

Spring 2024
BIOS 477/877
Bioinformatics and Molecular Evolution
Lecture 25

BIOS477/877 L25 - 1

1

TODAY'S TOPICS

- Protein structures
 - Protein structure databases (PDB, SCOP, CATH)
- Secondary structure prediction
 - Prediction statistics

BIOS477/877 L25 - 2

2

Prediction of protein structures

BIOS477/877 L24 - 3

3

Hierarchical nature of protein architecture

- **Primary structure:** amino acid sequence
- **Secondary structure:** alpha helices, beta sheet, etc.
 (hydrogen bonding pattern of main chain):
- **Tertiary structure:** (the assembly and interactions of the helices and sheets)
- **Quaternary structure:** (the assembly of the monomers)

BIOS477/877 L25 - 4

4

Proteins = polypeptide chains

Amino acid

BIOS477/877 L25 - 5

5

Proteins = polypeptide chains

- **Backbone (main chain):** atoms that participate in peptide bonds
 → ignores the side chains
 → 6 atoms (O, C, N, H, and C_α) of the peptide group lie in a plane

Backbone is a linked sequence of rigid planar peptide group

BIOS477/877 L25 - 6

6

Torsion angles

- The conformation of the backbone can be described by the torsion angles (or dihedral angles or rotation angles): ϕ and ψ
- The entire path of the backbone is known if ϕ and ψ are specified
- Some values of ϕ and ψ are more likely than others
- Due to steric interference between non-bonded atoms

BIOS477/877 L25 - 7

7

Secondary structure

➤ **Alpha helix:** $\phi = -60$ and $\psi = -45$

Hydrogen bonds stabilize the alpha helix

BIOS477/877 L25 - 8

8

Secondary structure

➤ **Beta sheets**

BIOS477/877 L25 - 9

9

Secondary structure

➤ **Beta turns**

- Proline and glycine occur frequently in beta turns

BIOS477/877 L25 - 10

10

Ramachandran plot

- The sterically allowed values for ϕ and ψ can be calculated (colored regions)
- α : right-handed alpha helix
- α_L : left-handed alpha helix
- ↑↑: parallel beta sheet
- ↑↓: antiparallel beta sheet
- C: collagen helix

BIOS477/877 L25 - 11

11

Secondary structure conformation

π -helix Right-handed α -helix 3_{10} -helix

	ϕ	ψ	H-bond pattern
Right handed α -helix	-57	-47	i + 4
π -helix	-57	-70	i + 5
3_{10} helix	-49	-26	i + 3
Parallel β -sheet	-119	113	
Antiparallel β -sheet	-139	135	

BIOS477/877 L25 - 12

12

SCOP: Structural Classification of Proteins

<http://scop.berkeley.edu/> (SCOPe: SCOP extended)

SCOPe

<http://scop.mrc-lmb.cam.ac.uk/> (SCOP 2, not available?)



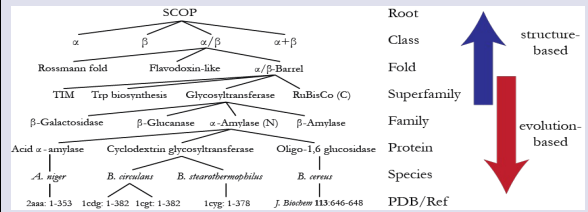
➤ Classification of protein domains of known structure according to their evolutionary and structural relationships

- **Protein domain:** fundamental unit of classification
 - an evolutionary unit observed in nature either in isolation or in more than one context in multidomain proteins
 - hierarchically classified into:
 - families, superfamilies, folds, and classes
- (in SCOP2, classification is network-like but not hierarchical)

BIOS477/877 L25 - 19

19

SCOP: Structural Classification of Proteins

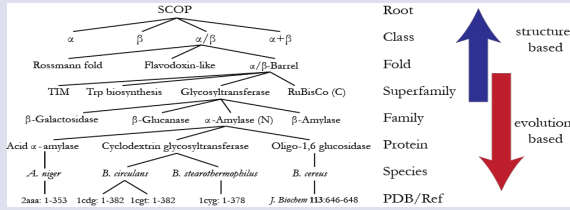


- **Family:** Proteins having a common evolutionary origin
 - Amino acid identities of 30 % and greater (detectable using BLAST, PSI-BLAST, HMMER, etc.)
- **Superfamily:** Families having a probable common evolutionary origin
 - Lower identities but whose structures and functional features suggest common evolutionary origin

BIOS477/877 L25 - 20

20

SCOP: Structural Classification of Proteins



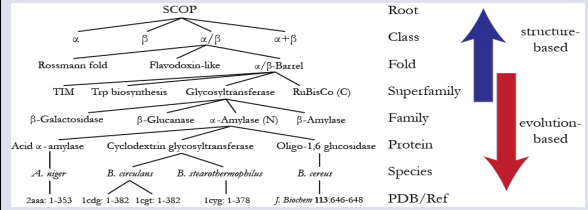
- **Fold:** Same major secondary structures in same arrangement with the same topological connections



