

# TMSEG: Novel prediction of transmembrane helices

Michael Bernhofer,<sup>1\*</sup> Edda Kloppmann,<sup>1,2</sup> Jonas Reeb,<sup>1</sup> and Burkhard Rost<sup>1,2,3,4</sup>

<sup>1</sup>Department of Informatics & Center for Bioinformatics & Computational Biology – i12, Technische Universität München (TUM), Boltzmannstr. 3, Garching/Munich 85748, Germany

<sup>2</sup>New York Consortium on Membrane Protein Structure, New York Structural Biology Center, New York, New York 10027

<sup>3</sup>Institute of Advanced Study (TUM-IAS), Lichtenbergstr. 2a, Garching/Munich 85748, Germany

<sup>4</sup>Institute for Food and Plant Sciences WZW – Weihenstephan, Alte Akademie 8, Freising, Germany

## ABSTRACT

Transmembrane proteins (TMPs) are important drug targets because they are essential for signaling, regulation, and transport. Despite important breakthroughs, experimental structure determination remains challenging for TMPs. Various methods have bridged the gap by predicting transmembrane helices (TMHs), but room for improvement remains. Here, we present TMSEG, a novel method identifying TMPs and accurately predicting their TMHs and their topology. The method combines machine learning with empirical filters. Testing it on a non-redundant dataset of 41 TMPs and 285 soluble proteins, and applying strict performance measures, TMSEG outperformed the state-of-the-art in our hands. TMSEG correctly distinguished helical TMPs from other proteins with a sensitivity of  $98 \pm 2\%$  and a false positive rate as low as  $3 \pm 1\%$ . Individual TMHs were predicted with a precision of  $87 \pm 3\%$  and recall of  $84 \pm 3\%$ . Furthermore, in  $63 \pm 6\%$  of helical TMPs the placement of all TMHs and their inside/outside topology was correctly predicted. There are two main features that distinguish TMSEG from other methods. First, the errors in finding all helical TMPs in an organism are significantly reduced. For example, in human this leads to 200 and 1600 fewer misclassifications compared to the second and third best method available, and 4400 fewer mistakes than by a simple hydrophobicity-based method. Second, TMSEG provides an add-on improvement for any existing method to benefit from.

Proteins 2016; 84:1706–1716.  
© 2016 Wiley Periodicals, Inc.

**Key words:** membrane protein; protein structure prediction; transmembrane helices;  $\alpha$ -helical integral membrane protein; transmembrane protein prediction; transmembrane helix prediction.

## INTRODUCTION

Transmembrane proteins (TMPs) are involved in numerous essential processes within living organisms such as signaling, regulation, and transport.<sup>1</sup> About 20–30% of all proteins within any organism have been estimated to be TMPs.<sup>2,3</sup> Many TMPs, especially G protein-coupled receptors (GPCRs), are primary drug targets<sup>4</sup> and therefore of high interest.

TMPs cross the membrane bilayer with either transmembrane helices (TMHs) or beta-strands. The latter are found in the outer membrane of Gram-negative bacteria, mitochondria, and chloroplasts. They make up only about 1–2% of all proteins in Gram-negative bacteria.<sup>5</sup> We concentrated on the more common class of helical TMPs and will refer to these as TMPs in the following. TMPs can cross the membrane only once (single-pass) or

multiple times (multi-pass). Due to the apolar and hydrophobic environment in the lipid bilayer, most of the amino acids found in TMHs are hydrophobic, and their orientation in the membrane (called TMP topology) can be discerned through Gunnar von Heijne's positive-inside rule.<sup>6,7</sup>

Additional Supporting Information may be found in the online version of this article.

Abbreviations used: 3D, three-dimensional; GPCR, G protein-coupled receptor; NN, (artificial) neural network; OPM, Orientations of Proteins in Membranes; PDB, Protein Data Bank; PDBTM, Protein Data Bank of Transmembrane Proteins; RF, random forest; TMH, transmembrane alpha-helix; TMP, transmembrane protein.

Grant sponsor: Alexander von Humboldt Foundation; Grant sponsor: National Institutes of Health (NIH); Grant number: U54 GM095315.

\*Correspondence to: Michael.Bernhofer@mytum.de

Received 22 May 2016; Revised 18 July 2016; Accepted 24 August 2016

Published online 26 August 2016 in Wiley Online Library (wileyonlinelibrary.com). DOI: 10.1002/prot.25155

Despite their immense importance, and despite crucial experimental advances,<sup>8–11</sup> <2% of the structures in the Protein Data Bank<sup>12</sup> (PDB) are TMPs.<sup>13–15</sup> As membrane regions are typically not visible in high-resolution structures, TMHs are assigned to PDB structures by expert resources, most prominently the Orientations of Proteins in Membranes<sup>16</sup> (OPM) database and the Protein Data Bank of Transmembrane Proteins<sup>17</sup> (PDBTM).

Recent advances in experimental structure determination have benefited from advanced computational predictions of TMHs from sequence.<sup>8,9</sup> In the last 25 years, many such tools have been developed, ranging from simple algorithms based solely on hydrophobicity scales (e.g., TopPred<sup>18</sup>) to sophisticated uses of hidden Markov models (e.g., TMHMM,<sup>19</sup> HMMTOP,<sup>20</sup> Phobius,<sup>21</sup> and PolyPhobius<sup>22</sup>), neural networks (e.g., PHDhtm,<sup>23,24</sup> and MEMSAT3<sup>25</sup>), and support vector machines (MEMSAT-SVM<sup>26</sup>). Arguably, the most important advance was the incorporation of evolutionary information from sequence profiles or multiple sequence alignments.<sup>23,24</sup> Consequently, almost all methods developed over the last decade are based on evolutionary information. A recent assessment applying strict evaluation measures showed that many methods perform well overall; the best are some recent methods.<sup>27</sup> Here, we show that a few simple ideas improve significantly over the state-of-the-art.

## MATERIAL AND METHODS

### Dataset TMP166: helical TMPs with known structures

We collected helical TMPs with known structures annotated in OPM<sup>16</sup> and PDBTM<sup>17</sup> (releases 2013\_07). Both databases use PDB<sup>12</sup> chain identifiers. We mapped those PDB chains to their UniProtKB<sup>28</sup> protein sequences using SIFTS.<sup>29</sup> We excluded all chimeric PDB chains, model structures, X-ray structures with >8 Å, and those for which some TMH residues did not map gapless to UniProtKB sequences. This gave 1087 PDB chains from 455 PDB structures (379 X-ray and 76 NMR structures).

