

Molecular BioSystems

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.



www.rsc.org/molecularbiosystems

1 **GPCRserver: an accurate and novel G protein-coupled receptor**
2 **predictor**

3

4 **Renxiang Yan^{a,*}, Xiaofeng Wang^b, Lanqing Huang^a, Jun Lin^a, Weiwen Cai^a and**
5 **Ziding Zhang^b**

6

7

8

9 *^aInstitute of Applied Genomics, School of Biological Sciences and Engineering,*

10 *Fuzhou University, Fuzhou 350002, China*

11 *^bState Key Laboratory of Agrobiotechnology, College of Biological Sciences, China*

12 *Agricultural University, Beijing 100193, China*

13

14

15

16 *Corresponding author (E: yanrenxiang@fzu.edu.cn; T/F: +86 591 22866273)

17

18 ***Running title: Prediction of GPCRs***

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

1 **ABSTRACT**

2 G protein coupled receptors (GPCRs), also known as seven-transmembrane domain
3 receptors, pass through the cellular membrane seven times and play diverse biological
4 roles in the cells such as signaling, transporting of molecules and cell-cell
5 communication. In this work, we develop a web server, namely GPCRserver, which is
6 capable of identifying GPCRs from genomic sequences, and locating their
7 transmembrane regions. The GPCRserver contains three modules: (1) Trans-GPCR
8 for transmembrane regions prediction by using sequence evolutionary profiles with
9 the assistance of neural network training, (2) SSEA-GPCR for identifying GPCRs
10 from genomic data by using secondary structure element alignment, and (3)
11 PPA-GPCR for identifying GPCRs by using profile-to-profile alignment. Our
12 predictor was strictly benchmarked and showed its favorable performance in the real
13 application. The web server and stand-alone programs are publicly available at
14 <http://genomics.fzu.edu.cn/GPCR/index.html>.

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

1 INTRODUCTION

2 G protein-coupled receptor (GPCR) is a major transmembrane (TM) protein type in
3 the cellular membrane and plays critical roles in a wide variety of biological processes,
4 including homeostasis modulation¹, cell growth² and transporting of small molecules
5³. GPCRs are also important to humans. The human genome encodes thousands of
6 GPCRs⁴, and, moreover, it is estimated that a large number of drugs in the market are
7 designed to regulate the mechanism involved in GPCRs⁵. GPCRs are referred to as
8 seven-TM receptors according to the fact that all existing GPCRs contain seven-TM
9 α -helices with loops connecting them. Determination of the three-dimensional (3D)
10 structures is a direct way to decipher their biological functions. Unfortunately, it is
11 very time-consuming, and requires amazing funding and extensive efforts to obtain
12 crystals of GPCRs. Compared with globular proteins, it is much more difficult to
13 determine 3D structures of GPCRs. Due to experimental difficulties, the existing
14 GPCR structures are very limited. For example, although there are more than 90000
15 protein structures deposited in the PDB database⁶, the existing 3D structures of
16 GPCRs in the PDB are only ~100 at the time of March, 2014, and the non-redundant
17 structures of GPCRs are much fewer. Considering the limitations of GPCR structural
18 determination using wet experiments, it is of great need to develop accurate and
19 high-throughput GPCR prediction methods.

20

21 Currently, there exist two major tasks to the computational study of GPCRs. One is to
22 identify GPCRs from genome-wide sequences; the other is to locate TM regions of
23 GPCR candidates. The low sequence similarities among some GPCRs, especially the
24 existence of orphan GPCRs, hampers their identification by classical
25 sequence-to-sequence alignments, such as BLAST⁷. Thus, the community needs
26 specific GPCR prediction and identification programs. The past two decades have
27 been witnessing exciting advances of a couple of such bioinformatics methods. In
28 general, a sliding window centered at the target residue is excised and fed into the
29 statistical learning algorithms to train the models. As one of the simplest forms, Gao
30 and Chess developed hydrophathy-curve algorithm to detect proteins with seven
31 hydrophobic stretches to screen potential GPCRs⁸. More sophisticated approaches
32 such as hidden Markov model (HMM)⁹ and Support Vector Machine (SVM)¹⁰ are
33 also used in the GPCR prediction. To develop the HMM-based methods, their designs
34 of topologies of the HMMs, number of states and their connection need to be fixed in
35 advance by taking insightful knowledge of known GPCRs. Once the topologies of
36 HMMs are fixed, the protein sequence/structural data are used to train the probability
37 of each transition of the HMMs. Phobius¹¹, TMHMM¹², GPCRHMM¹³ and
38 HMMTOP¹⁴ are hidden Markov model-based methods for GPCR TM region
39 prediction. PRED-GPCR by Papasaikas and his co-workers is a probabilistic method
40 that uses family-specific HMMs to determine to which GPCR family a target
41 sequence belongs¹⁵. Jones group proposed a SVM-based method for TM protein
42 topology prediction¹⁶. Meanwhile, a new set of conformational parameters for TM α
43 helices was developed by Gromiha¹⁷ and the parameters can be used to locate the TM
44 regions of GPCR. GPCRpred is also a SVM-based GPCR identification method by

1 clustering GPCRs into different families¹⁸. The TM region prediction programs can
2 be used for GPCR identification by scanning databases for proteins predicted to have
3 seven-TM helices. GPCR identification and TM region prediction have been widely
4 used in biological research. So for examples, Nowling *et al* screened GPCRs in the
5 genomes of three insect vectors using an ensemble procedure¹⁹; Takeda and his
6 co-workers identified a large number of potential GPCRs when searching human
7 proteome for proteins predicted to contain 6~8 TM helices²⁰. Meanwhile, there are
8 some other bioinformatics studies of GPCRs²¹⁻²³. In general, the performance of
9 statistical learning methods depends on the input features, learning algorithms and
10 optimized parameters. Developers are required to carefully tune the parameters of
11 training algorithms to obtain optimized performance.

12
13 In this work, we develop a predictor, which is capable of accurately identifying
14 GPCRs from genomic sequences as well as predicting their TM segments. The TM
15 regions of GPCRs are predicted by using sequence evolutionary profiles with the
16 assistance of neural network learning. Moreover, considering the secondary structure
17 topologies of GPCRs are conserved, protein secondary structure-based methods for
18 GPCRs identification may make sense and we therefore develop such a method.
19 Meanwhile, a novel profile-to-profile alignment algorithm is also developed to detect
20 GPCRs. As that clearly pointed out by Chou in his review²⁴ as well as that in several
21 closely related studies²⁵⁻²⁷, we can use the following procedure to establish a practical
22 and reliable bioinformatics predictor. Firstly, build a model by using effective
23 mathematical expressions that can truly reflect their intrinsic correlation with the
24 target to be predicted, and then construct or select reliable benchmark datasets to
25 train/test the models. Secondly, objectively evaluate the anticipated accuracy of the
26 new model and compare it with community popular methods. Last but not the least,
27 stand-alone programs and publicly available web servers for the models should be
28 developed to facilitate researchers to use new methods. We will describe the
29 procedure step-by-step in the following sections.

