

Spring 2025

BIOS 477/877

Bioinformatics and Molecular Evolution

Lecture 26

BIOS477/877 L26 - 1

1

TODAY'S TOPICS

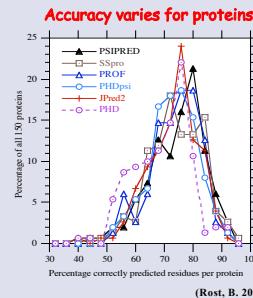
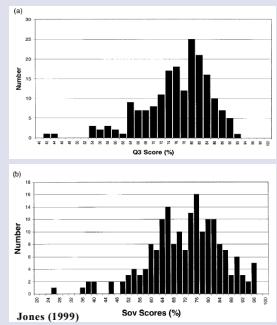
- Secondary structure prediction
- Transmembrane protein prediction
- Assignment 12 (due: May 9)

BIOS477/877 L26 - 2

2

PSIPRED prediction accuracy

Distribution of prediction accuracy from 187 PDB proteins (structure is known)



Q3=84.2% in ver 4 Buchan & Jones (2019)

BIOS477/877 L26 - 3

Sixty-five years of the long march in protein secondary structure prediction: the final stretch?

Yuedong Yang, Jianzhao Gao, Jihua Wang, Rhys Heffernan, Jack Hanson, Kuldip Paliwal and Yaoqi Zhou

Briefings in Bioinformatics, 19(3), 2018, 482–494

Data set	Method	TS115		CASP12		Server location
		Q3	P-value*	Q3	P-value*	
Jpred4		0.771 ^b	0.0007	0.751	0.04	http://www.compbio.dundee.ac.uk/jpred4/index.html
SPINE X		0.801	0.0002	0.769	0.006	http://spine-lab.gstt.org/
PSIPRED 3.3		0.809	0.01	0.799	0.19	http://spine-lab.gstt.org/psipred3.3/
SCORPION		0.817	0.45	0.805	0.44	Stand-alone version from http://bpcr.cs.odu.edu/c3scorpion/
SPIDER2		0.819	NA	0.798	NA	http://spine-lab.org/server/SPIDER2/
PORTER 4.0		0.820	0.17	0.798	0.67	http://distrill.ucd.ie/porterpalae/
DeepCNP		0.823	0.01	0.821	0.14	http://raptorx2.uchicago.edu/StructurePropertyPred/predict/

Note: *Paired t-test from SPIDR2.

SCORPION: Singlet/doublet/triplet prob. + PSSM

SPIDER2: Iterative learning of secondary structure, backbone torsion angles, and solvent accessible surface area using deep NN (three layers) with PSSM + physico-chemical properties

DeepCNP: Deep convolutional NN with PSSM + 78 AA related features

PSIPRED 4: Q3=84.2%; PSSM + 2 hidden layer NNs

BIOS477/877 L26 - 4

4

Protein secondary structure prediction using neural networks and deep learning: A review

Computational Biology and Chemistry 81 (2019) 1–8

Wafaa Wardah^a, M.G.M. Khan^a, Alok Sharma^{b,c}, Mahmood A. Rashid^{a,d,*}

The major periodically relevant state-of-the-art methods are shown along with the types of feature values they employed in their networks.

Neural network method	Accuracy (Q3)	Sq info	Evo info	Physico chem info
Qian & Sejnowski 1988 (Qian and Sejnowski, 1988)	64.3%	/	/	
PHD 1994 (Rost et al., 1994)	71.4%	/	/	
PSIPRED 1997 (Jones, 1999)	76.5%	/	/	
Jpred3 2008 (Cole et al., 2007)	81.5%	/	/	/
SPIDER3 2017 (Heffernan et al., 2017)	84%	/	/	/

Evolutionary or profile information = MSA

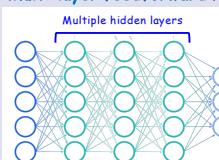
Limit of secondary structure prediction accuracy: estimated to be 88% (Rost 2001)

BIOS477/877 L26 - 5

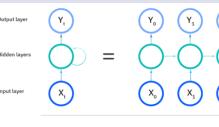
5

Deep Neural Networks

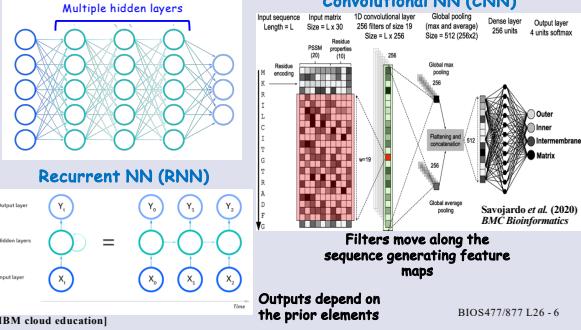
Multi-layer feedforward NN



Recurrent NN (RNN)

