assessment of GPCR structure modeling and ligand docking. The models were produced and submitted by participants before the experimental structures of the targets were released, and therefore, these models are appropriate for an objective evaluation of scoring functions. GPCR Dock 2008 consists of 206 models of the adenosine A2A receptor (PDB ID: 3EML) and GPCR Dock contains 166 models of the CXCR4 chemokine receptor (PDB ID: 3ODU) and 118 models of the human D3 dopamine receptor (PDB ID: 3PBL). The transmembrane helical regions were defined by STRIDE using the corresponding crystal structures. The loops of the soluble regions and the ligand were then accordingly removed to keep the seven $\alpha$-helices.

**Undersampling.** Undersampling is a problem for knowledge-based scoring functions and especially for transmembrane proteins because although the situation is slowly changing, there remains a dearth of membrane protein structures available and therefore certain data such as residue–residue and atom–atom type contacts remain rare. This produces low counts in the scoring matrices, and to address this issue, values in the scoring matrices which are based on fewer than two counts were discarded by setting the score for these contact types to zero. Therefore, these contacts were ignored when scoring protein structures providing no discrimination in the method. Furthermore, the smoothing density function used in the membrane burial propensity score and the SASA score reduces sampling errors which can occur especially in small data sets.

## ■ RESULTS AND DISCUSSION

**Analyses of Sequence and Structural Environment near Kinks.** When considering the populations of the amino

2930

dx.doi.org/10.1021/ci400236s | *J. Chem. Inf. Model.* 2013, 53, 2926–2936

acids, the results for the ratio of kinked to nonkinked windows were obtained and are shown in Figure 3 (a table containing all



**Figure 3.** Amino acid ratios of kinked to nonkinked windows from the database. The positions of the amino acid in the windows are labeled along the horizontal axis and the 20 amino acids are plotted on the vertical axis. The ratio is displayed if the amino acid is at least 2-fold over- (>2.0) or under-represented (<0.5) in kinked compared to nonkinked windows. A numerical scale is provided on the right side of the plot, and the squares are colored according to the scale.
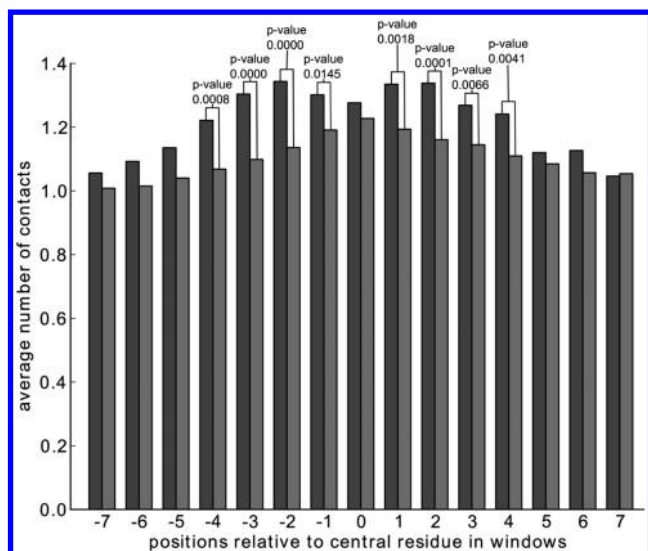
values can be found in Supporting Information Table S3). The kink center is defined as the center of the window which is the fifth residue in the window and labeled as 0 in the figure. As expected, Pro is predominant in the kink regions and it is more than 15-fold overrepresented at position 2 to 4 relative to the kink position and more than 2-fold overrepresented at position 0, 1, and 5. The bias of Pro toward the C-terminus of kinks is expected because of the loss of the hydrogen bond in the backbone due to the pyrollidine ring and steric clashes occur at residues N-terminal to the proline.[44] The spread over four positions suggests to us that it is not possible to define the exact kink position because of the variability of kinks and the shortfalls of the kink prediction algorithm, and even possibly due to the limitations in the methods for the determination of the structures. If a kink is to provide a change in "direction" for the helix, that the change occurs over a number of residues may be entirely appropriate. Proline was found in ∼40% of all kinked windows which is consistent with other reports.[45] However, 88% of all prolines in the transmembrane region are

within a kinked window and therefore the occurrence of proline in the transmembrane region is a strong indicator for kinks in this way. Although the difference of proline frequency between kinked and nonkinked is most evident, other residues are also over- or underrepresented near kinks (see Supporting Information Figure S3). The accumulation of charged residues is not surprising as polar residues might satisfy any broken hydrogen bond within the kink[39] or form hydrogen bonds to water molecules that often wedge kinks.[46] However, all residues which are slightly over- or underrepresented do not occur frequently in transmembrane helical regions and hence the observed frequencies may be possible arise from the limitation of the size of the data set. The preferences of Gly, Ser, or Thr observed by Hall et al.[39] were not observed here. They reported that Ser might have an impact on the stability of the proline-kinked helix by forming a hydrogen bond between the side-chain hydroxyl group of Ser and a backbone carbonyl oxygen.[47] The over- and underrepresented amino acids observed by Meruelo et al.[15] could not be completely matched. Furthermore, the analysis of Kneissl et al.[16] shows a much higher variation of at least 2-fold over- or under-represented amino acids. A reason why such variation might occur is that the kinks were assigned by manual inspection.