UniqueProt<sup>30</sup> reduced sequence-redundancy at HVAL > 0 (the HVAL depends on alignment length and the percentage of pairwise sequence identity<sup>31</sup>). At this threshold, no pair of proteins has >20% pairwise sequence identity for alignments of >250 residues (see Rost 1999<sup>32</sup> for precise definitions). The result of this is our final dataset consisting of 166 non-redundant TMPs (called TMP166, Supporting Information Table S1).

As the TMH annotations in OPM and PDBTM differed for some proteins, we associated TMH annotations from both databases with each sequence. The inside/outside topology of the non-transmembrane regions was assigned based on the ATOM coordinates and topology annotation from OPM (cf. Note Supporting Information S1 and Fig. S1). We considered re-entrant regions<sup>33,34</sup>

to be non-transmembrane due to their scarcity in the TMP166 dataset (only 15 proteins with one or two re-entrant regions each; Supporting Information Table S1).

### Dataset SP1441: proteins with and without signal peptides

As signal peptides are often confused with TMHs and vice versa,<sup>27</sup> a second dataset was derived from the SignalP4.1 dataset.<sup>35</sup> This dataset contained UniProtKB sequences of soluble proteins and TMPs with and without signal peptide annotations. Note that these TMPs have no inside/outside topology annotations and many of their TMH annotations are not supported by experimental evidence.

The SignalP4.1 dataset was redundancy reduced twice using UniqueProt. First, all proteins similar to any of those in the TMP166 dataset were removed at HVAL > 0. Second, the remaining proteins were redundancy-filtered at HVAL > 0. The final dataset contained 1441 protein sequences (299 TMPs and 1142 soluble proteins, called SP1441; Supporting Information Table S2). About 477 of those had signal peptide annotations (25 TMPs and 452 soluble proteins).

### Splitting the datasets

We split the combined TMP166 and SP1441 dataset into four subsets. We partitioned them in a way that all subsets have approximately the same distributions with respect to the number of soluble proteins and TMPs, protein sequences with and without signal peptides, and sequence lengths (Supporting Information Fig. S2).

We used the first three subsets to develop TMSEG in a three-fold cross-validation approach (cf. TMSEG training). The fourth split, the independent test set called BlindTest, was used only for the final performance evaluation, i.e., no parameter was optimized on that set. The BlindTest dataset contained 41 TMPs (from TMP166) with known structure and TMH annotations from OPM and PDBTM, and 285 soluble proteins from the SP1441 dataset. The 74 TMPs from the fourth split of SP1441 (Supporting Information Table S2) were not included in the BlindTest dataset, because they lack sufficient experimental annotations. However, we used them for the signal peptide prediction performance analysis, as we did not have curated signal peptide annotations for the TMPs from OPM and PDBTM.

### Human proteome

We retrieved the human proteome, 20,196 protein sequences, from UniProtKB/Swiss-Prot (release 2015\_03). We applied our TMSEG algorithm to the whole proteome to provide a summary of its TMP composition and to estimate run time.

**Table 1**  
Evaluation Measures

Measurement	Formula	Description
Precision (%)	$100 * \frac{\# \text{ of correctly predicted TMHs}}{\# \text{ of predicted TMHs}}$	Precision of TMH prediction
Recall (%)	$100 * \frac{\# \text{ of correctly predicted TMHs}}{\# \text{ of observed TMHs}}$	Recall of TMH prediction
$Q_{ok}$ (%)	$\frac{100}{N} * \sum_{i=1}^N x_i; x_i = \begin{cases} 1, & \text{if } p_i = r_i = 100\% \\ 0, & \text{else} \end{cases}$	Percentage of TMPs with correct TMH placement
$Q_{top}$ (%)	$\frac{100}{N} * \sum_{i=1}^N y_i; y_i = \begin{cases} 1, & \text{if } p_i = r_i = t_i = 100\% \\ 0, & \text{else} \end{cases}$	Percentage of TMPs with correct TMH placement and inside/outside topology
FPR (%)	$100 * \frac{\# \text{ of incorrectly predicted TMPs}}{\# \text{ of soluble proteins}}$	False positive rate of TMP prediction
Sensitivity (%)	$100 * \frac{\# \text{ of correctly predicted TMPs}}{\# \text{ of observed TMPs}}$	Sensitivity of TMP prediction

Listed are the evaluations measures used and how they were calculated. Precision and recall for the performance evaluation of the TMH prediction were computed by combining all TMHs within the dataset (i.e., not averaged over each protein).  $Q_{ok}$  and  $Q_{top}$  were calculated based on all TMPs, where  $N$  was the number of TMPs in the dataset,  $p_i$  and  $r_i$  were the TMH precision and recall for protein  $i$  within the dataset, and  $t_i = 100\%$  indicated a correctly predicted N-terminal inside/outside topology for protein  $i$ .

## Dataset New12

Our original datasets had been based on the PDB release from July 2013, when this work began. Shortly before submission of the work in February 2016, that is, 32 months later, we retrieved all TMPs added to OPM and PDBTM since July 2013. We removed all TMPs similar ( $HVAL > 0$ ) to proteins in datasets used previously (TMP166 and SP1441). Testing the pairwise similarity of the remaining TMPs we found that two pairs were similar ( $HVAL > 0$ ), but we decided to keep them due to their low HVAL. This resulted in 12 new TMPs (New12 dataset, Supporting Information Table S3) we used for additional testing. Although the statistical power of such a small set is very limited, these 12 constitute the entire addition of completely new structures from 2013/07 to 2016/02. Further, these or structurally related TMPs have most likely not been used to develop any method used for comparison.

## Evaluation

As per-protein scores (correct classification as TMP or non-TMP), we compiled the sensitivity (percentage of observed TMPs predicted as TMPs) and the false positive rate (FPR: percentage of soluble proteins predicted as TMPs, Table I). As per-TMH scores (correct identification and placement of TMHs), we compiled the precision (percentage of predicted TMHs that are correct), recall (percentage of observed TMHs predicted as TMHs),  $Q_{ok}$  and  $Q_{top}$ .  $Q_{ok}$  is the percentage of TMPs for which all TMHs are correctly predicted (Table I).  $Q_{top}$  requires in addition to  $Q_{ok}$  correct topology predictions (in/out: Table I). To resolve conflicts between OPM and PDBTM annotations, we chose whichever fit the

prediction best. Note that while sensitivity and recall have the same formula, we used sensitivity in conjunction with TMPs and recall with TMHs to better distinguish between those scores in the text.

Each TMH was considered correctly predicted, if predicted and observed TMH ends were within five residues (Supporting Information Fig. S3), and if predicted and observed TMH overlapped by at least half of the length of the longer of the two helices. These two criteria are more stringent than those that have commonly been used (typically: overlap  $>3-5$  residues anywhere between observed and predicted TMH<sup>36</sup>) and have recently led to re-evaluating TMH prediction methods.<sup>27</sup> None of our major conclusions changed upon applying values slightly different than five residues for the maximum allowed discrepancy between predicted and observed TMH ends (data not shown).