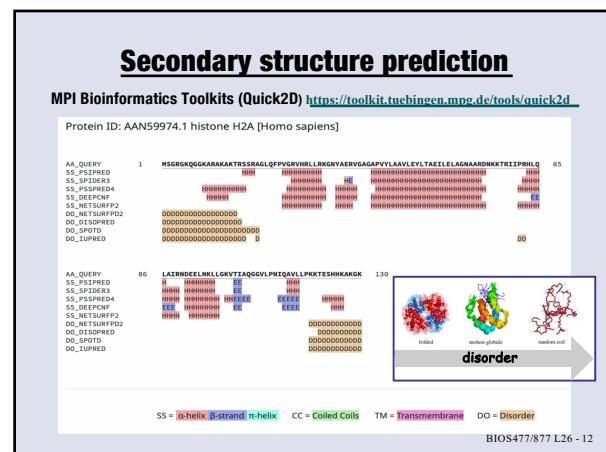
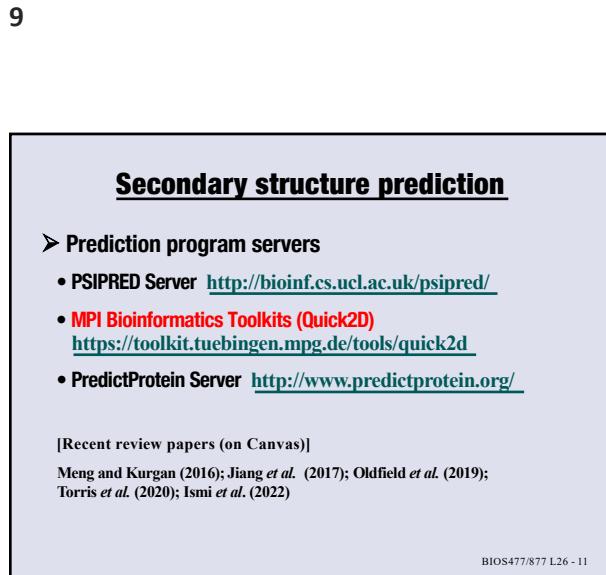
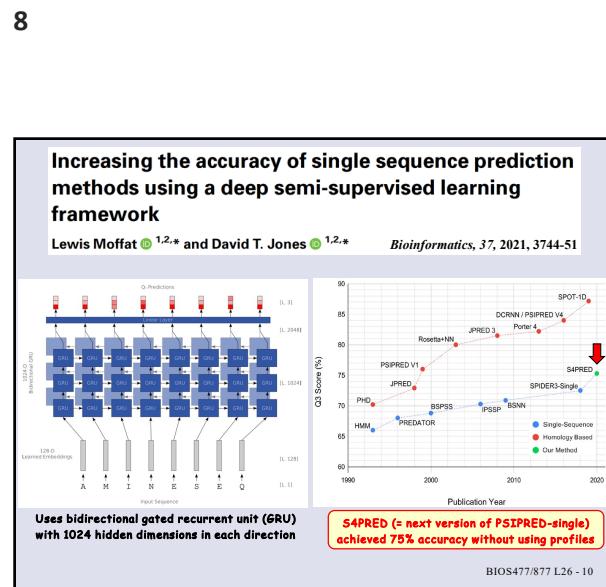
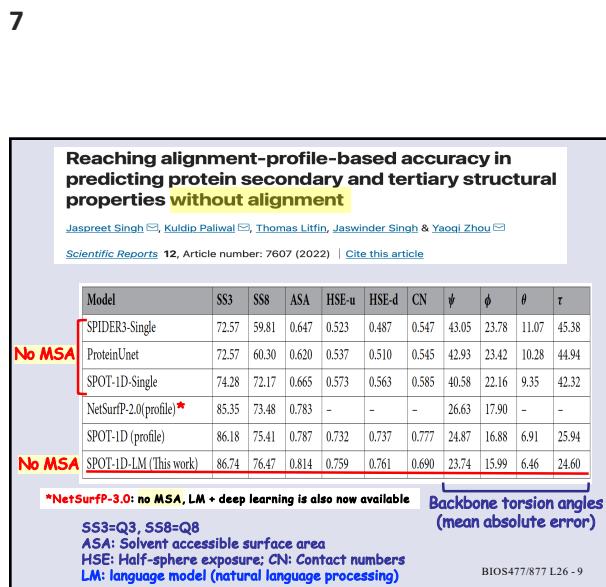
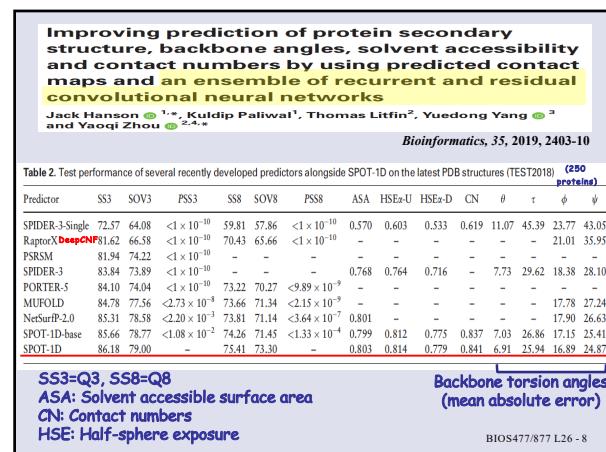
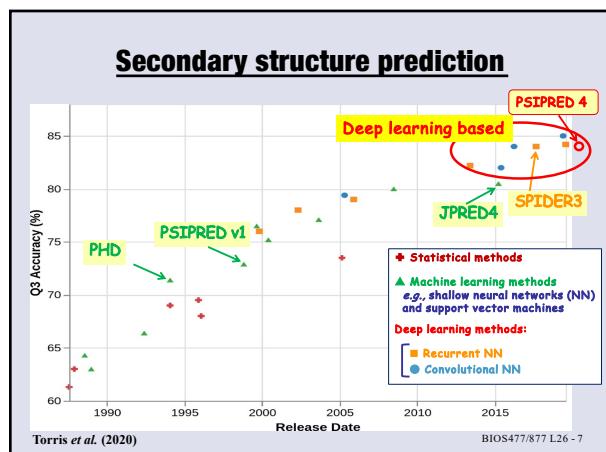