30 31 **2 MATERIALS AND METHODS**

32 **2.1 Datasets**

33 The benchmark datasets were constructed with the utilization of information in the
34 PDB²⁸, and UniProtKB²⁹. Firstly, we downloaded 55 structurally known GPCRs
35 from PDB database with timestamp of October, 2013. This dataset was named
36 GPCR_PDB55. At the same time, the Swiss-Prot²⁹ database of UniProtKB
37 ([ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/complete/
uniprot_sprot.dat.gz](ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/complete/uniprot_sprot.dat.gz)) was also downloaded in our local computers. We scanned all the
38 sequences in the Swiss-Prot database and there are 2222 GPCRs showing high
39 sequence similarity with sequences in the GPCR_PDB55 (BLAST e-value<0.01). We
40 further scanned the Swiss-Prot database and obtained 558 potential GPCRs, which
41 were not similar to the 2222 proteins at the sequence level. Among the 558 proteins,
42 256 ones have already been included in the GPCRDB³⁰ database
43 (<http://www.gpcr.org/7tm/>). Moreover, the remaining 302 proteins, which are
44

1 seven-TM proteins, out of the 558 proteins are probably GPCRs. We found there are
2 some annotations, such as ‘SIMILARITY: Belongs to the G-protein coupled receptor
3 4 family’, ‘SIMILARITY: Belongs to the G-protein coupled receptor Fz/Smo’, ‘DR
4 Pfam; PF10326; 7TM_GPCR_Str; 1’ and so on. Therefore, these 302 proteins are
5 most likely to be GPCRs. Therefore, the 558 proteins were regarded as GPCRs in our
6 benchmark. Further, we randomly selected 721 non-GPCRs membrane proteins from
7 Swiss-Prot database. These datasets were filtered by removing redundancies at 95%
8 sequence identity. Finally, we obtained 1697 train (GPCR_TRAIN1697), 492 test
9 (GPCR_TEST492) GPCRs and 504 non-GPCRs membrane proteins (MEM_504).
10 Details of removing redundancies are available in supplementary file 1. Meanwhile,
11 we collected 2014 non-GPCR proteins, covering 2014 SCOP protein families, from
12 SCOPe³¹ database, and the dataset of the 2014 proteins was named SCOP_2014. We
13 used GPCR_TEST492 to benchmark various GPCR TM location methods. The
14 performance of GPCR identification is assessed by methods’ abilities in classification
15 of GPCR/non-GPCR in the GPCR_TEST492, SCOP_2014 and MEM_504 datasets.
16 The datasets are available at <http://genomics.fzu.edu.cn/GPCR/dataset/>. It should be
17 clearly pointed out that the proteins in the GPCR_TEST492 share low similarity with
18 the proteins of the GPCR_TRAIN1697 dataset at the sequence level (BLAST
19 e-value>0.01).

20
21

22 **2.2 Trans-GPCR for TM region prediction**

23 Trans-GPCR is a neural network based method for TM region prediction. The neural
24 network algorithm used in this work was implemented utilizing Encog Java neural
25 network framework³², which can be downloaded from
26 <https://code.google.com/p/encog-java/>. The standard back propagation³³ and sigmoid
27 activation function were used. We trained the Trans-GPCR using a similar way to
28 PSIPRED³⁴. Briefly, two feed-forward back-propagation neural networks were
29 jointly used. In our work, the first neural network contains two hidden layers, whereas
30 the second neural network only contains a single hidden layer. The nodes in both two
31 hidden layers of the first neural network were set to 250; Meanwhile, the node
32 number in the hidden layer of the second neural network was set to 70. The
33 architectures and parameters of neural networks were optimized using the training
34 dataset. The input features of the first neural network are evolutionary sequence
35 profiles. The outputs of the first neural network are fed into the second neural network
36 that to refine the prediction. To obtain the sequence profiles, the target sequence is
37 iteratively threaded through NCBI³⁵ NR database for three repeats with an e-value
38 cutoff 0.001 for collecting multiple sequence alignments (MSAs) using PSI-BLAST³⁶.
39 The position specific scoring matrix/profile (PSSM) is generated by the option ‘-Q’.
40 The position specific frequency matrix/profile (PSFM) is calculated from the
41 generated MSA using Henikoff weight³⁷. In the Henikoff weight scheme, a residue in
42 each position is assigned a weight equal to $1/(t+s)$, where t is the number of different
43 residues in the column and s is the number of times the particular residue appears in
44 the column. The position-based weights (i.e. Henikoff weights) are then added for

1 each column and divided by the length of sequence. Then, we use following equation
2 to calculate the PSFM profile of each residue from a MSA

$$3 \quad f_{u,r} = \frac{\sum_{i=1}^N w_u^i \delta_{u,r}^i}{\sum_{i=1}^N w_u^i} \quad (1)$$

4 where $f_{u,r}$ is the amino acid frequency of residue r at column u ; N is the number of
5 sequences in the MSA; w_u^i is the Henikoff weight for column u of sequence i ; $\delta_{u,r}^i$ is
6 set to 1 if sequence i has residue r in column u and 0, otherwise. For unaligned
7 regions, only the target sequence itself is used to calculate the amino acid frequencies.

8
9 For each target residue, a sliding window containing $2n+1$ residues long (i.e. window
10 size = $2n+1$) fragment profiles centered at the target residue is excised from the
11 sequence profiles. The optimal window sizes of two neural networks were determined
12 by performance in the training dataset and were set to 21. There are two sets of
13 generated profiles, including PSFM and PSSM profiles. Using a similar way to Chen
14 *et al*³⁸, we also compute the Shannon entropy for each residue as

$$15 \quad Entropy = \sum_{r=1}^{20} -f_{u,r} \log(f_{u,r}) \quad (2)$$

16 where $f_{u,r}$ is calculated using Eq 1; r is the r th residue type. Meanwhile, there are
17 two-dimensional RW-GRMTP (relative weight of gapless real matches to
18 pseudocounts), which are the last two columns in the PSSM profile, of each residue
19 generated by PSI-BLAST. The RW-GRMTP represents the number of aligned
20 residues in that position. The RW-GRMTP information is also used as training
21 features. Considering some elements of the PSSM profile are negatives, we directly
22 scale the values to the range of 0~1 by using the standard logistic function as

$$23 \quad \frac{1}{1+e^{-x}} \quad (3)$$

24 where x is the element value of the PSSM profile. Again, we also compute the entropy
25 score for PSSM profile. For PSSM profile as well as PSFM profile, there are 20
26 residue frequencies and an entropy value. Additionally, an extra unit per amino acid
27 is used to indicate whether the residue spans either the N or C terminus of the protein
28 chain. For a given 21-residue window, input features for the first neural network are
29 window_size_1*(21+21+2+1), where 21 for PSSM, 21 for PSFM, 2 for RW-GRMTP
30 and an additional unit to indicate whether the residue spans either the N or C terminus.
31 The window_size_1 value of 21 is optimized by the performance in the training
32 dataset. Using a similar way to Chou *et al*²⁴, we can denote the input features for
33 position i of a protein as

$$34 \quad \{[PSSM(i+s,j)], [E(i+s)], [PSFM(i+s,j)], [\bar{E}(i+s)], RW-GRMTP(i+s)\}, j \in [0,20], s \in [-n,n] \quad (4)$$

35 where $PSSM(i+s,j)$ is for the scaled PSSM profile at the position $i+s$; j ranges from 0
36 to 19, in which [0,19] representing 20 amino acids and one additional bit that used to
37 indicate whether the residue spans either the N or C terminus of the protein chain; s is

1 a shift value and its value ranges from $-n$ to n (i.e. window size). $E(i+s)$ is the
 2 Shannon entropy for position $i+s$ calculated using scaled PSSM profile at position $i+s$.
 3 Similarly, $PSFM(i+s,j)$ is for the PSFM profile at the position $i+s$; $\bar{E}(i+s)$ is the
 4 Shannon entropy for position $i+s$ calculated using PSFM profile at position $i+s$;
 5 RW-GRMTP($i+s$) is the RW-GRMTP values (i.e. relative weight of gapless real
 6 matches to pseudocounts) at position $i+s$.

7
 8 The feature numbers for the second neural network are $window_size_2*(2+1)$, where
 9 2 denotes the outputs (e.g. prediction scores of TM/non-TM) of the first neural
 10 network and an additional unit to indicate whether the residue spans either the N or C
 11 terminus. The $window_size_2$ value of 21 is optimized using the same way as that of
 12 $window_size_1$. The average length of the TM regions is 22 in our training dataset.
 13 Meanwhile, lengths of the loops connecting the TM segments are diverse. Based on
 14 this observation, we transform the prediction of orphan residues, assigning a TM
 15 (non-TM) residue to non-TM (TM) region if its neighbor six residues (i.e. ± 3
 16 positions) are non-TM (TM).