Error rates for the evaluation measures were estimated by bootstrapping,<sup>37</sup> i.e., by re-sampling the population of proteins used for the evaluation 1000 times and calculating the sample standard deviation. Each of these sample populations contained 60% of the original proteins (picked randomly without replacement).

## State-of-the-art methods

We compared TMSEG to the best methods,<sup>27</sup> namely to PolyPhobius,<sup>22</sup> MEMSAT3,<sup>25</sup> and MEMSAT-SVM.<sup>26</sup> Like TMSEG, these methods also use evolutionary information to predict TMPs: MEMSAT3 and MEMSAT-SVM automatically generate position-specific scoring matrices (PSSMs) with PSI-BLAST, while PolyPhobius generates multiple sequence alignments (MSAs). To ensure equal conditions for all methods we ran them on our local machines and used the UniProt Reference Cluster with

90% sequence identity (UniRef90, release 2015\_03) as the homology search database, i.e., to generate the MSAs or PSSMs. While we used proteins completely unknown to TMSEG to assess its performance, some of the proteins used in our assessment might have been used to develop PolyPhobius, MEMSAT3, or MEMSAT-SVM. In this sense, our assessment was likely to over-estimate their performance, in particular with respect to TMSEG.

### Baseline performance

We also compared all methods to a simple baseline predictor similar to TopPred<sup>18</sup>: for all possible segments of 21 consecutive residues, we summed the Eisenberg-hydrophobicity<sup>38</sup> (EisenbergSum, Supporting Information Table S4). All non-overlapping segments with EisenbergSum  $\geq 4$  were predicted as TMHs, starting with the segments with the highest sum. The inside/outside topology was predicted based on the difference between arginine and lysine residues on either side of the TMHs, i.e., applying Gunnar von Heijne's positive-inside rule.<sup>6,7</sup>

### TMSEG input/output

TMSEG needs two input files to successfully run a prediction: a FASTA file with the protein sequence and a PSI-BLAST PSSM file for the input protein. The PSSM file is mandatory and used to include homology-based features that greatly increase the prediction accuracy.

Combining evolutionary information (e.g., PSSMs and MSAs) with machine learning has been the most important improvement in protein prediction and is commonly used in TMH and secondary structure prediction.<sup>24,27,39,40</sup> TMSEG incorporates evolutionary information through PSI-BLAST profiles<sup>41</sup> generated from UniRef90 (release 2015\_03). We used two sets of profiles: a training set with a stringent E-value cutoff of  $10^{-5}$  and five iterations for creating the profile, as well as a test set with a less strict E-value cutoff of  $10^{-3}$  and three iterations. We deactivated PSI-BLAST's low-complexity filter and enabled the option to calculate local optimal Smith-Waterman alignments in order to generate longer and more accurate alignments.

In addition, we used biophysical properties (charge, hydrophobicity, polarity; Supporting Information Table S4) and the overall amino acid composition. These features were calculated twice for each residue: once for all substitutions with a positive PSSM score and once based on all substitutions with a negative score.

The standard output gives a brief summary of the positions of the TMHs and signal peptide (if any) and the inside/outside topology. In addition, a raw output is available that also contains the unmodified output probabilities of the machine-learning tools.

### TMSEG algorithm

TMSEG combines several machine-learning tools and empirical filters. The machine-learning algorithms used are two random forests (RFs) and one neural network (NN), both of which are implementations from the WEKA Java package.<sup>42</sup> The output of these algorithms is further processed with empirically determined filters and thresholds. The TMSEG algorithm executes four separate steps (Fig. 1):

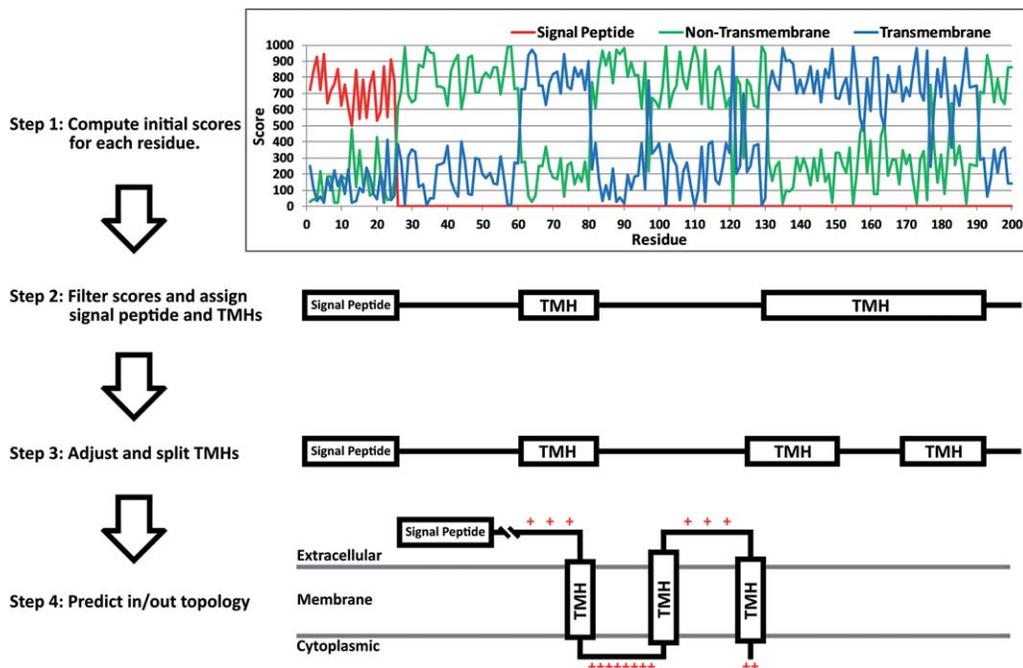
#### Step 1: initial per-residue prediction

An RF detects TMHs from the input sequence. This RF slides a window of 19 consecutive residues through the protein sequence, predicting whether or not the central residue in the window is in a TMH, signal peptide, or non-TM region, i.e., the probability of each residue for each state is calculated based on the residue itself and the nine residues left and right of it. For each of the 19 residue positions, we compute the PSSM profile. For the central nine residues in the window, we also compute the average Kyte-Doolittle<sup>43</sup> hydrophobicity, and the percentage of hydrophobic, charged, and polar residues (Supporting Information Table S4).

In addition to these local features, we compile global features: the distance of the residue to the N- and C-terminus, the length of the protein sequence, and the global amino acid composition. The RF assigns three values to each residue corresponding to the probability to be in a TMH, a signal peptide, or a non-TM region. Runtime is decreased by multiplication of the probabilities by 1000 and transformation into integers.