Convolutional NN (CNN)



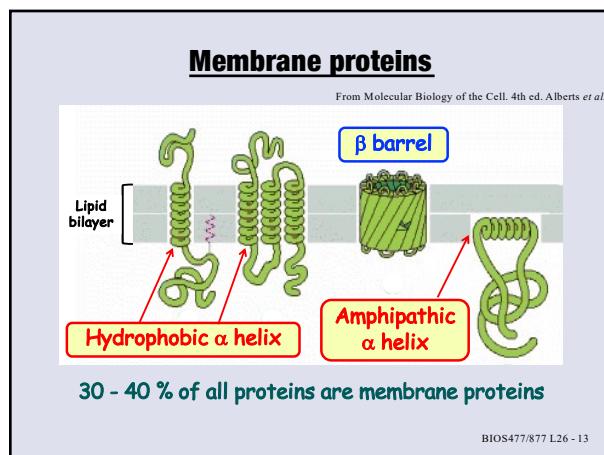
BIOS477/877 L26 - 6

6

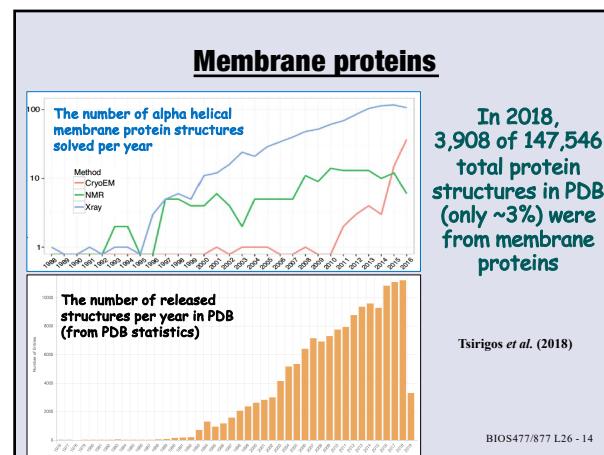


11

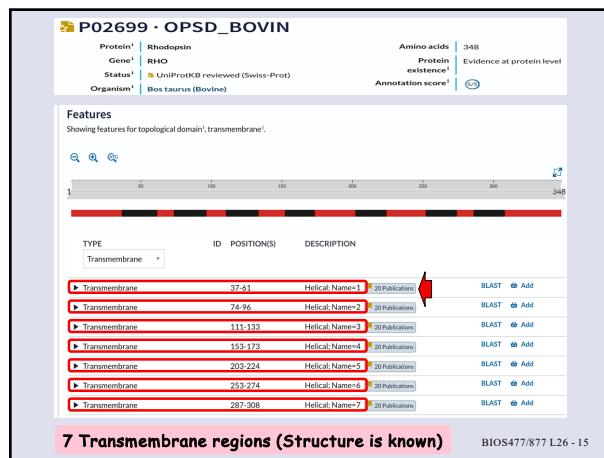
12



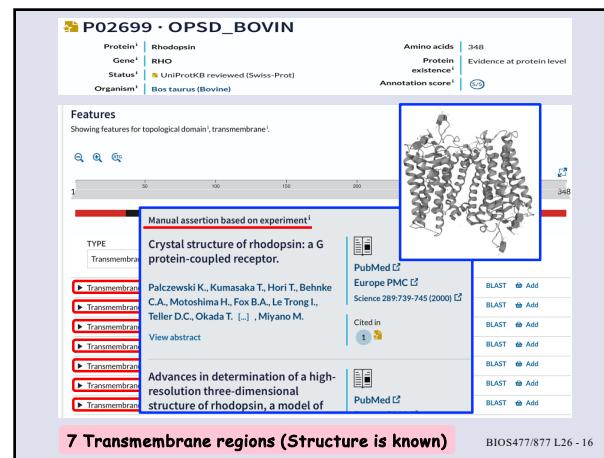
13



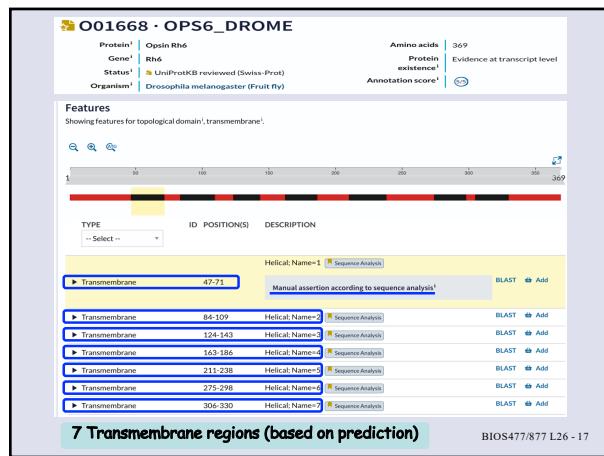
14



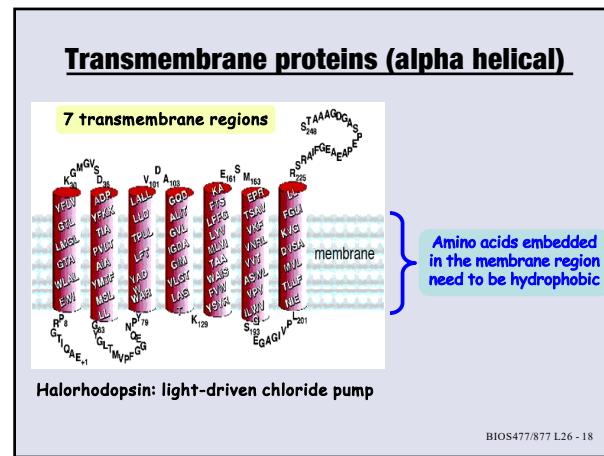
15

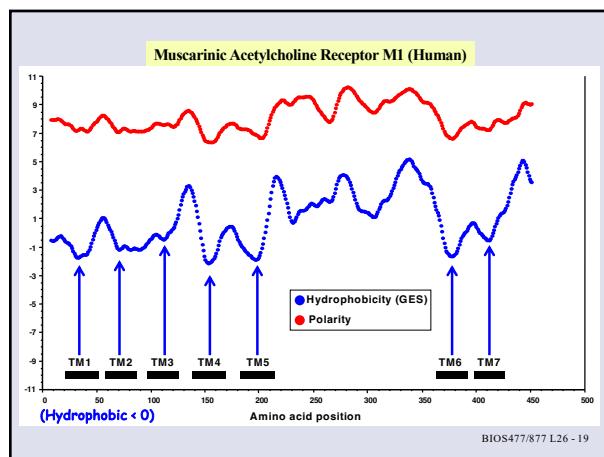


16

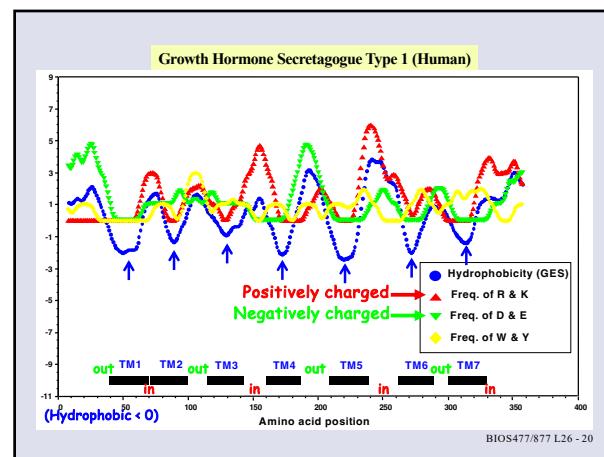


17

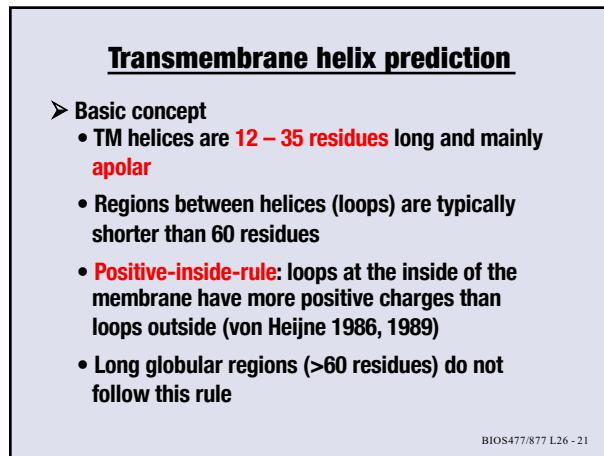




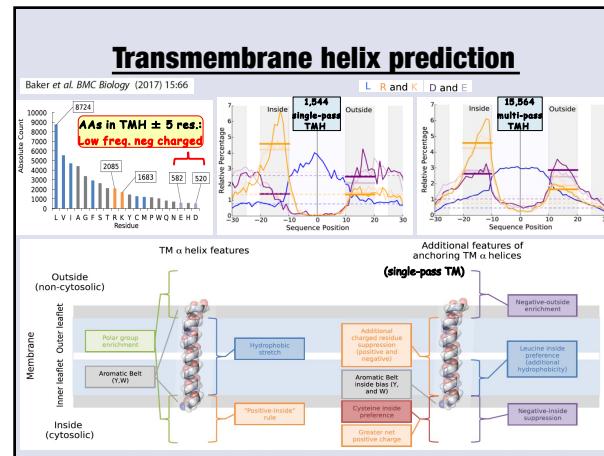
19



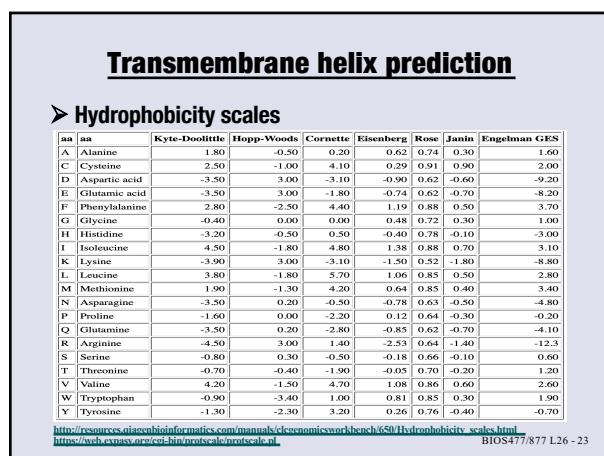
20