17 **2.3 GPCR identification**

18 **2.3.1 Trans-GPCR for GPCR identification**

19 Furthermore, Trans-GPCR not only predicts the TM regions of GPCRs but also can
 20 identify GPCRs. For a target sequence, Trans-GPCR determines whether it is GPCR
 21 by the following equation
 22

$$23 \quad TransGPCR_Score = \sum_{i=1}^N \max(NN(M) - NN(-), 0) \quad (5)$$

24 where $NN(M)$ and $NN(-)$ are the TM and non-TM prediction scores of residue i by two
 25 output nodes of the second neural network in Trans-GPCR method; N is length of
 26 target protein. We use $\max(NN(M) - NN(-), 0)$ to ensure that only predicted TM
 27 regions are summed (i.e. positive values). Here, we use a reliable parameter for
 28 position i of target protein as

$$29 \quad residue_reliable(i) = abs(NN(M) - NN(-)) \quad (6)$$

30
 31 where $residue_reliable(i)$ is a reliable index; abs is the absolute mathematic function;
 32 $NN(M)$ and $NN(-)$ are defined in Eq. 5. $residue_reliable(i)$ ranges [0-1], where a
 33 higher score corresponds to a more reliable prediction for residue i . It should be
 34 clearly pointed out that the parameter $TransGPCR_Score$ is to determine whether a
 35 protein is GPCR, whereas $residue_reliable(i)$ is a position-specific reliability index of
 36 prediction for position i of target protein.
 37

38 **2.3.2 SSEA-GPCR for GPCR identification**

39 Here, we also develop a GPCR identification algorithm by using secondary structure
 40 element alignment (SSEA). Since protein secondary structural topologies of GPCRs
 41

1 are more conserved than single sequences, SSEA is therefore able to identify GPCRs.
 2
 3 SSEA-GPCR method searches a target sequence against a GPCR and a non-GPCR
 4 databases. In this process, the top i SSEA similarity scores between GPCRs
 5 (non-GPCRs) are recorded (i.e. $SSEA_{max_gpcr_i}$ and $SSEA_{max_non_gpcr_i}$). In SSEA
 6 algorithm, the secondary structural string for each sequence is converted into
 7 secondary structure elements such that ‘H’ represents a helix element, ‘E’ denotes a
 8 strand element, and ‘C’ stands for a coil element. Meanwhile, the predicted secondary
 9 structural string was shortened and the length of each element was retained for the
 10 scoring of SSEA. Here, Needleman-Wunsch global alignment algorithm³⁹ was used
 11 with the gap penalties set to zeros. The alignment score of SSEA between two
 12 secondary structure elements with lengths L_i and L_j is defined as

$$14 \quad Score(i, j) = \begin{cases} \min(L_i, L_j) & \text{Match between two identical elements} \\ 0.5 \times \min(L_i, L_j) & \text{Match between } \alpha\text{-helix} / \beta\text{-strand and coil} \\ 0 & \text{Match between } \alpha\text{-helix and } \beta\text{-strand} \end{cases} \quad (7)$$

15
 16 where $\min(L_i, L_j)$ stands for the minimal length between L_i and L_j . The normalized
 17 SSEA alignment score is obtained by dividing by the length of the target sequence.
 18 Details of SSEA algorithm can refer to its original developer⁴⁰ or our previous work
 19⁴¹. For a target sequence, the $SSEA_gpcr$ prediction score is calculated using a simple
 20 K -nearest algorithm as

$$21 \quad SSEA_gpcr = \frac{\sum_{i=1}^K SSEA_{gpcr_top_i} - \sum_{i=1}^K SSEA_{non_gpcr_top_i}}{K} \quad (8)$$

22 where $SSEA_{gpcr_top_i}$ and $SSEA_{non_gpcr_top_i}$ are the top i prediction scores of searching
 23 target protein against GPCR and non-GPCR databases; The value of K is primarily
 24 optimized and set to 10. Here, the GPCRs in the training dataset are used as GPCRs
 25 database to calculate $SSEA_{max_gpcr}$. Meanwhile, we collected 3836 non-GPCR proteins,
 26 which cover 1061 folds, 1713 superfamilies and 3836 families, as a non-GPCR
 27 database from SCOPe database³¹. This dataset was named nonGPCRlib_3836. The
 28 nonGPCRlib_3836 is used when calculating $SSEA_{max_non_gpcr}$. The proteins in the
 29 SCOP_2014 dataset, which has been described in the Datasets section, share low
 30 similarity with proteins in the non-GPCR database (i.e. nonGPCRlib_3836) at the
 31 sequence level (BLAST e-value>0.01).

32 33 **2.3.3 PPA-GPCR for GPCR identification**

34 GPCRs constitute a large superfamily of proteins¹³. Therefore, profile-to-profile
 35 alignment, which represents one of useful methods to detect distant homologs, should
 36 be effective to identify potential GPCRs. Similar to SSEA-GPCR method,
 37 Needleman-Wunsch global alignment algorithm is also used and the penalties for
 38 ending gaps are set as zeros. The scoring function for profile-to-profile alignment is as
 39

$$Score(i,j)=PF(i,j)+w_1SS(i,j)+shift \quad (9)$$

where $PF(i,j)$ is an evolutionary profiles-based term. Evolutionary profiles are generated from MSAs, which represent the divergence of proteins in the same family, and contain important information to infer the protein features. The MSAs are obtained using the same way as that in Trans-GPCR. The values of gap opening, gap extension, w_1 and $shift$ were obtained by maximum of the sequence alignments to structural alignments⁴² of all-to-all pair-wises for the 55 structurally known GPCRs in the GPCR_PDB55 dataset. The values of gap opening, gap extension, w_1 and $shift$ were set to -7.1, -0.56, 0.7 and -0.9. The Profile similarity score is as

$$PF(i, j) = \frac{1}{2} \sum_{k=1}^{20} (PSFM(i, k)_q PSSM(j, k)_t + PSFM(j, k)_t PSSM(i, k)_q) \quad (10)$$

where $PSFM(i, k)_q$ represents the frequency of the k th amino acid at the i th position of PSFM profile for target protein; $PSSM(j, k)_t$ denotes the k th amino acid at the j th position of PSSM profile for template. Similarly, $PSSM(j, k)_t$ represents the frequency of the k th amino acid at the j th position of PSFM profile for template; $PSSM(i, k)_q$ denotes the k th amino acid at the i th position of PSSM profile for target protein. In our method, the similarity score for each pair of secondary structure profile columns is defined as the Pearson's correlation coefficient between them as

$$SS(i, j) = \frac{3 \sum_{k=1}^3 Q_{i,k} T_{j,k} - \sum_{k=1}^3 Q_{i,k} \sum_{k=1}^3 T_{j,k}}{\sqrt{3 \sum_{k=1}^3 Q_{i,k}^2 - (\sum_{k=1}^3 Q_{i,k})^2} \sqrt{3 \sum_{k=1}^3 T_{j,k}^2 - (\sum_{k=1}^3 T_{j,k})^2}} \quad (11)$$

where $Q_{i,k}$ is the possibility of k th (i.e. $k=1,2,3$ corresponding to α -helix (H), β -strand (E), and coil (C), respectively) secondary structure type at i th position of the target sequence. $T_{j,k}$ is the possibility of k th secondary structure type at j th position of the template sequence. The prediction possibilities of protein secondary structure are obtained by using PSIPRED. Similar to SSEA-GPCR, the normalized PPA-GPCR alignment score is also obtained by being divided by length of target sequence. Moreover, the estimated significant Zscore of PPA-GPCR alignment scores should be calculated. We use SCOPe_1187 dataset, which is constructed by randomly selecting one protein of each fold from SCOPe database, as a reference database to calculate mean and standard deviation of random scores. The Zscore is calculated as