#### Step 2: per-protein filter: TMP or soluble

The per-residue scores are filtered empirically. First to reduce short peaks of one or two residues, all per-residue scores are smoothed by compiling the median score over five consecutive residues and assigning it to the center residue. Next, each residue is assigned to the state with the highest score (TMH, signal peptide, or non-TM). To prevent over-prediction due to the under-sampling of signal peptide residues, we applied a penalty of 185 (that is, 18.5%) to non-TM and 60 (that is, 6%) to TMH residues. These penalties were optimized during cross-training to best balance over- and under-prediction. Finally, TMHs shorter than seven residues are changed into non-TM regions. If a signal peptide of at least four consecutive residues is identified within the first 40 N-terminal residues ending in residue at position  $i$ , TMSEG predicts a signal peptide from residue 1 to residue  $i$  ( $i \leq 40$ ). Signal peptide predictions outside the first 40 residues ( $i > 40$ ) are changed into non-TM, but do not invalidate signal peptides inside the first 40 residues.

**Figure 1**

TMSEG algorithm. The new method TMSEG has four steps of machine learning and optimization. Step 1: A random forest (RF) assigns a score to each residue for the three states transmembrane helix (TMH), signal peptide, and non-TM region. Step 2: The previous scores are smoothed (median over 5 residues), all residues are assigned to the state with the highest score, and short segments are removed. Step 3: A segment-based neural network (NN) adjusts the exact position of predicted TMHs, and their length, sometimes splitting TMHs, sometimes shifting, extending, or compressing them. Step 4: The inside/outside topology is predicted by another RF.

Initial predictions with fewer than four consecutive residues are changed into non-TM.

### Step 3: refinement of TMHs

In the third step an NN corrects the predicted TMHs. In contrast to the standard sliding window approach of the RF in Step 1, here we introduced a segment-based solution that used as input the following averages over the predicted TMHs: length of predicted TMH, amino acid composition, average hydrophobicity, as well as the percentages of hydrophobic and charged residues. The output of the NN is the predicted probability for the segment to be a TMH. Based on this probability, the predicted TMHs from Step 2 are adjusted.

First, TMHs  $\geq 35$  residues are split into two TMHs with at least 17 residues, if these two TMHs increase the overall probability. The minimum length of 35 residues for splitting long TMHs and of 17 residues for the resulting two TMHs were empirically chosen based on the overall performance during cross-training. Second, the start and end positions for each TMH are adjusted by shifting them by up to three residues in either direction. Shifts are accepted if they increase the overall probability. The maximum endpoint adjustment by three residues was empirically chosen based on the overall

performance during cross-training. In addition, the relatively long minimum TMH lengths to allow splitting and the relatively small shift of maximally three residues of the TMH ends allow TMSEG to maintain a short runtime.

### Step 4: topology prediction

Another RF predicts the inside/outside topology of the TMP, i.e., in which direction the TMHs cross the membrane. During this step the non-transmembrane regions are assigned to inside (e.g., cytoplasmic side of the membrane) or outside. This prediction is made for the entire protein. For each TMH, we consider up to 15 residues before and after the TMH, and eight residues at the TMH start and end (for TMHs  $< 16$  these residues overlap). As all predicted TMHs are assumed to cross the membrane, the in/out assignment is switched after each TMH. For each side, we compute as input to the RF the amino acid composition, the percentage of positively charged residues (we consider all arginine and lysine residues), and the absolute difference of positively charged residues between the two sides. Based on the RF output, one side is assigned to be inside (e.g., cytoplasmic), the other to be outside. Residues immediately after predicted signal peptides are assigned to outside (non-cytoplasmic) and all

**Table II**  
Per-Protein Distinction Between Helical TMPs and Other Proteins

Method	TMP sensitivity	TMP FPR	Topology correct	Misclassified in human	More mistakes than TMSEG in human
TMSEG	98 ± 2	3 ± 1	93 ± 4	558	-
PolyPhobius <sup>22</sup>	100 ± 0	5 ± 1	78 ± 7	770	212
MEMSAT3 <sup>25</sup>	100 ± 0	28 ± 2	93 ± 4	4313	3755
MEMSAT-SVM <sup>26</sup>	98 ± 2	14 ± 2	88 ± 5	2253	1695
Baseline	95 ± 3	31 ± 2	75 ± 7	5015	4457

Results are provided for all 41 TMPs and 285 soluble proteins in the BlindTest dataset. Error rates are the sample standard deviation based on bootstrapping (cf. Methods). Listed are the *TMP sensitivity* (percentage of correctly predicted helical TMPs), the *TMP FPR* (percentage of non-TMP proteins incorrectly predicted as TMP), *Topology correct* (percentage of proteins for which the topology (inside/outside) was correctly predicted; this differs from  $Q_{top}$  which requires topology and all TMHs to be predicted correctly), *Misclassified in human* (estimates the number of proteins misclassified for the entire human proteome), and *More mistakes than TMSEG in human* (estimates the number of proteins misclassified more by the method than by TMSEG). The estimates for the human proteome are based on two assumptions: (i) the error estimates on the BlindTest dataset hold true for the human proteome, (ii) the human proteome has 20,196 proteins, 4791 of which are TMPs (cf. Results section “Application to the human proteome”).

consecutive segments are assigned accordingly without any further prediction.

### TMSEG training

To reduce the risk of over-fitting, we split our combined TMP166 and SP1441 datasets into four even splits (cf. Supporting Information Tables S1 and S2). Note that the TMPs from the SP1441 dataset were used to train the random forest in the initial prediction (step 1) as they contain signal peptide annotations. They are, however, not used for the neural network (step 3) or the random forest in step 4, since they have no inside/outside topology annotations and many of their TMH annotations are not supported by experimental evidence.

The first of three splits was used to train, the second to cross-train, i.e., to optimize all other free parameters (e.g., the minimum TMH length), and the last to evaluate performance (test). This procedure was repeated three times, such that each protein had been used exactly once for training, cross-training and testing. The final parameters were frozen according to the overall best performance for all three rotations (on the test set). Given the frozen parameters, we applied the final method to the fourth split, the BlindTest dataset, which had not been used before.

Our careful four-fold split leading to three-fold development (each with training, cross-training, and testing), provided a double protection against overestimating performance. We decided about every detail in the final method before using the BlindTest dataset to evaluate TMSEG as presented here. Many developers use a two-fold split (training/testing), more careful ones the three-fold split (training/cross-training/testing), while the fourth split is occasionally introduced through pre-release data<sup>39</sup> like the New12 dataset that we generated.