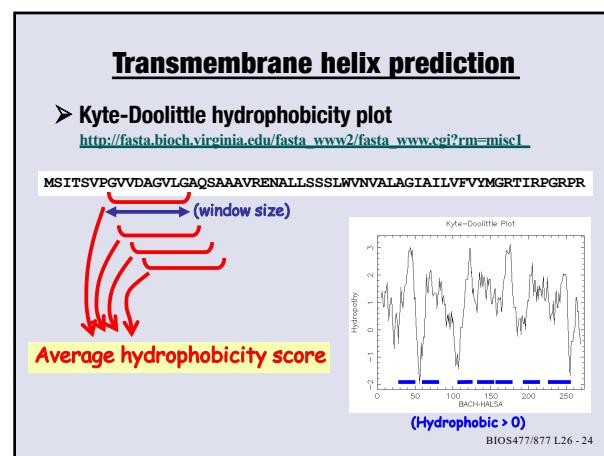
21



22



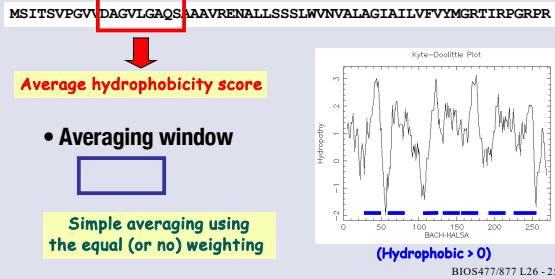
23



Transmembrane helix prediction

Kyte-Doolittle hydrophobicity plot

http://fasta.bioch.virginia.edu/fasta_www2/fasta_www.cgi?rm=misc1



25

Transmembrane helix prediction

TopPred (Heijne et al. 1992)

<https://github.com/C3BL-pasteur-fr/toppred>

→ Uses GES and other hydrophobicity scale

Averaging window



Trapezoid averaging (smaller weighting at both sides and more weighting at the window center)

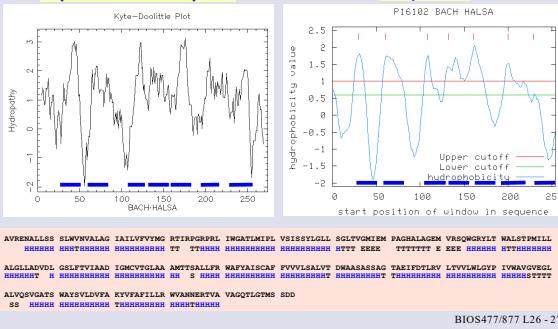