$$Zscore = \frac{raw - mean}{std} \quad (12)$$

where raw is the alignment score between a target and a specific template; $mean$ and std are the average and standard deviation of scores aligning target sequence to the 1187 proteins in the SCOPe_1187 dataset. There are two Zscores for any pair of target-template alignments. Here, we use a symmetrical Zscore similar to FFAS-3D⁴³ as

$$Zscore(q,t)=ave(Zscore_q,Zscore_t) \quad (13)$$

where $Zscore_q$ and $Zscore_t$ are the Zscores of the target and template proteins by searching SCOPe_1187 database using Eq. 12. Here, we use the average of $Zscore_q$ and $Zscore_t$ as the final value of the calibrated score. Note that $Zscore(q,t)$ is symmetrical with respect to two proteins. We also tested the minimum and maximum of the two Zscores, but the performance cannot be improved. For each target, we search it against a GPCR database, which is GPCR_TRAIN1697 in our benchmark. The maximum $Zscore(q,t)$ of the target and the templates (i.e. 1697 pair-wise alignment scores) in the GPCR_TRAIN1697 database is recorded and is named PPA_gpcr in this paper. Confidence intervals (CI) of PPA_gpcr are computed using the common assumption of a normal distribution by the following as

$$\left[\mu - Z \frac{SD}{\sqrt{n}}, \mu + Z \frac{SD}{\sqrt{n}} \right] \quad (14)$$

where μ and SD are mean and standard deviation of PPA_gpcr scores; n is sample size; Z is the critical value and the value of Z is 1.96 in a 95% confidence level.

2.3.4 Combined methods

The combined methods can be constructed by using complementary algorithms with improved performance. When combining the top four methods (HMMTOP, TMHMM, Phobius and Trans-GPCR) for TM/non-TM region prediction, we use two bits to denote their prediction for each residue (i.e. [1, 0] for TM and [0, 1] for non-TM predictions). To combine the four methods, the corresponding bit values are simply added. For example, [1, 0], [1, 0], [1, 0] and [0, 1] are added and the result is [3, 1]. The combined prediction for a residue is TM if the value of the first bit is bigger than that of the second bit, and non-TM, otherwise. The combined method for TM/non-TM prediction is named TM-Combined in this paper. Similarly, we also combined the methods (Trans-GPCR, SSEA-GPCR and PPA-GPCR) for GPCR identification (Iden-Combined) using a weighted score as

$$Iden-Combined = w_1 PPA_gpcr + w_2 Trans_gpcr + w_3 SSEA_gpcr \quad (15)$$

where $Iden-Combined$ is the combined prediction score; w_1 , w_2 and w_3 are weighted to balance the three terms. Considering the value ranges of the three terms, the values of w_1 , w_2 and w_3 are primarily optimized and set to 0.1, 0.0067 and 1, respectively.

2.3.5 Amino acid distribution of TM/non-TM regions

It is also interesting to mine the amino acid distribution of TM/non-TM regions in the GPCRs. The formula for calculating the composition of i th residue is defined as

$$composition(i) = \frac{\sum_{k=1}^N \delta_k}{N} \quad (16)$$

where i stands for composition of i th residue; δ_k is set to 1 if the position k of

1 sequence is i th residue and 0, otherwise; N is the total number of residues in the
2 TM/non-TM regions.

3 **2.4 Performance assessment**

4 When the benchmark is performed over the test dataset in the TM/non-TM region
5 prediction, the overall performance of different methods is evaluated with respect to
6 four parameters: Accuracy (Ac), Sensitivity (Sn), Specificity (Sp) and Matthew
7 correlation coefficient (Mcc). The TM (non-TM) residues of GPCRs are considered
8 positives (negatives). The equations for these parameters are as follows

$$9 \quad Ac = \frac{tp + tn}{tp + fn + tn + fp} \quad (17)$$

$$10 \quad Sn = \frac{tp}{tp + fn} \quad (18)$$

$$11 \quad Sp = \frac{tn}{tn + fp} \quad (19)$$

$$12 \quad Mcc = \frac{tp \times tn - fp \times fn}{\sqrt{(tp + fp)(tp + fn)(tn + fn)(tn + fp)}} \quad (20)$$

13 where tp , fp , fn and tn are the numbers of true positives, false positives, false
14 negatives and true negatives, respectively. The performance of GPCR identification
15 can be measured by receiver operating characteristic (ROC) curves⁴⁴. The ROC
16 curves plot true-positive rate (instances) as a function of false-positive rate (instances)
17 for all possible thresholds of prediction scores by various methods. The set of four
18 equations (Eqs. 17-20) is used for single-label systems. For multi-label systems,
19 which are more frequent in system biology^{45, 46}, a completely different set of metrics
20 as defined in⁴⁷ is needed.

21 **3 RESULTS AND DISCUSSIONS**

22 **3.1 The performance of TM region prediction**

23 Among the resulting measures, Ac and Mcc are the most comprehensive parameters to
24 assess the prediction performance. The neural network model of Trans-GPCR was
25 intensively trained on the GPCR_TRAIN1697 dataset and generated the results of
26 $Ac=0.940$ and $Mcc=0.877$. Further, the performance of TM region location was tested
27 on the GPCR_TEST492 dataset. HMMTOP, TMHMM, Memast and Phobius
28 programs were installed in our local computers and the proteins were directly fed into
29 them. The prediction results of the TM regions for various methods were summarized
30 in Table 1. HMMTOP, TMHMM, Memast and Phobius generated Ac (Mcc) scores of
31 0.927 (0.804), 0.934 (0.823), 0.912 (0.766) and 0.935 (0.826), respectively.
32 Trans-GPCR generated a slightly lower Ac and Mcc values than that of HMMTOP,
33 TMHMM and Phobius. Although the these methods were benchmarked on the same
34 dataset, it should be pointed out that proteins in the test dataset of Trans-GPCR share
35 low similarity with the proteins in the training dataset (BLAST e-value>0.01).
36 Meanwhile, the TM regions of some proteins in Swiss-Prot database are annotated by
37 using TMHMM, Memast and Phobius (see <http://www.uniprot.org/manual/transmem>
38

1 for details). The complementarity of these methods is given in Figure 1 using
2 VennDiagram package⁴⁸. For example, HMMTOP, TMHMM, Phobius, Memast and
3 Trans-GPCR methods correctly distinguish 1197, 1130, 787, 885 and 1003 residues
4 that can not be correctly distinguished by other methods. In Figure 2, two *Mcc* values
5 of each protein by two methods correspond to a point. We calculated the statistical
6 significances of them using the student t-test (Table 2). The p-values of *Mcc* scores for
7 the methods were lower 0.01 although both HMMTOP, TMHMM and Phobius were
8 HMM-based algorithms. The different and complementary methods can be combined
9 to generate improved performance. This is demonstrated by the TM-combined method,
10 which generated the highest *Ac* (0.935) and *Mcc* (0.828) values in the
11 GPCR_TEST492 dataset. The increase in sensitivity using TM-combined may be
12 ascribed to that TM-combined measure is a consensus method by considering scores
13 of the top four methods. But TM-combined method did not generate higher *Mcc* value
14 in the GPCR_TRAIN1697, and this may be because the proteins of
15 GPCR_TRAIN1697 were used to train Trans-GPCR method. Therefore, it is very
16 difficult for TM-combined method to generate better performance. Meanwhile, we
17 also calculated the Pearson's correlation coefficient (*Pcc*) between them (Figure 2). As
18 seen from the data above, we can know that the benchmarked five methods were
19 significantly different (p-value<0.01). The most significant methods were TMHMM
20 and Memast (p-value<2.2e-16). To better understand the prediction error generation, it
21 is important to know the misclassification rates between TM/non-TM. As can be seen
22 from Table 3, the largest misclassification state is TM to non-TM, which is consistent
23 for the five predictors.