The overall results show that proline plays an important role in kink formation and perhaps the role is not entirely about creating the kink just at the "unattainable hydrogen bond". However, it is not clear whether or not other residues are involved in kink formation or the occurrence of these residues is biased by their low helical propensity in a nonpolar environment. Therefore, analysis in the future on larger data sets are necessary as more unique structures become available to reduce the bias to certain residue types.

The mismatch between membrane thickness and the length of the stretch of hydrophobic residues in the protein might cause kink formation,[48] and while it is well-known that often proline plays an important role, the main cause of kink formation is still unknown and is dependent on more than the identity of one amino acid. This might suggest that global effects of the protein structure dominate the cause, but it remains unclear whether global or local effects are the main reasons for kink formation. Many studies have focused on the local effects by analyzing the sequence adjacent to the kinks, whereas the global effects have largely been ignored. Therefore, we decided to consider the more global aspects which might affect on kink formation. The solvent accessible surface area of kinked windows and nonkinked windows were compared, and the number of contacts to neighboring helices were analyzed. The SASA was calculated by NACCESS[35] using the relative exposure of the side chains. The SASA in the kinked windows was significantly ($p$-value 0.0001) lower (27.09% accessibility; std dev 1.46%) than the SASA in the nonkinked windows (30.12% accessibility; std dev 1.20%). This result shows that residues in kinked windows are more buried than residues in nonkinked windows which is also in agreement with the observations of Kneissel et al.[16] To analyze these findings further, we calculated the average number of contacts to neighboring helices in kinked and nonkinked windows at each position (Figure 4). Residues in kinked windows have significantly more contacts with residues in neighboring helices at positions −4 to +4 around the kink position, which might suggest that kinks are stabilized by neighboring helices. However, the statistical differences do not necessarily indicate a physical significance. The statistical differences have evolu-
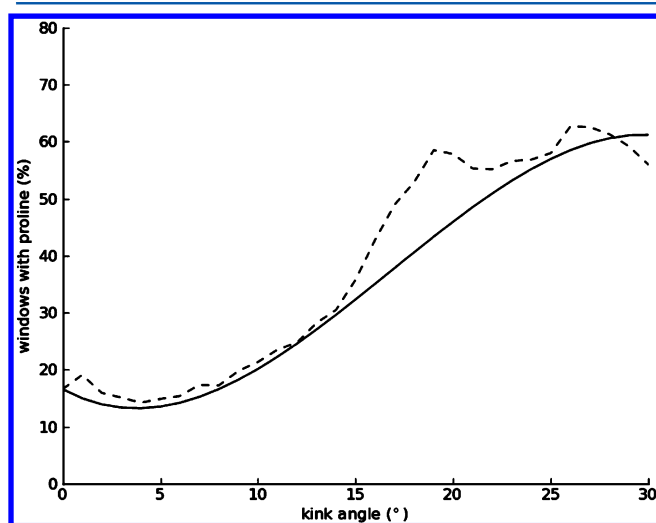
**Figure 4.** Position specific number of contacts to neighboring helices. The positions relative to the kink position are labeled on the *x*-axis. 0 is the kink position, positions with negative prefix are N-terminal to the kink, and the positions with positive prefix are C-terminal relative to the kink. On the *y*-axis, the average number of contacts to residues in neighboring helices is plotted. Only *p*-values of 0.05 or smaller are shown. Residue contacts in kinked windows are plotted in the dark gray bars, and residue contacts in nonkinked windows are shown in light gray.

tionary significance in that kinks might allow flexibility of highly packed helices to make adjustments to imposed changes. In structural considerations, it is not possible to determine if the observed higher number of contacts in kinked transmembrane helices in comparison to the nonkinked transmembrane helices is a causal link for the kink formation.

**Proline Preference for Larger Kinks.** The preference of proline for a larger kink was described previously.[49] Our data shows the same trend (Figure 5). The number of windows containing proline will increase with the kink angle size in the



**Figure 5.** Comparison of kink size to frequency of windows containing proline. The dotted curve displays the raw data and the continuous curve displays a line fitted to the raw data. The kink angle size is depicted on the *x*-axis and the percentage of kinked windows containing at least one proline is plotted on the *y*-axis.

windows. This result is perhaps not surprising because the backbone hydrogen bonds between proline at position i and the carbonyl groups at positions i-3 and i-4 cannot be formed with proline within the helix.[44] Moreover, the close proximity of the proline ring and the backbone carbonyl group at position i-4 represent a steric constraint to be avoided. The absent hydrogen bond provides flexibility in the helix and the kink has no restriction so that large kinks can prevail.