## RESULTS AND DISCUSSION

The novel TMSEG method introduced here distinguishes between proteins with transmembrane helices

(TMHs) and soluble proteins. For all helical transmembrane proteins (TMPs), it predicts the placement of the TMHs, and their orientation in the membrane, i.e., their inside/outside topology. We established sustained performance through cross-validation with two levels of blind testing. We compared our new methods to others, including the best at predicting TMPs,<sup>27</sup> namely PolyPhobius<sup>22</sup> and MEMSAT-SVM.<sup>26</sup> Furthermore, we analyzed MEMSAT3<sup>25</sup> because it excels at the inside/outside topology prediction,<sup>44</sup> and SignalP4.1 as the leading method for signal peptide identification.<sup>35</sup> In addition, we compared to a simple hydrophobicity-based prediction similar to TopPred.<sup>18</sup>

### Outstanding per-protein distinction between TMPs and other proteins

TMSEG correctly identified 40 of the 41 TMPs in the BlindTest dataset (98 ± 2% sensitivity) and incorrectly predicted 8 of 285 soluble proteins as TMPs (3 ± 1% false positive rate: FPR). TMSEG performed similar to PolyPhobius (100% sensitivity and 5 ± 1% FPR) and significantly better than MEMSAT3 and MEMSAT-SVM (Table II).

Although signal peptides can be confused with TMHs due to the similarity of their signal, only one of the 8 mistakes of predicting soluble proteins as TMPs originated from incorrectly predicting a signal peptide as a TMH. This shows that training on a dataset containing signal peptides helped significantly to reduce false positive predictions. PolyPhobius, which also includes a sophisticated signal peptide prediction, did not confuse any signal peptides with TMHs. However, MEMSAT-SVM, MEMSAT3, and the Baseline predictor had 13, 41, and 69 predicted TMHs, respectively, that overlapped by at least half their length with annotated signal peptides. Overall, TMSEG was able to reliably detect signal peptides and to not predict them as TMHs (Supporting Information Table S5).

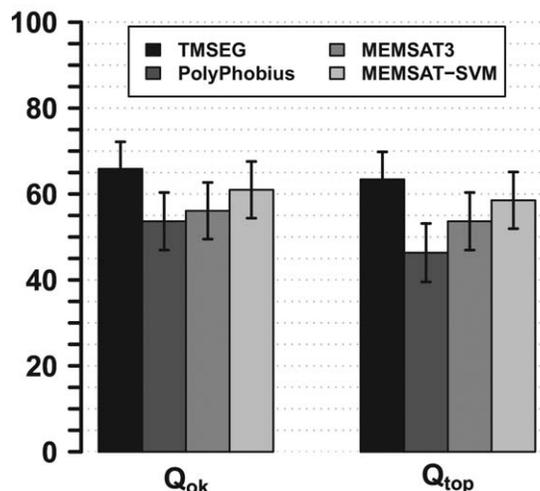
We used the 74 TMPs from the fourth subset of the SP1441 dataset (cf. Supporting Information Table S2) to further test the prediction of signal peptides and TMHs. For these proteins, TMSEG and PolyPhobius incorrectly predicted several single-pass TMPs as soluble proteins, because they confused their TMHs near the N-terminus with signal peptides (Supporting Information Table S5). This trend did not occur with the TMPs from the TMP166 dataset (evident by their high sensitivity values; Table II). An explanation might be that TMPs with TMHs within the first 40 residues are more prevalent in the SP1441 dataset, which makes this misclassification more likely to happen. Although these misclassification rates would lower our previous sensitivity estimates for TMSEG and PolyPhobius (at least for single-pass TMPs with their TMH near the N-terminus), we hesitate to generalize the results to everyday applicability since the SP1441 dataset is biased (it was generated to develop the signal peptide predictor SignalP4.1) and contains many TMPs with a TMH near the N-terminus. Further, only 2 of the 9 TMHs that were incorrectly predicted as SPs had experimental evidence.

While all methods reached high sensitivity, they differed vastly in their false positive rates, i.e., soluble proteins incorrectly considered to contain TMHs (Table II). By translating the error rates, the number of proteins that would be misclassified in the entire human proteome can be estimated using two reasonable assumptions: (i) the error estimates for all methods based on the 326 non-redundant proteins (41 TMPs and 285 soluble proteins) in the BlindTest dataset hold true for the (redundant) human proteome, (ii) the human proteome has 20,196 proteins and 4791 of those are TMPs (cf. Section below “Application to the human proteome”). Under these assumptions, TMSEG achieves 97% per-protein accuracy and misclassifies only about 558 human proteins. The second best method, PolyPhobius, makes 770 mistakes (212 more than TMSEG) and MEMSAT-SVM as the third best method already misclassifies 2253 proteins (1695 more than TMSEG, Table II). In fact, TMSEG is almost 8.8-times superior to the Baseline predictor, PolyPhobius over 6.5-times better, and MEMSAT-SVM 2.2-times better than the Baseline predictor (Supporting Information Table S6).

### Best overall per-TMH prediction

Overall, TMSEG achieved a sustained level of precision ( $87 \pm 3\%$ ) and recall ( $84 \pm 3\%$ ) for the TMHs, that is,  $87 \pm 3\%$  of all predicted TMHs were at the correct position and  $84 \pm 3\%$  of all observed TMHs had been accurately predicted [Supporting Information Fig. S4(A,B)]. These values were second to no other method, however, only slightly above the second best method MEMSAT-SVM ( $85 \pm 3\%$  precision at  $83 \pm 3\%$  recall). All other methods had scores below 80%. For  $66 \pm 6\%$  of all TMPs, TMSEG predicted all observed TMHs at their

## Performance on BlindTest



**Figure 2**

TMSEG compared favorably to state-of-the-art. Results are provided for all 41 TMPs in the BlindTest dataset. Error bars are the sample standard deviation based on bootstrapping (cf. Methods). Shown is on the left the percentage of proteins for which all TMHs were predicted correctly ( $Q_{0k}$ , Table I) and on the right the percentage of proteins with correctly predicted TMHs and inside/outside topology ( $Q_{top}$ , Table I; note that  $Q_{0k} \geq Q_{top}$  by definition).

correct positions, i.e.,  $Q_{0k} = 66 \pm 6\%$  (Fig. 2). MEMSAT-SVM followed as second best with  $Q_{0k} = 61 \pm 7\%$  (Fig. 2). Nevertheless, given the small datasets, the top performance of TMSEG remained within one standard deviation of all compared methods, except the baseline hydrophobicity prediction (Fig. 2: error bars).