24

25 3.2 Benchmark of GPCR identification

26 The performance of GPCR identification was compared via ROC analysis. As can be
27 seen from Figure 3, PPA-GPCR generated the best performance, resulting in an AUC
28 score of 0.990. Trans-GPCR and SSEA-GPCR generated AUC scores of 0.978 and
29 0.955. Because the performance at low false positive rates is more important in
30 real-world application, therefore, we paid more attention to the comparison of
31 different methods' performance at < 1% false positive rates (Figure 3B and Table 4).
32 As shown in Table 4, SSEA-GPCR correctly recognized 193 GPCRs before including
33 36 false positives, whereas Trans-GPCR can detect 306 GPCRs. The distribution of
34 profile-to-profile alignment scores (i.e. *PPA_gpcr* measure) in the three types of
35 proteins (i.e. GPCR, non-GPCR membrane proteins, and globular proteins) are
36 presented in Figure 4. The confident interval (CI) values of *PPA_gpcr* for GPCRs,
37 non-GPCR membrane proteins and globular proteins were [13.53, 14.47], [6.77, 7.27]
38 and [2.48, 2.61] (Table 5), respectively. There is no overlap among these intervals,
39 suggesting PPA-GPCR method can be used to distinguish GPCRs in a reasonable
40 result. PPA-GPCR detects the more GPCRs (385 hits) than Trans-GPCR and
41 SSEA-GPCR methods at the same false positives cutoff. When we used the
42 *Iden-Combined* measure to identify GPCRs and it identifies the most GPCRs at the
43 1% false positive rate (Figure 3 and Table 4). Despite the lack of sequence homology

1 between some GPCRs, all GPCRs share similar conserved secondary structural
2 topologies and have the homologous relationships. Therefore, SSEA-GPCR and
3 PPA-GPCR should be effective to detect them. Our benchmark results also support
4 this point of view.
5

6 **3.3 Significances of prediction scores and decision making**

7 It is very necessary to estimate the significances of predictions when developing new
8 probabilistic models. We estimated the significant scores of Trans-GPCR,
9 SSEA-GPCR and PPA-GPCR from the test dataset. In the Trans-GPCR method, we
10 designed two output nodes in two neural networks to represent the prediction scores
11 of TM/non-TM regions. The difference of the two nodes of the second neural network
12 for target residue is represented by the measure *residue_reliable(i)*. The larger
13 *residue_reliable(i)* score is, the more significant and reliable for target residue. In our
14 benchmark result, if the *residue_reliable(i)* > 0.911, it can generate a prediction result
15 with less than a 1% false positive rate. Meanwhile, we also tested the
16 *TransGPCR_Score*, *SSEA_gpcr* and *PPA_gpcr* scores, which are the parameters to
17 identify GPCRs, in the benchmark dataset to obtain their reliable cutoffs. In our
18 benchmark, if *TransGPCR_Score* is larger than 84.834, the prediction result is at less
19 than a 5% false positive rate. At the same false positive rate control, *SSEA_gpcr* and
20 *PPA_gpcr* should be larger than 0.094 and 7.545, respectively. The prediction scores
21 and corresponding false positive rates were summarized in Table 6. A question should
22 be discussed here is that how to determine whether a protein is GPCR using these
23 methods. We suggest combining the three methods to make decisions. If proteins are
24 predicted to have less than 1% false positive rates by the three methods, the proteins
25 should be regarded as candidates for being GPCRs with high confidences. It is easy to
26 distinguish GPCRs and globular proteins. However, it may be difficult to distinguish
27 GPCRs from some non-GPCR membrane proteins according to the fact that some of
28 them have similar topologies and exist in the similar biological environments. For
29 such cases, maybe researchers can use the TM helices number and *PPA_gpcr* scores
30 to determine whether a protein is GPCR or not. Alternatively, users can resort to
31 combined methods (i.e. TM-Combined and Iden-Combined) to make decisions. If
32 *Iden-Combined* score higher than 1.589 by the combined method, the prediction is at
33 less than a 1% false positive rate. Some hard targets may need further literature survey.
34 In our web server (see supplementary file 2 for details), we provide the prediction
35 scores by Trans-GPCR, SSEA-GPCR, PPA-GPCR and Iden-Combined for each job.
36 To provide a real application example, we conduct our method on the proteome of
37 *Homo sapiens* (see supplementary file 3 for details).

39 **3.4 Lengths and amino acid distribution of TM/non-TM regions**

40 We calculated the mean lengths for TM/non-TM regions, but did not find significant
41 differences between TM segments of different GPCRs in the training dataset. The
42 lengths of TM helices are in the range of 6 to 30 amino acids and the average length

1 of TM helices is 22. The Beta-1 adrenergic receptor (Swiss-Prot entry: Q9TT96)
2 contains the longest (30 residues) TM segment in the sixth TM in our training dataset.
3 Although Q7P0A1, Q6BKW6 and Q60880 proteins contain TM segments longer than
4 30, the segments of them were annotated as two independent parts. For example, the
5 220-261 of protein Q7P0A1 is TM region. But 220-240 and 241-261 of this long
6 region were annotated as two independent parts in the Swiss-Prot database. For such
7 regions, we also counted them as two segments. Meanwhile, putative olfactory
8 receptor 10J6 protein (Swiss-Prot entry: Q8NGY7) contains the minimum length TM
9 regions (6 residues) in our training dataset. The length of loops connecting TM helices
10 is more diverse. The protein Q4LBB6 contains the longest loop (843 residues), which
11 connects the fifth and sixth TMs.

12

13 The amino acid compositions in the TM, non-TM regions and differences of them are
14 shown in Figure 7, in which the similarities and differences of the 20 amino acid
15 residues in the TM/non-TM regions were brought out. Residues with positive scores
16 of differences suggest their preference in the TM regions while those with negative
17 scores show their dominance in the non-TM regions. As can be seen from Figure 7,
18 the most differences of amino acids are R (ARG), K (LYS), E (GLU), L (LEU), V
19 (VAL) and I (ILE). Among them, L, V and I are aliphatic amino acids; R, K, and E are
20 charged amino acids. Interestingly, L, V and I are enriched in the TM regions whereas
21 R, K, and E are enriched in non-TM regions. Meanwhile, C (CYS) and G (GLY) show
22 subtle difference in the amino acids composition. The amino acid compositions
23 differences in the TM and non-TM regions can be regarded as conformational
24 parameters of amino acids in TM regions. Similarly, Gromiha developed a set such
25 conformational parameters in a different way in 1999. Pearson's correlation
26 coefficient between our parameters and those developed by Gromiha is 0.932 (see
27 supplementary file 4 for details), suggesting both sets of parameters can be used to
28 represent the preferences of amino acids in the TM regions although they are
29 calculated using different ways. Meanwhile, we also tested the performance of
30 secondary structure prediction by PSIPRED on GPCRs, and PSIPRED shows an
31 overall Q3 accuracy of 76.6% (see supplementary file 5 for details).

32

33 **4 CONCLUSIONS**

34 In this work, we developed a practical predictor for GPCR TM region prediction
35 (Trans-GPCR), and GPCR identification (Trans-GPCR, SSEA-GPCR and
36 PPA-GPCR). Our predictor has been intensively benchmarked and has been
37 demonstrated its favorable performance in the real application.