**Performance of the Scoring Function for the 5-fold Cross Validation Data Set.** The performance of the final scoring function was evaluated on the 5-fold cross validation data set by deriving the performance for each of the individual 5 test sets (see Supporting Information Table S4). The ⟨*z*-score⟩ of the scoring function with optimized weights is 4.91 (std dev 4.17). The average Pearson correlation coefficient is 0.78 (std dev 0.2), and most of the crystal structures were ranked first. The outlier in chain A of the cytochrome b6f complex (PDB ID: 2E74) shows a bad performance for all three measurements. The structure was ranked last, the Pearson correlation coefficient is 0.12, and the *z*-score is −2.75, which means that the crystal structure is scored worse than the decoys. However, the RMSD of the decoys to the crystal structure is 0.3 Å, which means that the decoys are almost identical to the crystal structure, and therefore, it is difficult for the scoring function to distinguish between crystal and decoy structures. On the other hand, the performance of the scoring function on the carnitine transporter (PDB ID: 2WSW) is encouraging. In this case, the crystal structure ranked first, the *z*-score is 12.44 and the Pearson correlation coefficient is 0.96. The average RMSD of the decoys to the crystal structure is 0.46 Å, which also place it into the category of difficult targets for the scoring function. Overall, the analyses show that the scoring function performs reasonably on the 5-fold cross validation data set.

**GPCR Ranking and Correlation Coefficient.** In theory, the scoring function can be used for modeling of transmembrane proteins without the focus on kinks, but the scoring components were chosen to optimize the performance of the scoring function to correctly predict kinks in transmembrane proteins. Therefore, in order to validate the ability of the scoring function to distinguish between structures with native fold and misfolded structures with different kinks sizes, the *z*-score for the independent test data set was calculated and the ranks for the crystal structure as well as for the decoys were determined (Table 2). The table shows that the scoring function is able to discriminate between the crystal structure and the decoys. An ideal scoring function should produce a high correlation between the score and the quality of the model to the native structures. The RMSD metric was chosen as a quality measure over other structural similarity measurements because the majority of the decoys are close to the native structure (<3.5 Å) and no loops are present in the structures. Poor modeling of the loops has the potential to have a huge impact on the overall RMSD in a structure which is otherwise reasonably accurate. The correlation between the score and the RMSD to the crystal structure was calculated using the Pearson's correlation. Pearson's correlation was chosen because a correlation between score and quality of the models was considered to be linear. The Pearson's correlation coefficient for the GPCRs shows that the score is mostly well-correlated to the RMSD (see Table 2 and Supporting Information Figure S4) and therefore the score can be successfully used to generate a ranked list for membrane protein kink prediction. The Pearson's correlation coefficient is between 0.46 and 0.72 with

**Table 2. Rank, $z$-Score, and Pearson's Correlation Coefficient for the 14 Different GPCR Crystal Structures within the Independent Test Set[a]**
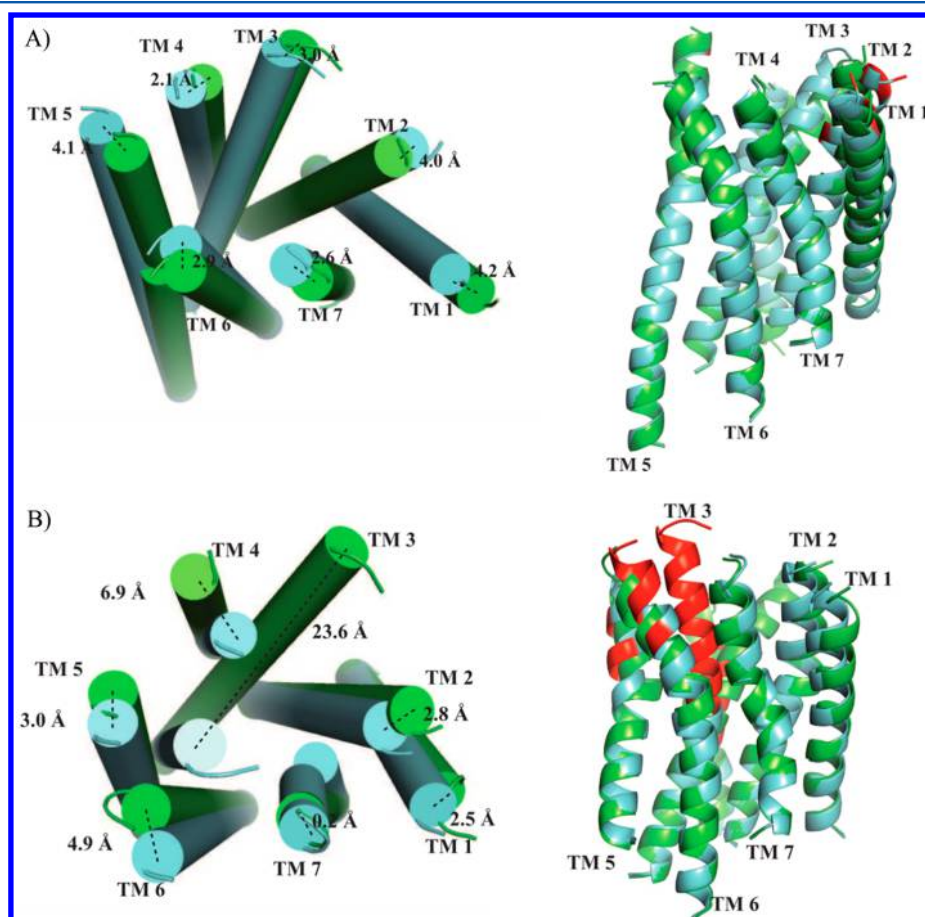
| PDB ID | rank | $z$-score | Pearson's correlation coefficient | RMSD (Å) from best scored decoy to crystal structure |
|---|---|---|---|---|
| 1U19 | 1 | 2.92 | 0.63 | 3.51 |
| 2RH1 | 1 | 2.42 | 0.66 | 0.83 |
| 2VT4 | 1 | 3.12 | 0.46 | 2.06 |
| 3UON | 1 | 3.17 | 0.53 | 3.48 |
| 4EA3 | 1 | 2.44 | 0.63 | 2.92 |
| 4DAJ | 2 | 2.62 | 0.50 | 2.68 |
| 3PBL | 3 | 1.93 | 0.71 | 1.72 |
| 4DKL | 5 | 2.25 | 0.67 | 2.56 |
| 3RZE | 6 | 2.23 | 0.61 | 4.58 |
| 3V2Y | 9 | 1.97 | 0.72 | 2.15 |
| 4EJ4 | 9 | 1.94 | 0.66 | 1.71 |
| 3EML | 14 | 2.03 | 0.57 | 2.59 |
| 3ODU | 19 | 1.83 | 0.46 | 4.36 |
| 4DJH | 20 | 1.69 | 0.63 | 1.56 |