When comparing the performance on TMP subsets based on the number of TMHs, the performance got worse the more TMHs a protein had [Supporting Information Fig. S4(C,D)]. This might be misunderstood to imply that prediction methods perform better in placing the TMHs in single-pass TMPs than in, e.g., GPCRs (with 7 TMHs). However, this simple numerical comparison ignores the difference in the difficulty of the task: The Baseline predictor reached a high value in  $Q_{0k}$  for single-pass TMPs, but failed to predict all TMHs correctly for any TMP with  $>5$  TMHs [Supporting Information Fig. S4(C)]. In fact, when we simply compiled performance for the subset of proteins for which the Baseline predictor failed, we found similar values for proteins with one TMH, those with 2–5, and those with  $>5$  TMHs (Supporting Information Fig. S5).

In contrast, it surprised us that even for the trivial cases, i.e., those for which the Baseline predictor had all TMHs correct, the more advanced methods failed for some of them. This suggests that the large number of different features used by the more advanced methods sometimes interfere with and obscure a strong

hydrophobicity signal. Indeed, only 11 of the 19 trivial TMPs were correctly predicted by all four other methods. However, TMSEG still performed best with  $Q_{ok} = 89 \pm 6\%$ , followed by MEMSAT3 and MEMSAT-SVM with  $Q_{ok} = 84 \pm 7\%$  (data not shown).

### Best inside/outside topology prediction

TMSEG and MEMSAT3 correctly placed the N-terminus as inside (e.g., cytoplasmic) or outside (e.g., extracellular), i.e., correctly predicted the topology, for  $93 \pm 4\%$  of all TMPs (Table II). When taking into account the global topology and correct TMH placement (i.e.,  $Q_{top}$ ), TMSEG performed better than all other methods reaching  $Q_{top} = 63 \pm 6\%$  (Fig. 2). This is five percentage points higher than the second best method, MEMSAT-SVM (albeit still within one standard deviation). Most advanced methods predicted the topology correctly for almost all proteins for which they correctly predicted all TMHs ( $Q_{top}$  almost identical to  $Q_{ok}$  for all methods, except for the Baseline predictor in Fig. 2).

### Application to the human proteome

We applied TMSEG to predict all helical TMPs in the human proteome (20,196 proteins from UniProtKB/Swiss-Prot). TMSEG predicted a total of 5157 TMPs, almost half of these (2300 = 45%) were predicted with one TMH. Given the sensitivity and false positive rate of TMSEG ( $98 \pm 2$  and  $3 \pm 1\%$ , respectively; Table II), we estimate that 462 TMPs were incorrectly predicted (over-predicted) and 96 were missed (under-predicted). In total, we thus misclassified 558 proteins, and our corrected estimate was that humans have about 4791 TMPs, i.e., about 24% of all proteins cross the membrane. While TMSEG misclassified about 558 human proteins, the mistake in the estimate of this percentage appeared to be less than a per mille, that is,  $\pm 0.01\%$ . However, our error estimate might be too simplistic due to the high number of single-pass TMPs for which the error rates are much higher than for proteins with more TMPs.

Confirming previous observations,<sup>2,3</sup> we also observed two peaks of predicted TMPs for proteins with 7 TMHs (819 proteins) and 12 TMHs (189 proteins). These likely represent G protein-coupled receptors (GPCRs) and transporter proteins. Applying UniqueProt to the 5157 predicted TMPs, we found around 500 non-redundant TMPs of which 320 are single-pass TMPs.

### Latest experimental structures confirmed our estimates

The 12 new TMPs (New12 dataset) that have recently been added to the PDB constituted the only dataset with truly identical conditions for all methods assessed. The New12 dataset allowed us to confirm the outstanding performance of our new method TMSEG. TMSEG and

PolyPhobius correctly identified 10 of the 12 TMPs ( $83 \pm 10\%$  sensitivity), while MEMSAT3, MEMSAT-SVM, and the Baseline predictor identified 11 ( $92 \pm 7\%$  sensitivity). However, TMSEG correctly predicted every TMH of those 10 TMPs, resulting in a  $Q_{ok} = 83 \pm 10\%$ , compared to  $Q_{ok} = 58 \pm 13\%$  for PolyPhobius, MEMSAT3, and MEMSAT-SVM (Baseline predictor  $Q_{ok} = 50 \pm 13\%$ ). TMSEG also performed best taking into account the topology prediction and reached  $Q_{top} = 66 \pm 12\%$ , compared to a  $Q_{top} = 58 \pm 13\%$  for MEMSAT3 and MEMSAT-SVM, and  $Q_{top} = 50 \pm 13\%$  for PolyPhobius and the Baseline predictor.

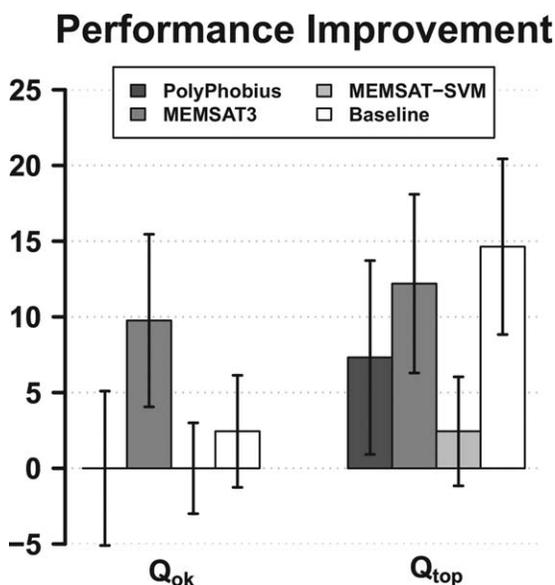
### Comparisons complicated by small datasets

The two small datasets available for evaluation (BlindTest with 41 TMPs and New12 with 12 TMPs) implied high standard errors for many performance estimates. Especially standard errors for the TMH-segment based scores are so high (up to 16 percentage points, Supporting Information Fig. S4) that comparisons between methods hardly provide statistically significant differences on the TMH-segment level. Nevertheless, TMSEG seemed to perform on par with any existing method. Note that the differences in the distinction between helical TMPs and other proteins in the BlindTest dataset were statistically significant even in considering TMSEG as slightly better than the second best PolyPhobius (Table II).

Further, we could not use a single gold standard, because OPM and PDBTM differed in their TMH annotations: comparing the OPM annotations to the PDBTM annotations (that is, “predicting” one with the other) yielded  $Q_{ok} = 56 \pm 7\%$ . In other words, if we considered one of those experiment-based annotations as the prediction of the other, the average performance would be similar to that of TMSEG and the other methods. When using only OPM or PDBTM annotations to evaluate the prediction performance, TMSEG still performed excellently (Supporting Information Fig. S6). However, this was also the only comparison in which one other method reached a numerically higher value for a dataset than TMSEG, namely MEMSAT-SVM on the PDBTM annotations. Overall, all predictions agreed more with OPM than with PDBTM annotations (Supporting Information Fig. S6).