38

39 Objectively speaking, our predictor has strengths and limitations compared to some
40 other methods. The most obvious strength is its potential application to identify
41 GPCRs that show little sequence similarity to known GPCRs but with similar
42 topologies or homologous relationships. However, the qualities of both GPCR
43 identification and their TM regions location are relied on the input profiles, and it may
44 create problems if there are false homologous sequences imbedded in the MSAs that

1 used to calculate sequence profiles. This is one obvious limitation/disadvantage of our
2 predictor.

3 Anyway, our server should be useful based on its performance in the benchmark.
4 Although our predictor is a solely computational tool, we also hope that the
5 development of such novel methods will be helpful to accelerate the exploration of
6 the sequence-structure-function landscape in GPCRs.
7

8 ACKNOWLEDGMENT

9 We would like to thank Shaoyu Su in Fujian Science and Technology Information
10 Institute for critical reading the manuscript. This work was supported by Start-Up
11 Fund of Fuzhou University (510046), National Natural Science Foundation of China
12 (31301537) and Science Development Foundation of Fuzhou University
13 (2013-XY-17).
14

15 REFERENCES

- 16 1. G. G. Hazell, C. C. Hindmarch, G. R. Pope, J. A. Roper, S. L. Lightman, D. Murphy, A. M.
17 O'Carroll and S. J. Lolait, *Frontiers in neuroendocrinology*, 2012, **33**, 45-66.
- 18 2. R. T. Dorsam and J. S. Gutkind, *Nature reviews. Cancer*, 2007, **7**, 79-94.
- 19 3. F. Giordano, S. Simoes and G. Raposo, *Proceedings of the National Academy of Sciences of the*
20 *United States of America*, 2011, **108**, 11906-11911.
- 21 4. D. K. Vassilatis, J. G. Hohmann, H. Zeng, F. Li, J. E. Ranchalis, M. T. Mortrud, A. Brown, S. S.
22 Rodriguez, J. R. Weller, A. C. Wright, J. E. Bergmann and G. A. Gaitanaris, *Proceedings of the*
23 *National Academy of Sciences of the United States of America*, 2003, **100**, 4903-4908.
- 24 5. J. P. Overington, B. Al-Lazikani and A. L. Hopkins, *Nature reviews. Drug discovery*, 2006, **5**,
25 993-996.
- 26 6. H. M. Berman, *Acta crystallographica*, 2008, **64**, 88-95.
- 27 7. S. F. Altschul, W. Gish, W. Miller, E. W. Myers and D. J. Lipman, *Journal of molecular biology*,
28 1990, **215**, 403-410.
- 29 8. Q. Gao and A. Chess, *Genomics*, 1999, **60**, 31-39.
- 30 9. L. Rabiner, *Proceedings of the IEEE*, 1989, **77**, 257-286.
- 31 10. L. C. Chang CC, *Computer Program*, 2001.
- 32 11. L. Kall, A. Krogh and E. L. Sonnhammer, *Journal of molecular biology*, 2004, **338**, 1027-1036.
- 33 12. A. Krogh, B. Larsson, G. von Heijne and E. L. Sonnhammer, *Journal of molecular biology*, 2001,
34 **305**, 567-580.
- 35 13. M. Wistrand, L. Kall and E. L. Sonnhammer, *Protein Sci*, 2006, **15**, 509-521.
- 36 14. G. E. Tusnady and I. Simon, *Journal of molecular biology*, 1998, **283**, 489-506.
- 37 15. P. K. Papasaïkas, P. G. Bagos, Z. I. Litou and S. J. Hamodrakas, *SAR and QSAR in environmental*
38 *research*, 2003, **14**, 413-420.
- 39 16. T. Nugent and D. T. Jones, *BMC bioinformatics*, 2009, **10**, 159.
- 40 17. M. M. Gromiha, *Protein engineering*, 1999, **12**, 557-561.
- 41 18. M. Bhasin and G. P. Raghava, *Nucleic acids research*, 2004, **32**, W383-389.
- 42 19. W. Yang, K. Wang and W. Zuo, *International journal of bioinformatics research and*
43 *applications*, 2013, **9**, 207-219.
- 44 20. S. Takeda, S. Kadowaki, T. Haga, H. Takaesu and S. Mitaku, *FEBS letters*, 2002, **520**, 97-101.

- 1 21. D. W. Elrod and K. C. Chou, *Protein engineering*, 2002, **15**, 713-715.
 2 22. K. C. Chou, *Journal of proteome research*, 2005, **4**, 1413-1418.
 3 23. X. Xiao, J. L. Min, P. Wang and K. C. Chou, *PloS one*, 2013, **8**, e72234.
 4 24. K. C. Chou, *Journal of theoretical biology*, 2011, **273**, 236-247.
 5 25. W. Chen, P. M. Feng, H. Lin and K. C. Chou, *Nucleic acids research*, 2013, **41**, e68.
 6 26. Y. Xu, X. J. Shao, L. Y. Wu, N. Y. Deng and K. C. Chou, *PeerJ*, 2013, **1**, e171.
 7 27. X. Xiao, J. L. Min, P. Wang and K. C. Chou, *Journal of theoretical biology*, 2013, **337**, 71-79.
 8 28. J. L. Sussman, D. Lin, J. Jiang, N. O. Manning, J. Prilusky, O. Ritter and E. E. Abola, *Acta*
 9 *crystallographica*, 1998, **54**, 1078-1084.
 10 29. E. Boutet, D. Lieberherr, M. Tognolli, M. Schneider and A. Bairoch, *Methods in molecular*
 11 *biology (Clifton, N.J)*, 2007, **406**, 89-112.
 12 30. F. Horn, E. Bettler, L. Oliveira, F. Campagne, F. E. Cohen and G. Vriend, *Nucleic acids research*,
 13 2003, **31**, 294-297.
 14 31. N. K. Fox, S. E. Brenner and J. M. Chandonia, *Nucleic acids research*, 2014, **42**, D304-309.
 15 32. J. Heaton, 2008, 1-429.
 16 33. E. R. David, E. H. Geoffrey and J. W. Ronald, in *Neurocomputing: foundations of research*, eds.
 17 A. A. James and R. Edward, MIT Press 1988, pp. 696-699.
 18 34. D. T. Jones, *Journal of molecular biology*, 1999, **292**, 195-202.
 19 35. K. D. Pruitt, T. Tatusova, W. Klimke and D. R. Maglott, *Nucleic acids research*, 2009, **37**,
 20 D32-36.
 21 36. S. F. Altschul, T. L. Madden, A. A. Schaffer, J. Zhang, Z. Zhang, W. Miller and D. J. Lipman,
 22 *Nucleic acids research*, 1997, **25**, 3389-3402.
 23 37. S. Henikoff and J. G. Henikoff, *Journal of molecular biology*, 1994, **243**, 574-578.
 24 38. Z. Chen, Y. Wang, Y. F. Zhai, J. Song and Z. Zhang, *Molecular bioSystems*, 2013, **9**, 2213-2222.
 25 39. S. B. Needleman and C. D. Wunsch, *Journal of molecular biology*, 1970, **48**, 443-453.
 26 40. T. Przytycka, R. Aurora and G. D. Rose, *Nature structural biology*, 1999, **6**, 672-682.
 27 41. R. X. Yan, Z. Chen and Z. Zhang, *BMC bioinformatics*, 2011, **12**, 76.
 28 42. Y. Zhang and J. Skolnick, *Nucleic acids research*, 2005, **33**, 2302-2309.
 29 43. D. Xu, L. Jaroszewski, Z. Li and A. Godzik, *Bioinformatics (Oxford, England)*, 2014, **30**, 660-667.
 30 44. T. Fawcett, *Pattern Recognition Letters*, 2006, **27**, 861-874.
 31 45. K. C. Chou, Z. C. Wu and X. Xiao, *PloS one*, 2011, **6**, e18258.
 32 46. K. C. Chou, Z. C. Wu and X. Xiao, *Molecular bioSystems*, 2012, **8**, 629-641.
 33 47. K. C. Chou, *Molecular bioSystems*, 2013, **9**, 1092-1100.
 34 48. H. Chen and P. C. Boutros, *BMC bioinformatics*, 2011, **12**, 35.