[a]The average Pearson's correlation coefficient is 0.6, whereas 3V2Y has the highest with 0.72 and 2VT4 the lowest with 0.46. The RMSD for the comparison of the crystal structure and the best scored decoy structure for the eight structures which were not ranked first in the independent test set are also shown.

an average of 0.60. The rank for each crystal structure is also shown in Table 2. In all test cases, the crystal structure is ranked within the best 2% of the scored structures. Six of the 14 structures were ranked first and only three structures were not within the 1% of the best scored structures. The $z$-score ranks in a range between 1.69 and 3.17 with an average of 2.33. Interestingly, the Pearson's correlation coefficient is not always highest in the case in which the $z$-score or the rank of the structures is the highest. For example, the turkey $\beta_1$ adrenergic receptor (2VT4) is ranked first and the $z$-score of 3.12 is one of the highest in the data set; however the Pearson correlation coefficient is with 0.46 the lowest.
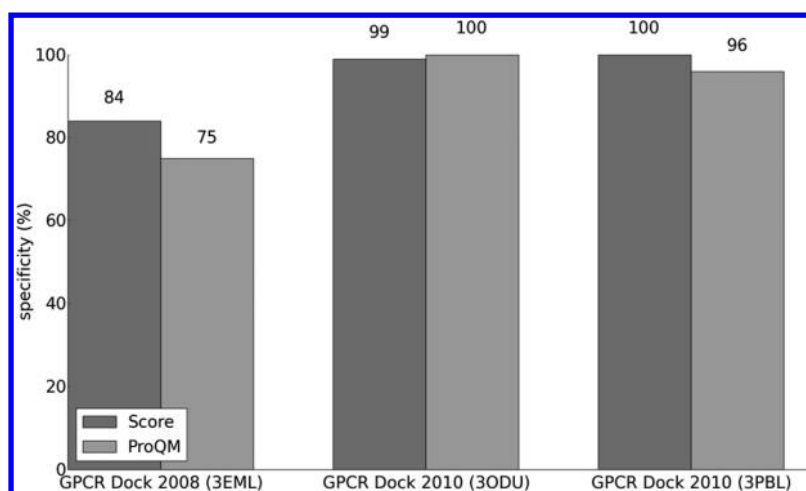
**Comparison of the Crystal Structure to the Best Scored Structure.** Nine GPCR structures were not ranked first using the scoring function, yet the scoring function was able to rank the nine native structures among the top 20 out of 1000. This is significantly better than random ($p$-value $\approx$ 0). The RMSDs of the highest scored decoy to the corresponding crystal structure vary between 0.83 and 4.58 Å with an average of 2.62 Å (std dev 1.07) as displayed in Table 2.

Superposition of the native structure to the best model selected by the scoring function for the human kappa opioid receptor (PDB ID: 4DJH)[50] and the human histamine H1 receptor (PDB ID: 3RZE)[51] are shown in Figure 6. These structures demonstrate the strength and weakness of the scoring function in the cases in which the crystal structure is not ranked first by showing the results for the best scored



**Figure 6.** (A) Superposition of the crystal structure of the human kappa opioid receptor (4DJH) (green) and the best scored decoy structure in blue and (B) the crystal structure of the human histamine H1 receptor (3RZE) in green and the best scored decoy (blue). Regions found to be >3.5 Å from the crystal structure are displayed in red.