### Performance best with diverse alignments

TMSEG strongly depends on the evolutionary information taken from PSI-BLAST PSSMs. We recommend using a sufficiently large search database (e.g., UniRef90) to generate the PSSMs. Additionally, redundancy reduction might help (e.g., at 90% pairwise sequence identity as in UniRef90).



**Figure 3**

TMSEG applied to refine other methods. The TMSEG algorithm iteratively refines performance through four consecutive steps. Here, we applied Steps 3 and 4 as post-filters to other methods (dataset and error bars as in Fig. 2). Given is the improvement of  $Q_{ok}$  and  $Q_{top}$  (cf. Table I for definitions) of the prediction method by applying TMSEG, i.e.,  $Q(\text{method} + \text{TMSEG}) - Q(\text{method})$ . Note that PolyPhobius (first bar on the left) and MEMSAT-SVM (third bar on the left) showed, on average, no improvement in  $Q_{ok}$ .

Alignments built from smaller search-databases (e.g., UniRef50 and Swiss-Prot) only slightly lowered the per-protein performance: the sensitivity never dropped below  $90 \pm 4\%$ , while the false positive rate remained at or below  $3 \pm 1\%$ . However, the TMH-based precision and recall values dropped substantially (Supporting Information Fig. S7). Thus, for sequences that produce no PSI-BLAST hits, we recommend using a larger search database or—in the rare case that the protein is a true singleton—a method that is independent of evolutionary information, e.g., Phobius.<sup>21,27</sup>

### Re-entrant membrane helices not predicted correctly

Our dataset contained only few re-entrant helices, insufficient to learn their prediction (Supporting Information Table S1). Therefore, we considered re-entrant helices as non-TM during training to avoid later interference with the inside/outside topology prediction. Due to the lack of data, we could not reliably assess how well TMSEG distinguishes TMHs from re-entrant membrane helices: The BlindTest dataset included only seven re-entrant regions (OPM and PDBTM annotations combined). TMSEG incorrectly predicted five of seven as TMHs; two of these five were predicted as two separate TMHs; thus, the overall inside/outside topology was not

influenced. MEMSAT-SVM, the only tested method that predicts re-entrant helices, identified five of the seven as re-entrant, predicted one as a TMH, and missed the last. When considering re-entrant regions as TMHs,  $Q_{ok}$  remained the same for TMSEG and PolyPhobius and dropped by 2–5 percentage points for MEMSAT-SVM, MEMSAT3, and the Baseline predictor.

### TMSEG easily combined with other methods

Due to the modularity of TMSEG (i.e., its four separate steps, Fig. 1), it can be used to refine other methods. This includes the adjustment of the TMHs as well as the inside/outside topology prediction. We used the TMH predictions of the reference methods, and applied Steps 3 and 4 of TMSEG to their prediction (Fig. 2). Applying TMSEG as refinement improved the performance for most methods (Fig. 3; Supporting Information Fig. S8). While the improvement was small for the TMH placement ( $Q_{ok}$ ), TMSEG improved most methods by over eight percentage points in  $Q_{top}$  (correct TMHs and topology).

### Runtime estimation

We estimated the runtime by applying TMSEG to the human proteome (20,196 proteins). As the time to run PSI-BLAST differs depending on the database size, we decided to use pre-computed PSSMs to measure only the time needed by TMSEG. Given those PSI-BLAST profiles, the prediction for the entire human proteome took about 90 min (Intel Core i7-3632QM 2.2 GHz, 8GB RAM; no multithreading), which corresponds to three to four protein sequences per second.

## CONCLUSION

In our hands, our new method TMSEG almost always outperformed existing state-of-the-art prediction methods (Table II, Fig. 2). However, due to the small datasets, many improvements on the per-TMH level remained too small for the large margin of statistical significance (standard errors up to 16 percentage points, Supporting Information Fig. S4). Most importantly, TMSEG achieved the significantly best per-protein classification in the distinction between helical TMPs and all other proteins. For instance, for the prediction of all human proteins, this implied about 558 incorrectly predicted proteins. This number might appear high; however, no method tested reached such a low level, e.g., PolyPhobius misclassified about 200 more proteins than TMSEG and MEMSAT-SVM fared about four times worse (corresponding to >2000 incorrect predictions).

The highest per-protein performance resulted from a combined prediction of TMHs, non-TM regions, and signal peptides. In order to predict re-entrant helices,

another state would have to be introduced; as is, TMSEG predicted five of seven re-entrant helices in our dataset as TMHs. The sustained high levels of per-segment predictions resulted from our new segment-focused algorithm. Another major advantage of our new concept is that it can be used to improve the predictions of most other TMH prediction methods.

### Availability and speed

Other than its top performance, using TMSEG may also be recommended due to its speed and because it might help to improve over the method that you run locally. The method is easily and freely available: online through the PredictProtein<sup>45</sup> webserver ([www.predictprotein.org](http://www.predictprotein.org)), and as standalone Debian package from the Rostlab Debian repository ([www.rostlab.org/owiki](http://www.rostlab.org/owiki)) and GitHub ([www.github.com/Rostlab/TMSEG](http://www.github.com/Rostlab/TMSEG)). A tutorial on how to use PSI-BLAST and TMSEG can be found in the Rostlab Wiki ([www.rostlab.org/owiki/index.php/TMSEG](http://www.rostlab.org/owiki/index.php/TMSEG)).

### ACKNOWLEDGMENTS

Thanks to Tim Karl for technical and to Inga Weise (both TUM) for administrative assistance. Thanks to all authors who made their methods openly available and provided us with versions to run on our own machines. Last but not least, thanks to all who practice open science and deposit their data into public databases and those who maintain these excellent databases.