• Uses “positive-inside-rule” to predict the transmembrane topology

BIOS477/877 L26 - 26

26

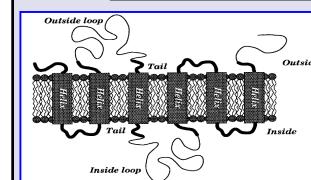
Transmembrane helix prediction

Kyte-Doolittle plot

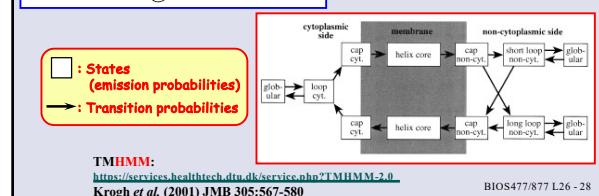


27

Transmembrane helix prediction

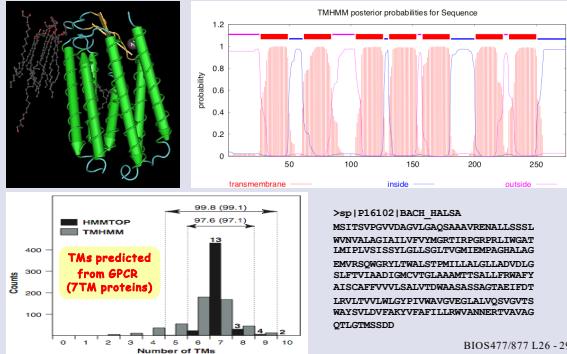


HMMTOP:
Tusnady & Simon (1998) JMB 283: 489-506



28

Transmembrane helix prediction



29

Evaluation of transmembrane helix predictions in 2014

Jonas Reeb,¹ Edda Kloppmann,^{1,2,*} Michael Bernhofer,¹ and Burkhard Rost^{1,2,3,4}
Proteins 83: 473-484