**Figure 7.** Comparison of the proposed scoring function with ProQM using the GPCR Dock 2008 and GPCR Dock 2010 data sets.

structure with the lowest RMSD to the native structure and the results for the best scored structure with the highest RMSD to the native structure. As shown in Figure 6A, the best scored structure for 4DJH is similar to the crystal structure. Only eight residues have a distance of 3.5 Å or higher to the corresponding residue in the native structure. The most significant differences between the best scored structure and the native structure are observed in helix 1 and helix 2, neither of which are involved in forming the binding site.

Figure 6B shows the superposition of 3RZE crystal structure with the corresponding model we derived. The relatively high RMSD between the two structures is mainly because of helix 3 which kink is highly different to the crystal structure kink resulting into a helix which axis top is in a distance of 23.6 Å from the crystal structure. 39 residues in the decoy structure have a distance of more than 3.5 Å to the corresponding crystal structure residues and 17 of the 39 residues are within helix 3. The movement of helix 3 in toward the interior of the structure allows a tighter helix packing, resulting in a better score, which is possible because of the hollow in the binding site which is not occupied by any ligand. However, the integrity of the native binding site in the structure has been lost, and therefore, the model is not useful for drug design, which can be easily filtered out in the kink sampling step by removing models without the binding site in a structure-based drug design study. However, the purpose of the sampling was the proof of concept and not to generate models for drug design

We believe that the results demonstrate the power of the scoring functions for kink prediction. However, care is still required for selecting an appropriate model because the scoring function takes not into account ligands and possible binding sites.

**GPCR Dock 2008 and GPCR Dock 2010 Data Sets Scoring Results.** The GPCR Dock 2008 data set consists of 206 models of the A2A adenosine receptor. Many models are similar because most participants built homology models based on the structure of $\beta_2$-adrenergic receptor, which is also reflected by the low standard deviation (0.58 Å) of the RMSD in the transmembrane helices (average RMSD: 2.65 Å). The scoring function implemented here ranks the crystal structure at position 34 of 206 models, the Pearson's correlation coefficient is 0.26 and the z-score is 0.81.

The GPCR Dock 2010 contains 166 models of the CXCR4 chemokine receptor (CXCR4) and 118 models of the human

D3 dopamine receptor (D3). The average RMSDs in the transmembrane helices to the crystal structures are 3.74 (std dev 1.37) and 2.48 Å (std dev 2.88), respectively. The CXCR4 crystal structure is ranked second, the Pearson's correlation coefficient is 0.39, and the z-score is 1.75. In the third data set, the D3 crystal structure is ranked first, the Pearson's correlation coefficient is 0.84, and the z-score is 1.46.

**Comparison of the Scoring Function with ProQM.** In order to compare the proposed scoring function with a well-established membrane specific scoring function, the learning-based model quality assessment program ProQM was applied to the GPCR Dock 2008 and GPCR Dock 2010 data sets. ProQM is able to select the crystal structures among all of the models with a specificity of 75.3, 100, and 95.8%, which is slightly lower than the specificity achieved for the scoring function presented here (83.6, 98.8, and 100%) (Figure 7). In the comparison, the z-score for ProQM is also lower in all three data sets and the Pearson's correlation coefficient is lower in two of the data sets (see Supporting Information Figure S5). Only the Pearson's correlation coefficient for the GPCR Dock 2008 data set is slightly higher for ProQm (0.36) than for the proposed scoring function (0.26).

**Statistical Evaluation.** The scoring function was carefully evaluated in a two-step procedure to ensure the quality of the score and to remove any possible bias in the result. The 5-fold cross validation procedure supports the reliability of the score and the choice of the scoring components as well as the weighting coefficients. In an additional step, the scoring function was evaluated using an independent test set of 14 000 structures derived from the 14 GPCRs. The 14 GPCR crystal structures were removed from the procedure in which the score was generated to ensure a truly independent test set with minimal bias to the GPCRs. Although only 14 GPCR structures were available, GPCRs are one of the most important targets for drug design. Furthermore, moving only a small number of structures from the training set to the independent test set ensures that the training set is sufficiently diverse to have adequate information to derive a meaningful score.

## ■ CONCLUSION

To the best of the authors' knowledge, this is the first time that a method was specifically developed for kink modeling. Methods for kink position prediction or experimentally derived information can be used to find the position of a kink and then

the modeling approach here presented can be applied to sample and select appropriate conformations. We have laid the groundwork for a predictive method; however, it is likely that the modeling results will be improved over time as more structures become available. In spite of the limited data set, we were able to develop a modeling method for introducing kinks, which can be used in many modeling approaches to improve the quality of models. Furthermore, the performance of the scoring function proposed in this study was compared to ProQM, a widely used model quality assessment program, showing that the scoring function at hand performs slightly better but effectively both approaches perform well in this stringent test. As there is a high prevalence of kinking in transmembrane domains at 60%, of all helices at the current time we anticipate methods for the prediction of the transmembrane structures will require the incorporation of algorithms as the one we have reported here.