35 TABLES

36 **Table 1. Performance of TM region prediction of various methods on datasets.**

Method ^a	Ac	Sn	Sp	Mcc
Benchmark result on GPCR_TRAIN1697				
HMMTOP	0.910	0.896	0.919	0.814
TMHMM	0.907	0.890	0.920	0.809
Memsat	0.892	0.906	0.882	0.780
Phobius	0.903	0.894	0.909	0.801
Trans-GPCR ^b	0.940	0.930	0.948	0.877
TM-Combined	0.935	0.943	0.930	0.867

Benchmark result on GPCR_TEST492

HMMTOP	0.927	0.865	0.947	0.804
TMHMM	0.934	0.874	0.954	0.823
Memsat	0.912	0.848	0.932	0.766
Phobius	0.935	0.884	0.951	0.826
Trans-GPCR	0.923	0.833	0.952	0.791
TM-Combined	0.935	0.901	0.946	0.828

1 ^aAll residues of test dataset were used to count true positives (TP), true negatives (TN), false positives (FP) and
2 false negatives (FN) measures.

3 ^bTrans-GPCR was intensively trained on GPCR_TRAIN1697 dataset. Proteins in the GPCR_TEST492 dataset
4 share low similarity with proteins in GPCR_TRAIN1697 (BLAST e-value>0.01). Therefore, benchmark of
5 Trans-GPCR on GPCR_TRAIN1697 dataset does not make a lot of sense. We just want to know how much
6 performance decrease when tested Trans-GPCR on the GPCR_TEST492 compared with that of
7 GPCR_TRAIN1697.

8 **Table 2. The student t-test p-values of the five methods of *Mcc* scores**

Method ^a	HMMTOP	TMHMM	Memsat	Phobius	Trans-GPCR
HMMTOP		3.594e-05	9.885e-13	1.443e-06	3e-4
TMHMM			2.2e-16	0.6975	4.947e-13
Memsat				2.2e-16	1.588e-08
Phobius					2.2e-16
Trans-GPCR					

9

10 **Table 3. Misclassification rates in the benchmark dataset**

Native	Predicted	HMMTOP	TMHMM	Memsat	Phobius	Trans-GPCR
M	-	0.134	0.126	0.151	0.115	0.166
-	M	0.052	0.045	0.067	0.048	0.047

11 ^aHere 'M' and '-' represent transmembrane and non-transmembrane residues. Misclassification rate is calculated
12 using equation $E(i)/N(i)$, where $E(i)$ is the number of misclassified state i and $N(i)$ is the total number of state i in
13 the benchmark dataset.

14

15 **Table 4. ROC table (≤ 36 false positives) for different methods**

Methods	Receiver operator characteristics (≤ 36 false positives ^a)							
	12	16	20	24	28	32	36	Auc ^b
Trans-GPCR	133	165	233	266	289	293	306	0.978
SSEA-GPCR	120	139	160	173	188	192	193	0.955
PPA-GPCR	319	346	354	356	374	382	385	0.990
Iden-Combined	343	381	388	411	431	444	461	0.993

16 ^a Here, false positives correspond to those non-GPCRs predicted as GPCRs.

17 ^b The Auc score represents the corresponding area under a ROC curve.

18

19

20

21 **Table 5. Mean, standard deviation and confidence intervals (CI) at a 95% level**

Methods ^a	Mean	Standard deviation	CI
----------------------	------	--------------------	----

GPCRs	14.00	5.32	[13.53, 14.47]
Membrane proteins	7.02	2.88	[6.77, 7.27]
Globular proteins	2.54	1.44	[2.48, 2.61]

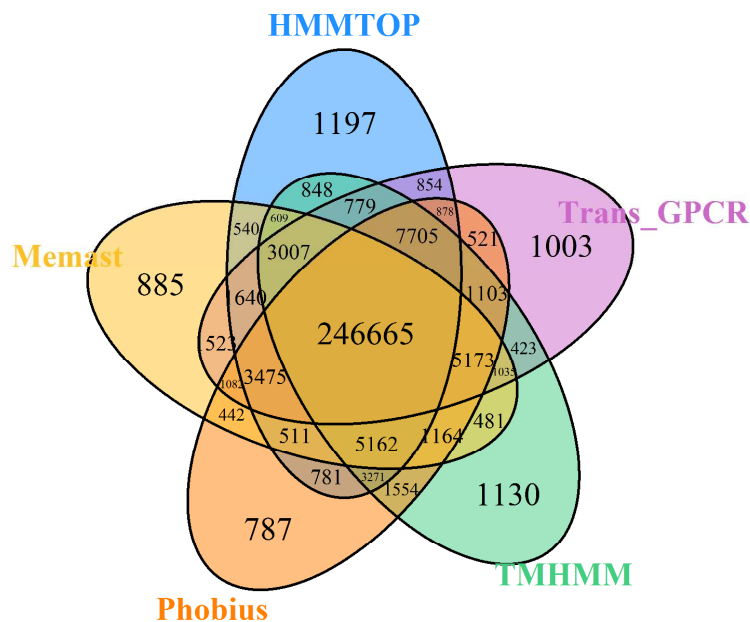
1

2 **Table 6. Cutoffs of different methods at 95% and 99% confident levels**

Methods	95% level	99% level
TransGPCR_Score ^a	84.834	112.295
residue_reliable ^a	0.000	0.911
SSEA_gpcr	0.094	0.139
PPA_gpcr	7.545	9.664
Iden-Combined	1.354	1.589

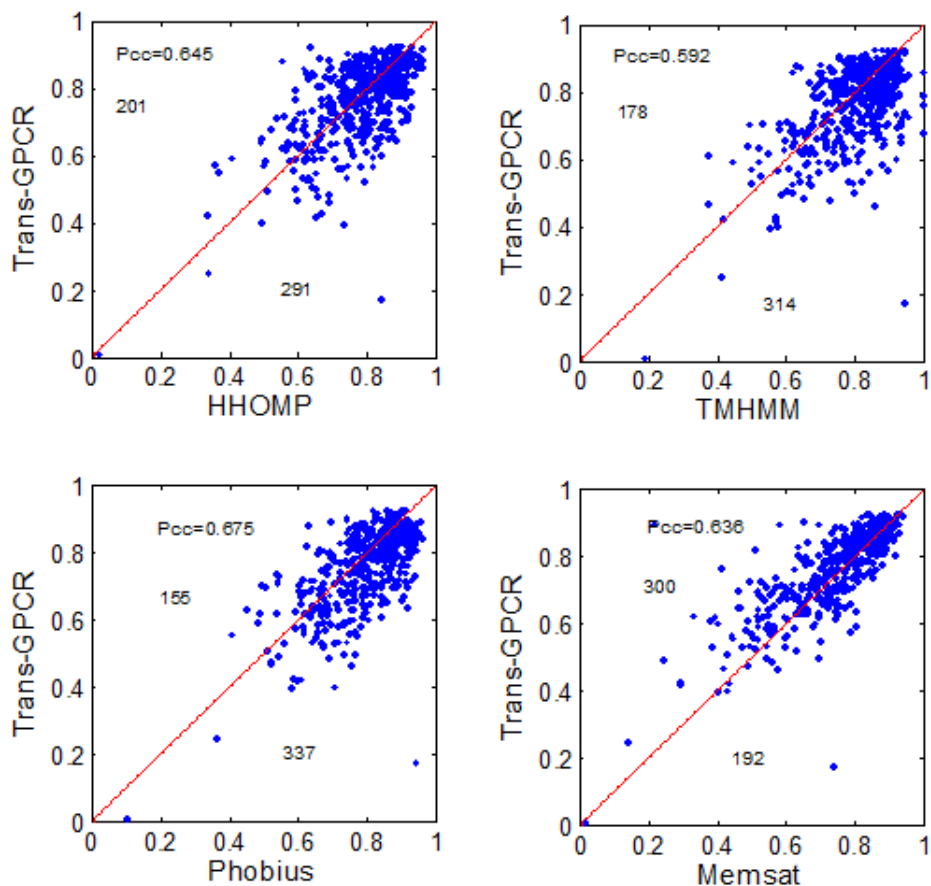
3 ^aTransGPCR_Score is a measure to determine whether a protein is GPCR, whereas residue_reliable(i) is a
 4 parameter to describe the reliability in position *i* of a protein (i.e. TM or non-TM residue).