Table 1
Transmembrane Helix Prediction Methods

Name	Year	Method	Evolutionary information	Signal peptides	Topology
TopPred	1994	—	No	No	Yes
PHDhtm	1995	NN	Yes	No	Yes
HMMTOP	2001	HMM	No	No	Yes
TMHMM2	2001	HMM	No	No	Yes
SOSUI	2002	—	No	No	No
Phobius	2004	HMM	No	Yes	Yes
PolyPhobius	2005	HMM	Yes	Yes	Yes
MEMSAT3	2007	NN	Yes	Yes	Yes
PhiBis	2008	DBN	No	Yes	Yes
SCAMPI	2008	HMM	No	No	Yes
SPOTOPUS	2008	NN+HMM	Yes	Yes	Yes
MEMSAT-SVM	2009	SVM	Yes	Yes	Yes

Signal peptide prediction can be incorporated with TM prediction
→ False positive rates can be reduced

BIOS477/877 L26 - 30

30

Phobius
A combined transmembrane topology and signal peptide predictor

Phobius

<https://phobius.sbc.su.se/> or
<https://www.ebi.ac.uk/idispatcher/pfa/phobius/>

Signal peptide model

(Käll et al. 2004 J Mol Biol 338: 1027-1036; Käll et al. 2007 Nucl Acids Res 35: W429-W432)

BIOS477/877 L26 - 31

31

Evaluation of transmembrane helix predictions in 2014

Jonas Reeb,¹ Edda Kloppmann,^{1,2,*} Michael Bernhofer,¹ and Burkhard Rost^{1,2,3,4}
Proteins 83: 473-484

Table II
Discrimination Between TMs and Soluble Proteins As Well As TMs and Signal Peptides

	(a) no SPs			(b) Sol [no TM]			(c) TMP			
	Sol [no TM]	Sol [no TM]	Sol [no TM]	Euk (5106) FPR	Gram- (356) FPR	Gram+ (911) FPR	Euk (1297) FPR	Gram- (400) FPR	Gram+ (204) FPR	Euk (332) Sens
TopPred2 ^a	62	53	52	97	99.8	98	100	98	99	99
PHDhtm ^b	11	11	19	31	17	29	97	89	96	97
HMMTOP ^c	21	17	15	56	77	25	63	93	97	97
TMHMM2 ^d	1	1	1	20	25	63	93	97	97	97
SOSUI ^e	3	1	62	32	44	97	99	99	99	99
Phobius	3	1	1	4	2	13	99	99	99	99
PolyPhobius	6	5	5	5	5	5	97	97	97	97
MEMSAT3	8	4	3	60	47	67	100	100	100	100
Philius	2	1	1	3	1	12	93	93	93	93
SCAMPI ^f	30	27	21	92	95	97	100	100	100	100
SPECTOPUS	11	24	25	13	3	13	97	97	97	97
MEMSAT-SVM	6	5	6	25	6	24	99	99	99	99

(a) Soluble proteins without signal peptides (SP)
(b) Soluble proteins with SP, (c) Transmembrane proteins with SP

FPR: false positive rate; Sens: Sensitivity

BIOS477/877 L26 - 32

32

Evaluation of transmembrane helix predictions in 2014

Jonas Reeb,¹ Edda Kloppmann,^{1,2,*} Michael Bernhofer,¹ and Burkhard Rost^{1,2,3,4}
Proteins 83: 473-484

Table I
Transmembrane Helix Prediction Methods

Name	Year	Method	PSSM	Evolutionary information	Signal peptides	Topology
TopPred2	1994	—	No	No	No	Yes
PHDhtm	1995	NN	Yes	No	Yes	Yes
HMMTOP	2001	HMM	No	No	Yes	Yes
TMHMM2	2001	HMM	No	No	Yes	Yes
SOSUI	2002	—	No	No	No	No
Phobius	2004	HMM	No	Yes	Yes	Yes
PolyPhobius	2005	HMM	Yes	Yes	Yes	Yes
MEMSAT3	2007	NN	Yes	Yes	Yes	Yes
Philius	2008	DBN	No	Yes	Yes	Yes
SCAMPI	2008	HMM	No	No	Yes	Yes
SPECTOPUS	2008	NN+HMM	Yes	Yes	Yes	Yes
MEMSAT-SVM	2009	SVM	Yes	Yes	Yes	Yes

PSSM+SVM (binary classifier)
- TM helix or not
- inside or outside loop
- Re-entrant helix or not
- TM or globular protein

BIOS477/877 L26 - 33

33

MEMSAT-SVM

BMC Bioinformatics

Research article
Open Access
Transmembrane protein topology prediction using support vector machines
Timothy Nugent and David T Jones* *BMC Bioinformatics* 2009, 10:159 doi:10.1186/1471-2105-10-159

Table 2: Benchmark results for the SVM-based method ('MEMSAT-SVM') against a selection of leading topology predictors