Comparing several studies is always difficult because the data on which the observations are based might vary, depending on the choice of the nonredundancy cutoffs, or the number of available structures at the date the database was constructed. This is more especially true for the small data sets like transmembrane proteins data sets. Furthermore, the cutoff for categorizing windows into kinked or nonkinked will impact the outcome. However, strong signature of the proline over-representation within the kink windows or the number of contacts to neighboring helices will probably persist over time. Other observations such as the under- or over-representation of charged residues should be taken with a grain of salt.

Our results indicate that there are significant differences in the sequence as well as the structural environment between kinked and nonkinked transmembrane helices. The results show that kinks are generated by local as well as global effects. However, we cannot necessary link cause and effect. For example, it is not clear if kinks are caused by local effects and subsequently stabilized by neighboring helices forming long-range contacts to kinked transmembrane helices or if kinks are generated by the long-range contacts of neighboring helices in the first place. Furthermore, proline is the only residue which is currently clearly overrepresented near kinks. While this indicates that proline is associated with factors beneficial for kink formation, proline is not present in all kinks and therefore is not a necessary condition for kink formation. We believe that both local and global effects are important for kink formation and therefore the sequence alone is not sufficient for reliable kink predictions especially for predicting the kink size. The origin of sequence preferences other than for proline are not clear and the associations themselves currently seem to vary with the criteria for choosing the data sets. Further analysis with larger data sets of more unique structures is required, as they are necessary to achieve the thresholds for significance to determine if other amino acids are involved in kink formation.

## ASSOCIATED CONTENT

### Ⓢ Supporting Information

GPCR kink size distribution, definition of donor/acceptor hydrogen atoms in amino acids, membrane burial propensity score distribution along the membrane normal, SASA score distribution along the membrane normal, kink position and size in each GPCR crystal structure, amino acid ratios of kinked to nonkinked windows from the database, performance of the scoring function using the 5-fold cross validation data set, scatter plot of RMSD versus score for each GPCR in the independent test set, comparison of ProQM with the proposed scoring function, and 102 structures used for the kink analysis and scoring generating. This material is available free of charge via the Internet at http://pubs.acs.org

## AUTHOR INFORMATION

### Corresponding Author
*E-mail: bret.church@sydney.edu.au. Tel.: +61 2 9036 6569.
### Notes
The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

## ABBREVIATIONS

OPM, Orientations of Proteins in Membranes database; GPCR, G-protein coupled receptor; RMSD, root-mean-square-deviation; vdW, van der Waals; SASA, solvent accessible surface area;

## REFERENCES

(1) Krogh, A.; Larsson, B.; von Heijne, G.; Sonnhammer, E. L. Predicting Transmembrane Protein Topology with a Hidden Markov Model: Application to Complete Genomes. *J. Mol. Biol.* **2001**, *305*, 567−580.

(2) Granseth, E.; Daley, D. O.; Rapp, M.; Melén, K.; von Heijne, G. Experimentally Constrained Topology Models for 51,208 Bacterial Inner Membrane Proteins. *J .Mol. Biol.* **2005**, *352*, 489−494.

(3) Bordoli, L.; Kiefer, F.; Arnold, K.; Benkert, P.; Battey, J.; Schwede, T. Protein Structure Homology Modeling Using SWISS-MODEL Workspace. *Nat. Protoc.* **2009**, *4*, 1−13.

(4) Chothia, C.; Lesk, A. M. The Relation Between the Divergence of Sequence and Structure in Proteins. *EMBO J.* **1986**, *5*, 823−826.

(5) Forrest, L. R.; Tang, C. L.; Honig, B. On the Accuracy of Homology Modeling and Sequence Alignment Methods Applied to Membrane Proteins. *Biophys. J.* **2006**, *91*, 508−517.

(6) Pellegrini-Calace, M.; Carotti, A.; Jones, D. T. Folding in Lipid Membranes (FILM): a Novel Method for the Prediction of Small Membrane Protein 3D Structures. *Proteins* **2003**, *50*, 537−545.

(7) Barth, P.; Wallner, B.; Baker, D. Prediction of Membrane Protein Structures with Complex Topologies Using Limited Constraints. *Proc. Natl. Acad. Sci. U.S.A.* **2009**, *106*, 1409−1414.

(8) Heim, A. J.; Li, Z. Developing a High-Quality Scoring Function for Membrane Protein Structures Based on Specific Inter-Residue Interactions. *J. Comput. Aided. Mol. Des.* **2012**, *26*, 301−309.

(9) Ray, A.; Lindahl, E.; Wallner, B. Model Quality Assessment for Membrane Proteins. *Bioinformatics* **2010**, *26*, 3067−3074.

(10) Bright, J. N.; Shrivastava, I. H.; Cordes, F. S.; Sansom, M. S. P. Conformational Dynamics of Helix S6 from Shaker Potassium Channel: Simulation Studies. *Biopolymers* **2002**, *64*, 303−313.