### REFERENCES

1. von Heijne G. The membrane protein universe: what's out there and why bother? *J Intern Med* 2007;261:543–557.
2. Liu J, Rost B. Comparing function and structure between entire proteomes. *Protein Sci* 2001;10:1970–1979.
3. Fagerberg L, Jonasson K, von Heijne G, Uhlen M, Berglund L. Prediction of the human membrane proteome. *Proteomics* 2010;10:1141–1149.
4. Overington JP, Al-Lazikani B, Hopkins AL. How many drug targets are there? *Nat Rev Drug Discov* 2006;5:993–996.
5. Bigelow HR, Petrey DS, Liu J, Przybylski D, Rost B. Predicting transmembrane beta-barrels in proteomes. *Nucleic Acids Res* 2004;32:2566–2577.
6. von Heijne G. The distribution of positively charged residues in bacterial inner membrane proteins correlates with the trans-membrane topology. *EMBO J* 1986;5:3021–3027.
7. von Heijne G, Gavel Y. Topogenic signals in integral membrane proteins. *Eur J Biochem* 1988;174:671–678.
8. Punta M, Love J, Handelman S, Hunt JE, Shapiro L, Hendrickson WA, Rost B. Structural genomics target selection for the New York consortium on membrane protein structure. *J Struct Funct Genomics* 2009;10:255–268.
9. Love J, Mancía F, Shapiro L, Punta M, Rost B, Girvin M, Wang DN, Zhou M, Hunt JE, Szyperski T, Gouaux E, MacKinnon R, McDermott A, Honig B, Inouye M, Montelione G, Hendrickson WA. The New York Consortium on Membrane Protein Structure (NYCOMPS): a high-throughput platform for structural genomics of integral membrane proteins. *J Struct Funct Genomics* 2010;11:191–199.
10. Caffrey M. A comprehensive review of the lipid cubic phase or in meso method for crystallizing membrane and soluble proteins and complexes. *Acta Crystallogr F Struct Biol Commun* 2015;71:3–18.
11. Moraes I, Evans G, Sanchez-Weatherby J, Newstead S, Stewart PD. Membrane protein structure determination - the next generation. *Biochim Biophys Acta* 2014;1838:78–87.
12. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. The Protein Data Bank. *Nucleic Acids Res* 2000;28:235–242.
13. Kloppmann E, Punta M, Rost B. Structural genomics plucks high-hanging membrane proteins. *Curr Opin Struct Biol* 2012;22:326–332.
14. White SH. Biophysical dissection of membrane proteins. *Nature* 2009;459:344–346.
15. White SH. The progress of membrane protein structure determination. *Protein Sci* 2004;13:1948–1949.
16. Lomize MA, Lomize AL, Pogozheva ID, Mosberg HI. OPM: orientations of proteins in membranes database. *Bioinformatics* 2006;22:623–625.
17. Kozma D, Simon I, Tusnady GE. PDBTM: Protein Data Bank of transmembrane proteins after 8 years. *Nucleic Acids Res* 2013;41:D524–529.
18. von Heijne G. Membrane protein structure prediction. Hydrophobicity analysis and the positive-inside rule. *J Mol Biol* 1992;225:487–494.
19. Krogh A, Larsson B, von Heijne G, Sonnhammer EL. Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *J Mol Biol* 2001;305:567–580.
20. Tusnady GE, Simon I. The HMMTOP transmembrane topology prediction server. *Bioinformatics* 2001;17:849–850.
21. Käll L, Krogh A, Sonnhammer EL. A combined transmembrane topology and signal peptide prediction method. *J Mol Biol* 2004;338:1027–1036.
22. Käll L, Krogh A, Sonnhammer EL. An HMM posterior decoder for sequence feature prediction that includes homology information. *Bioinformatics* 2005;21:i251–257.
23. Rost B, Casadio R, Fariselli P, Sander C. Transmembrane helices predicted at 95% accuracy. *Protein Sci* 1995;4:521–533.
24. Rost B, Casadio R, Fariselli P. Refining neural network predictions for helical transmembrane proteins by dynamic programming. *Proc Int Conf Intell Syst Mol Biol* 1996;4:192–200.
25. Jones DT. Improving the accuracy of transmembrane protein topology prediction using evolutionary information. *Bioinformatics* 2007;23:538–544.
26. Nugent T, Jones DT. Transmembrane protein topology prediction using support vector machines. *BMC Bioinformatics* 2009;10:159
27. Reeb J, Kloppmann E, Bernhofer M, Rost B. Evaluation of transmembrane helix predictions in 2014. *Proteins* 2015;83:473–484.
28. UniProt C. UniProt: a hub for protein information. *Nucleic Acids Res* 2015;43:D204–212.
29. Velankar S, Dana JM, Jacobsen J, van Ginkel G, Gane PJ, Luo J, Oldfield TJ, O'Donovan C, Martin MJ, Kleywegt GJ. SIFTS: Structure Integration with Function, Taxonomy and Sequences resource. *Nucleic Acids Res* 2013;41:D483–489.
30. Mika S, Rost B. UniqueProt: Creating representative protein sequence sets. *Nucleic Acids Res* 2003;31:3789–3791.
31. Sander C, Schneider R. Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins* 1991;9:56–68.
32. Rost B. Twilight zone of protein sequence alignments. *Protein Eng* 1999;12:85–94.
33. Granseth E, Viklund H, Elofsson A. ZPRED: predicting the distance to the membrane center for residues in alpha-helical membrane proteins. *Bioinformatics* 2006;22:e191–196.
34. Papaloukas C, Granseth E, Viklund H, Elofsson A. Estimating the length of transmembrane helices using Z-coordinate predictions. *Protein Sci* 2008;17:271–278.
35. Petersen TN, Brunak S, von Heijne G, Nielsen H. SignalP 4.0: discriminating signal peptides from transmembrane regions. *Nat Methods* 2011;8:785–786.

36. Chen CP, Kernytsky A, Rost B. Transmembrane helix predictions revisited. *Protein Sci* 2002;11:2774–2791.
37. Efron B, Tibshirani RJ. *An introduction to the bootstrap*. Chapman & Hall (New York) 1993.
38. Eisenberg D. Three-dimensional structure of membrane and surface proteins. *Annu Rev Biochem* 1984;53:595–623.
39. Rost B, Sander C. Prediction of protein secondary structure at better than 70% accuracy. *Journal of Molecular Biology* 1993;232:584–599.
40. Rost B, Liu J, Nair R, Wrzeszczynski KO, Ofran Y. Automatic prediction of protein function. *Cellular and Molecular Life Sciences* 2003;60:2637–2650.
41. Altschul SE, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 1997;25:3389–3402.
42. Hall M, Frank E, Holmes G, Pfahringer B, Reutemann P, Witten IH. The WEKA Data Mining Software: An Update. *SIGKDD Explorations* 2009;11:10–18.
43. Kyte J, Doolittle RF. A simple method for displaying the hydropathic character of a protein. *J Mol Biol* 1982;157:105–132.
44. Rath EM, Tessier D, Campbell AA, Lee HC, Werner T, Salam NK, Lee LK, Church WB. A benchmark server using high resolution protein structure data, and benchmark results for membrane helix predictions. *BMC Bioinformatics* 2013;14:111
45. Yachdav G, Kloppmann E, Kajan L, Hecht M, Goldberg T, Hamp T, Honigschmid P, Schafferhans A, Roos M, Bernhofer M, Richter L, Ashkenazy H, Punta M, Schlessinger A, Bromberg Y, Schneider R, Vriend G, Sander C, Ben-Tal N, Rost B. PredictProtein—an open resource for online prediction of protein structural and functional features. *Nucleic Acids Res* 2014;42:W337–343.