Method	Algorithm	Correct helix count	Correct helix locations	Correct N-terminal	FP helix	FN helix	Correct SP topology	Correct RE topology	Correct topology
MEMSAT-SVM	SVM	95%	91%	91%	4%	5%	93%	64%	89%
OCTOPUS	NN + HMM	86%	81%	84%	14%	2%	21%	73%	79%
MEMSAT	NN	84%	76%	84%	8%	8%	57%	64%	76%
ENSEMBLE	NN + HMM	77%	76%	79%	18%	5%	7%	55%	67%
PHOBIOUS	HMM	75%	76%	79%	9%	16%	93%	36%	63%
HMMTOP	HMM	77%	76%	78%	18%	6%	29%	64%	63%
PRODIV	HMM	79%	64%	76%	19%	0%	18%	57%	57%
SVMTOP	SVM	66%	64%	66%	22%	22%	0%	55%	53%
TMHMM	HMM	75%	68%	72%	14%	20%	29%	55%	53%
PhDiem	NN	73%	54%	55%	23%	30%	29%	18%	45%

<http://bioinf.cs.ucl.ac.uk/psipred/>

BIOS477/877 L26 - 34

TOPCONS2

Fast sequence
MSA profile
Predicted topologies
Signal pep: M: TM region
i: inside or outside

TOPCONS HMM

Consensus topology

<http://topcons.net/>

BIOS477/877 L26 - 35

35

Prediction statistics (general)

Accuracy = $(TP+TN)/(TP+FN+FP+TN)$
Error = 1 - Accuracy
Sensitivity = $TP/(TP+FN) [=TP/(actual\ positives)]$
Specificity = $TN/(TN+FP) [=TN/(actual\ negatives)]$
True positive rate = Sensitivity (or Recall)
False positive rate = $FP/(TN+FP) = 1 - \text{Specificity}$
Precision (or Positive predictive value) = $TP/(TP+FP)$

- Mathews correlation coefficients: CC (or MCC)
$$CC = \frac{(TP \cdot TN - FP \cdot FN)}{\sqrt{(TN+FN)(TN+FP)(TP+FN)(TP+FP)}}$$
- If everything is correct, $CC = 1$, if all predictions are wrong, $CC = -1$
- F-measure (F or F_1): the harmonic mean of precision and recall
$$F = 2TP/(2TP+FN+FP) \quad 0 < F < 1$$
- ROC (receiver operating characteristic) plot
ROC plot can be obtained using different threshold scores for identifying positive and negative predictions
X: If prediction is perfect ($TP=1, FP=0$)
The worst case scenario (random choice prediction)

BIOS477/877 L26 - 36

36

CCTOP

Three steps:

- 1) Remove cleavable parts (e.g., signal peptide predicted by SignalP 4.0).
- 2) TMP filtering (distinguish TM and globular proteins) using a simple voting for the results of Phobius, Scampi, and TMHMM.
- 3) Topology prediction using ten methods (constrained HMM using HMMTOP model)

<http://cctop.ttk.hu/>

Method	Sensitivity	Specificity	MCC	Acc _{Tpg}	Acc _{Tpl}
CCTOP	98	98	97.7	84	81
TopCons	97	97.7	97.2	79	75
ScampiMsa	97	96.7	96.8	76	72
Pro	96	97.3	96.5	75	70
Prodip	96	94.8	95.5	75	69
Octopus	93	98.2	95.7	71	66
HMMTOP	95	94.7	94.8	69	64
Phobius	95	97	95.9	71	64
TMHMM	93	97.2	94.9	66	59
MetaTM	94	97	96.0	67	58
Phobius	93	97.3	95.1	62	56
MemSat	94	98	95.8	66	53
MemBrain	92	97.4	94.7	62	0

BIOS477/877 L26 - 37

37

MEMBRAIN3

<http://www.csbio.situ.edu.cn/bioinf/MemBrain/>

Journal of Molecular Biology (2020) 432, 1279–1296

Table 1. Performance comparison on the test set under the TMH definition with tails^a.

Algorithm	Residue-precision	Residue-recall	Helical-precision	Helical-recall	V _p	V _{top}
MEMSAT3	0.885	0.713	0.629	0.561	11	11
MEMSAT-SVM	0.890	0.701	0.591	0.547	10	7
TMpred	0.884	0.694	0.495	0.443	8	5
Phobius	0.886	0.750	0.579	0.563	12	10
PolyPhobius	0.916	0.731	0.73	0.590	11	7
SCAMPi	0.919	0.729	0.641	0.611	13	10
CD-TOMT	0.909	0.668	0.527	0.456	6	6
SPOTCUPUS	0.913	0.721	0.522	0.524	10	10
TMHMM2	0.890	0.719	0.552	0.497	12	10
TOPCONS	0.918	0.722	0.573	0.530	12	10
COOP	0.921	0.671	0.511	0.477	10	10
TopGraph	0.873	0.717	0.504	0.464	12	6
TMSEG	0.909	0.727	0.578	0.550	15	13
MemBrain 2.0	0.917	0.737	0.570	0.547	12	10
MemBrain 3.0	0.892	0.900	0.608	0.819	21	20

Precision: # correct TM / # predicted TM Recall: # correct TM / # true TM

BIOS477/877 L26 - 38

Precision: # correct TM / # predicted TM Recall: # correct TM / # true TM

BIOS477/877 L26 - 38

TMPSS: A Deep Learning-Based Predictor for Secondary Structure and Topology Structure Prediction of Alpha-Helical Transmembrane Proteins