(11) Topiol, S.; Sabio, M. X-Ray Structure Breakthroughs in the GPCR Transmembrane Region. *Biochem. Pharmacol.* **2009**, *78*, 11−20.

(12) Ago, H.; Kanaoka, Y.; Irikura, D.; Lam, B. K.; Shimamura, T.; Austen, K. F.; Miyano, M. Crystal Structure of a Human Membrane Protein Involved in Cysteinyl Leukotriene Biosynthesis. *Nature* **2007**, *448*, 609−612.

(13) Miyano, M.; Ago, H.; Saino, H.; Hori, T.; Ida, K. Internally Bridging Water Molecule in Transmembrane Alpha-Helical Kink. *Curr. Opin. Struct. Biol.* **2010**, *20*, 456−463.

(14) Rigoutsos, I.; Riek, P.; Graham, R. M.; Novotny, J. Structural Details (kinks and non-α conformations) in Transmembrane Helices Are Intrahelically Determined and Can Be Predicted by Sequence Pattern Descriptors. *Nucleic Acids Res.* **2003**, *31*, 4625−4631.

(15) Meruelo, A. D.; Samish, I.; Bowie, J. U. TMKink: a Method to Predict Transmembrane Helix Kinks. *Protein Sci.* **2011**, *20*, 1256–1264.

(16) Kneissl, B.; Mueller, S. C.; Tautermann, C. S.; Hildebrandt, A. String Kernels and High-Quality Data Set for Improved Prediction of Kinked Helices in α-Helical Membrane Proteins. *J. Chem. Inf. Model.* **2011**, *51*, 3017–3025.

(17) Vaidehi, N.; Bhattacharya, S. Multiscale Computational Methods for Mapping Conformational Ensembles of G-Protein-Coupled Receptors. *Adv. Protein Chem. Struct. Biol.* **2011**, *85*, 253–280.

(18) Barth, P.; Schonbrun, J.; Baker, D. Toward High-Resolution Prediction and Design of Transmembrane Helical Protein Structures. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104*, 15682–15687.

(19) Hopkins, A. L.; Groom, C. R. The Druggable Genome. *Nat. Rev. Drug Discov.* **2002**, *1*, 727–730.

(20) Jaakola, V.-P.; Griffith, M. T.; Hanson, M. A.; Cherezov, V.; Chien, E. Y. T.; Lane, J. R.; IJzerman, A. P.; Stevens, R. C. The 2.6 Angstrom Crystal Structure of a Human A2A Adenosine Receptor Bound to an Antagonist. *Science.* **2008**, *322*, 1211–1217.

(21) Warne, T.; Serrano-Vega, M. J.; Baker, J. G.; Moukhametzianov, R.; Edwards, P. C.; Henderson, R.; Leslie, A. G. W.; Tate, C. G.; Schertler, G. F. X. Structure of a Beta1-Adrenergic G-Protein-Coupled Receptor. *Nature* **2008**, *454*, 486–491.

(22) Haga, K.; Kruse, A. C.; Asada, H.; Yurugi-Kobayashi, T.; Shiroishi, M.; Zhang, C.; Weis, W. I.; Okada, T.; Kobilka, B. K.; Haga, T.; et al. Structure of the Human M2Muscarinic Acetylcholine Receptor Bound to an Antagonist. *Nature* **2012**, *482*, 547–551.

(23) Lomize, M. A.; Lomize, A. L.; Pogozheva, I. D.; Mosberg, H. I. OPM: Orientations of Proteins in Membranes Database. *Bioinformatics* **2006**, *22*, 623–625.

(24) Li, W.; Godzik, A. Cd-Hit: a Fast Program for Clustering and Comparing Large Sets of Protein or Nucleotide Sequences. *Bioinformatics* **2006**, *22*, 1658–1659.

(25) Zhang, Y.; Skolnick, J. TM-Align: a Protein Structure Alignment Algorithm Based on the TM-Score. *Nucleic Acids Res.* **2005**, *33*, 2302–2309.

(26) Larkin, M. A.; Blackshields, G.; Brown, N. P.; Chenna, R.; McGettigan, P. A.; McWilliam, H.; Valentin, F.; Wallace, I. M.; Wilm, A.; Lopez, R.; et al. Clustal W and Clustal X Version 2.0. *Bioinformatics* **2007**, *23*, 2947–2948.

(27) Sali, A.; Blundell, T. L. Comparative Protein Modelling by Satisfaction of Spatial Restraints. *J. Mol. Biol.* **1993**, *234*, 779–815.

(28) Siew, N.; Elofsson, A.; Rychlewski, L.; Fischer, D. MaxSub: An Automated Measure for the Assessment of Protein Structure Prediction Quality. *Bioinformatics* **2000**, *16*, 776–785.

(29) Sippl, M. J. Recognition of Errors in Three-Dimensional Structures of Proteins. *Proteins* **1993**, *17*, 355–362.