5

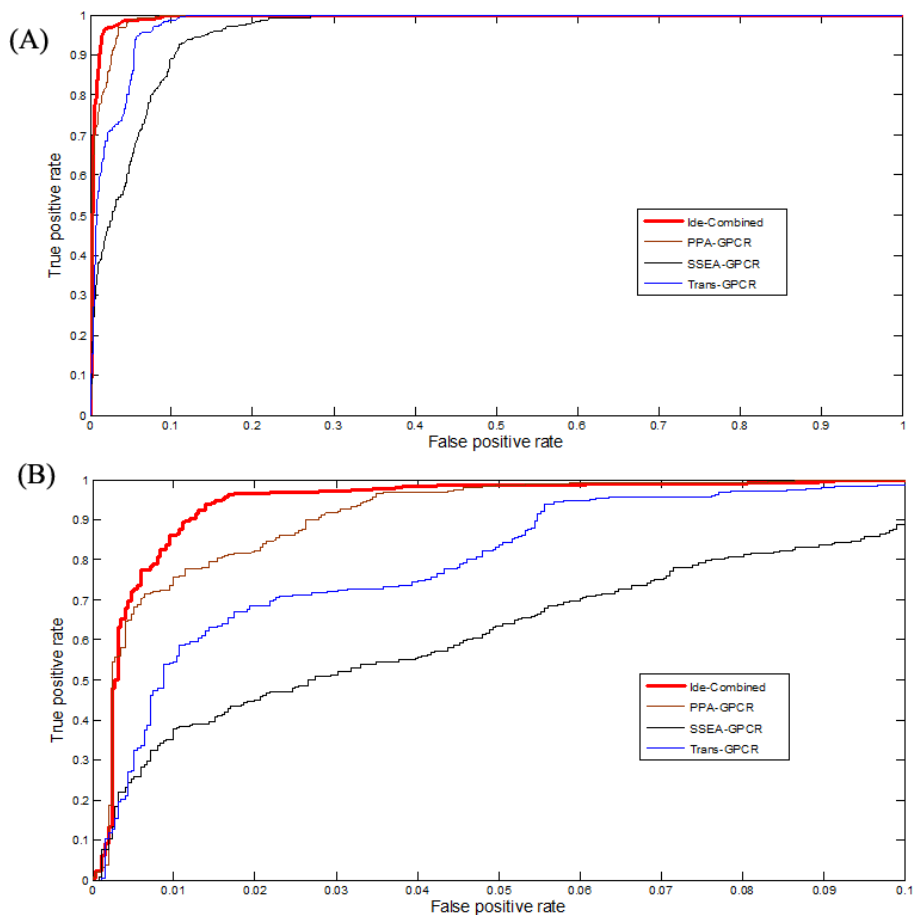
6 **Figures**

7

8 **Figure 1. Venn diagram showing the complementarity of various methods.** For
 9 the samples correctly distinguished by two or more methods, they correspond
 10 to the number in the overlapped regions.



1
 2 **Figure 2. All-to-all comparisons of *Mcc* scores between methods on the**
 3 **GPCR_TEST492 dataset.** The number in each panel denotes the number of
 4 proteins/points in upper and lower triangles, respectively. Meanwhile, Pearson's
 5 correlation coefficient (*Pcc*) values are also given.

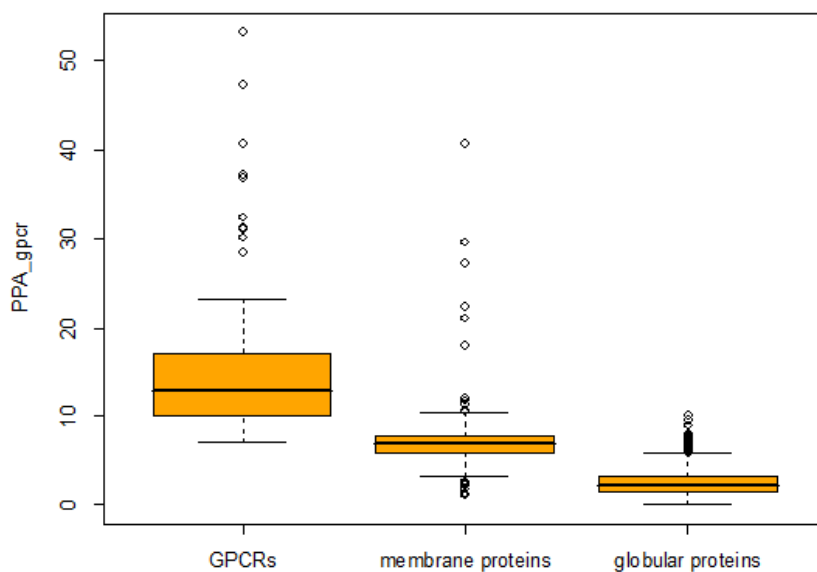


1

2

Figure 3. Comparison of ROC curves for different methods.

3

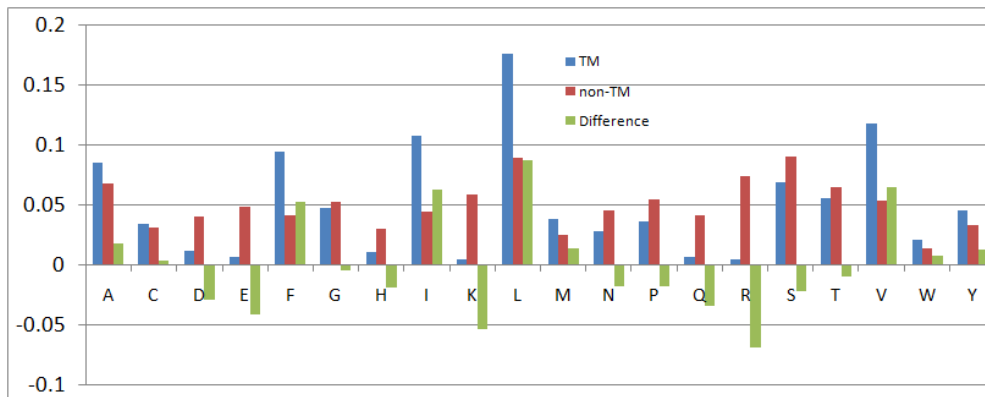


4

5

Figure 4. Boxplot of *PPA_gpcr* scores in the three types of proteins. Here, membrane proteins denote the non-GPCR membrane proteins.

6



1

2 **Figure 5. Amino acid composition of the 20 amino acid residues in TM regions**
 3 **(blue bars), non-TM regions (red bars) and differences between them (green**
 4 **bars).**

5

6 **Supplementary files**

7 **Supplementary file 1:** Removing redundancies of datasets

8 **Supplementary file 2:** The web server for GPCR prediction

9 **Supplementary file 3:** Proteome-wide GPCR identification in *Homo sapiens*

10 **Supplementary file 4:** Correlation of conformational parameters for TM helices

11 **Supplementary file 5:** PSIPRED for protein GPCR secondary structure prediction

12