Zhe Liu^{1,2}, Yingli Gong³, Yihang Bao⁴, Yuanzhuan Guo⁴, Han Wang^{**} and Guan Ning Lin^{1,2*}

Secondary structure prediction

Method	Trans SS Q3	Non-trans SS Q3
SSpro5 (with templates)	0.90	0.89
PSIPRED 4	0.94	0.79
RaptorX-Property	0.95	0.80
Porter 5	0.95	0.81
DeepCNF	0.91	0.80
Spider3	0.95	0.80
SPOT-ID	0.95	0.81
MUFOLD-SSW	0.94	0.81
JPred4	0.90	0.75
TMSS	0.97	0.78

TM topology prediction

Method	ACC	MCC
HMMTOP 2	0.84	0.64
OCTOPUS	0.87	0.71
TOPCONS	0.88	0.72
Phobius	0.87	0.71
PolyPhobius	0.88	0.72
SCAMPi	0.87	0.70
SPOTCUPUS	0.87	0.71
TMSS	0.90	0.76

Trans, transmembrane region; Non-trans, non-transmembrane region. Bold fonts represent the best experimental results.

BIOS477/877 L26 - 39

39

TMbed: transmembrane proteins predicted through language model embeddings

Michael Bernhofer^{1,2*} and Burkhard Rost^{1,3,4}

β-TMP (57)

Recall (%)	FPR (%)	Recall (%)	FPR (%)	Recall (%)	FPR (%)	
TMbed	93.8 ± 7.5	0.1 ± 0.1	97.5 ± 0.7	0.5 ± 0.2	99.5 ± 0.2	2.8 ± 1.2
DeepTMHMM	77.9 ± 12.7	0.1 ± 0.1	95.8 ± 1.3	0.5 ± 0.2	99.5 ± 0.2	5.9 ± 2.2
TMSEG	—	—	96.5 ± 1.0	2.3 ± 0.3	97.7 ± 0.3	3.5 ± 1.0
TOPCONS2 ¹	—	—	94.2 ± 1.3	2.6 ± 0.3	97.4 ± 0.3	5.8 ± 1.3
OCTOPUS ¹	—	—	94.2 ± 1.9	9.1 ± 0.7	90.9 ± 0.7	5.8 ± 1.9
Phobius ¹	—	—	92.5 ± 1.4	2.6 ± 0.2	97.4 ± 0.2	7.5 ± 1.4
PolyPhobius ¹	—	—	97.2 ± 1.1	5.3 ± 0.4	94.7 ± 0.4	2.8 ± 1.1
SPOTCUPUS ¹	—	—	97.5 ± 1.6	17.2 ± 0.8	82.8 ± 0.8	2.5 ± 1.6
SCAMPi (MSA)	—	—	94.2 ± 1.6	5.6 ± 0.3	94.4 ± 0.3	5.8 ± 1.6
CCTOP ²	—	—	96.1 ± 2.1	3.7 ± 0.6	96.3 ± 0.6	3.9 ± 2.1
HMM-TM (MSA) ³	—	—	97.3 ± 1.6	21.4 ± 0.5	78.6 ± 0.5	2.7 ± 1.6
BOCTOPUS ²	84.0 ± 13.3	4.2 ± 0.5	—	—	95.8 ± 0.5	16.0 ± 13.3
BetAware-Deep	85.1 ± 9.3	4.7 ± 0.3	—	—	95.3 ± 0.3	14.9 ± 9.3
PRED-TMBB ⁴	88.8 ± 12.1	7.1 ± 0.4	—	—	92.9 ± 0.4	11.2 ± 12.1
PROfmb	91.9 ± 9.0	6.1 ± 0.5	—	—	93.9 ± 0.5	8.1 ± 9.0

<https://github.com/BernhoferM/TMbed?tab=readme-ov-file>

BIOS477/877 L26 - 41

41

Transmembrane protein databases

UniTmp

Topology Data Bank of Transmembrane Proteins

Human Transmembrane Proteome

TOPDB

TOPDOM

PDBTM

PDBBank

MemProtMD

<https://www.unitmp.org>

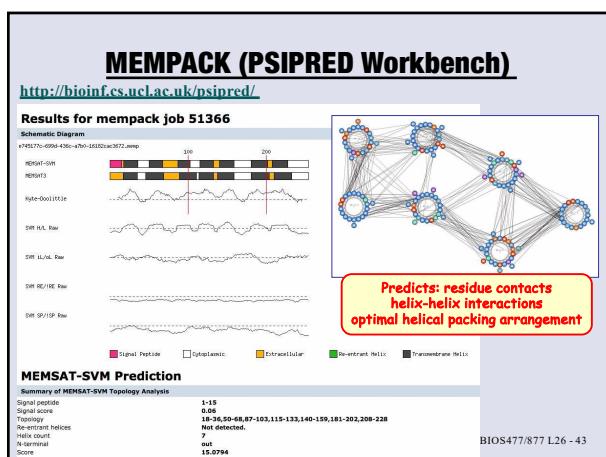
<http://memprotmd.biochem.ox.ac.uk>

A database of membrane proteins embedded in lipid bilayers

More reviewed in: Sun et al. (2023)

BIOS477/877 L26 - 42

7



43