(30) Bowie, J. U.; Lüthy, R.; Eisenberg, D. A Method to Identify Protein Sequences That Fold into a Known Three-Dimensional Structure. *Science* **1991**, *253*, 164–170.

(31) Li, A. J.; Nussinov, R. A Set of van Der Waals and Coulombic Radii of Protein Atoms for Molecular and Solvent-Accessible Surface Calculation, Packing Evaluation, and Docking. *Proteins* **1998**, *32*, 111–127.

(32) Wallner, B.; Elofsson, A. Can Correct Protein Models Be Identified? *Protein Sci.* **2003**, *12*, 1073–1086.

(33) Yau, W. M.; Wimley, W. C.; Gawrisch, K.; White, S. H. The Preference of Tryptophan for Membrane Interfaces. *Biochemistry* **1998**, *37*, 14713–14718.

(34) Ulmschneider, M. B.; Sansom, M. S. P.; Di Nola, A. Properties of Integral Membrane Protein Structures: Derivation of an Implicit Membrane Potential. *Proteins* **2005**, *59*, 252–265.

(35) Hubbard, S.; Thornton, J. *Naccess*; 1993.

(36) Gao, J.; Li, Z. Comparing Four Different Approaches for the Determination of Inter-Residue Interactions Provides Insight for the Structure Prediction of Helical Membrane Proteins. *Biopolymers* **2009**, *91*, 547–556.

(37) Frishman, D.; Argos, P. Knowledge-Based Protein Secondary Structure Assignment. *Proteins* **1995**, *23*, 566–579.

(38) Kahn, P. C. Defining the Axis of a Helix. *Comput. Chem.* **1989**, *13*, 185–189.

(39) Hall, S. E.; Roberts, K.; Vaidehi, N. Position of Helical Kinks in Membrane Protein Crystal Structures and the Accuracy of Computational Prediction. *J. Mol. Graph. Model.* **2009**, *27*, 944–950.

(40) Krivov, G. G.; Shapovalov, M. V; Dunbrack, R. L. Improved Prediction of Protein Side-Chain Conformations with SCWRL4. *Proteins* **2009**, *77*, 778–795.

(41) Hinsen, K. The Molecular Modeling Toolkit: A New Approach to Molecular Simulations. *J. Comput. Chem.* **2000**, *21*, 79–85.

(42) Michino, M.; Abola, E.; Brooks, C. L.; Dixon, J. S.; Moult, J.; Stevens, R. C. Community-Wide Assessment of GPCR Structure Modelling and Ligand Docking: GPCR Dock 2008. *Nat. Rev. Drug Discov.* **2009**, *8*, 455–463.

(43) Kufareva, I.; Rueda, M.; Katritch, V.; Stevens, R. C.; Abagyan, R. Status of GPCR Modeling and Docking as Reflected by Community-Wide GPCR Dock 2010 Assessment. *Structure* **2011**, *19*, 1108–1126.

(44) Von Heijne, G. Proline Kinks in Transmembrane α-Helices. *J. Mol. Biol.* **1991**, *218*, 499–503.

(45) Langelaan, D. N.; Wieczorek, M.; Blouin, C.; Rainey, J. K. Improved Helix and Kink Characterization in Membrane Proteins Allows Evaluation of Kink Sequence Predictors. *J. Chem. Inf. Model.* **2010**, *50*, 2213–2220.

(46) Miyano, M.; Ago, H.; Saino, H.; Hori, T.; Ida, K. Internally Bridging Water Molecule in Transmembrane Alpha-Helical Kink. *Curr. Opin. Struct. Biol.* **2010**, *20*, 456–463.

(47) Weber, M.; Tome, L.; Otzen, D.; Schneider, D. A Ser Residue Influences the Structure and Stability of a Pro-Kinked Transmembrane Helix Dimer. *Biochim. Biophys. Acta* **2012**, *1818*, 2103–2107.

(48) Park, S. H.; Opella, S. J. Tilt Angle of a Trans-Membrane Helix Is Determined by Hydrophobic Mismatch. *J. Mol. Biol.* **2005**, *350*, 310–318.

(49) Cordes, F. S.; Bright, J. N.; Sansom, M. S. P. Proline-Induced Distortions of Transmembrane Helices. *J. Mol. Biol.* **2002**, *323*, 951–960.

(50) Wu, H.; Wacker, D.; Mileni, M.; Katritch, V.; Han, G. W.; Vardy, E.; Liu, W.; Thompson, A. A.; Huang, X.-P.; Carroll, F. I.; et al. Structure of the Human κ-Opioid Receptor in Complex with JDTic. *Nature* **2012**, *485*, 327–332.

(51) Shimamura, T.; Shiroishi, M.; Weyand, S.; Tsujimoto, H.; Winter, G.; Katritch, V.; Abagyan, R.; Cherezov, V.; Liu, W.; Han, G. W.; et al. Structure of the Human Histamine H1 Receptor Complex with Doxepin. *Nature* **2011**, *475*, 65–70.

2936

dx.doi.org/10.1021/ci400236s | *J. Chem. Inf. Model.* 2013, 53, 2926–2936