BIOS477/877 L25 - 21

21

SCOP: Structural Classification of Proteins



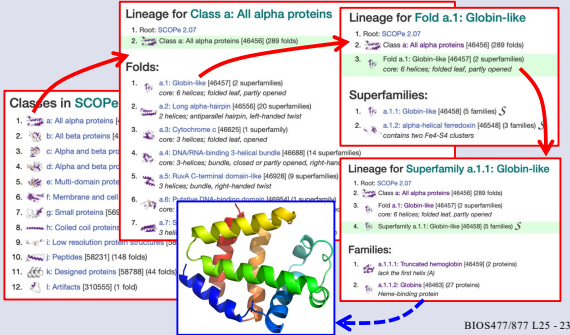
- **IUPRs (Intrinsically Unstructured Protein Region):** non-globular structures (IUPRs and Types are introduced in SCOP 2)
 - Folds and IUPRs are grouped into 5 structural Classes:
 - 1) all α , 2) all β , 3) α and β (α/β); interspersed α and β ,
 - 4) α plus β ($\alpha+\beta$); segregated α and β), and 5) small proteins
 - Folds and IUPRs are also grouped based on 4 protein Types: soluble, membrane, fibrous, or intrinsically disordered

BIOS477/877 L25 - 22

22

SCOP: Structural Classification of Proteins

<http://scop.berkeley.edu/>



23

SCOP: Structural Classification of Proteins

"The SCOP classification of proteins has been constructed manually by visual inspection and comparison of structures, but with the assistance of tools"

SCOPe 2.08-stable statistics: 106976 PDB entries (released/updated prior to 2021-07-28), 344851 Domains, 1 Literature reference.

Class	Number of folds	Number of superfamilies	Number of families
a: All alpha proteins	290	519	1089
b: All beta proteins	180	375	993
c: Alpha and beta proteins (a/b)	148	247	1003
d: Alpha and beta proteins (a+b)	396	580	1387
e: Multi-domain proteins (alpha and beta)	74	74	128
f: Membrane and cell surface proteins and peptides	69	131	204
g: Small proteins	100	141	280
Totals	1257 (26 new)	2067 (42 new)	5084 (88 new)

BIOS477/877 L25 - 24

24

CATH Protein Structure Classification

<http://www.cathdb.info/>

25

CATH Protein Structure Classification

- **Class, C-level:** determined according to the secondary structure composition and packing within the structure
→ mainly- α , mainly- β , α - β , and few secondary structures
- **Architecture, A-level:** the overall shape of the domain structure
→ determined by the orientations of the secondary structures
→ done by manually
- **Topology (Fold family), T-level:** depending on both the overall shape and connectivity of the secondary structures
→ done by the structure comparison algorithm (SSAP)
- **Homologous Superfamily, H-level:** depending on sequence similarity and structural comparison
→ sequence identity $\geq 20 \sim 35\%$
- **Sequence families, S-level, or Domain:** based on sequence similarity

BIOS477/877 L25 - 26

26

CATH Protein Structure Classification

<http://www.cathdb.info/>

27

CATH Protein Structure Classification

CATH Superfamily 1.50.10.100
Chondroitin AC/alginase lyase

28

Secondary structure prediction

29

Secondary structure prediction

➤ Various amino acids stabilize or destabilize alpha helix

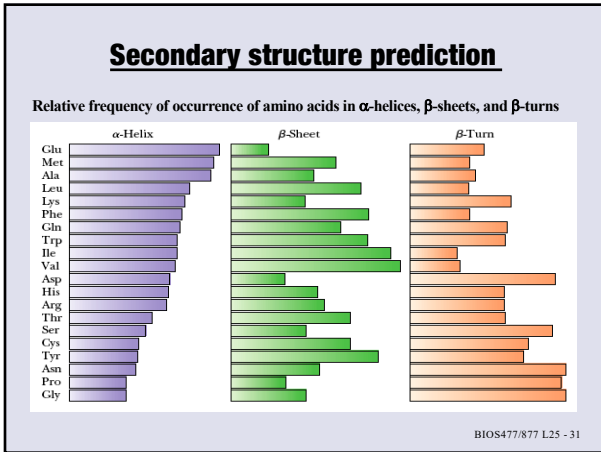
Table 6.1 Helix-Forming and Helix-Breaking Behavior of the Amino Acids		
Amino Acid	Helix Behavior*	
A	Ala	H (I)
C	Cys	Variable
D	Asp	Variable
E	Glu	H
F	Phe	H
G	Gly	I (B)
H	His	H (I)
I	Ile	H (C)
K	Lys	Variable
L	Leu	H
M	Met	H
N	Asn	C (I)
P	Pro	B
Q	Gln	H (I)
R	Arg	H (I)
S	Ser	C (B)
T	Thr	Variable
V	Val	Variable
W	Tyr	H (C)
Y	Tyr	H (C)

*H = helix former; I = indifferent; B = helix breaker; C = random coil; () = secondary tendency.

H: helix forming amino acid

Proline: helix breaker

30



31

Secondary structure prediction

➤ **Chou-Fasman method** Chou and Fasman (1974)

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

- The first generation, simple statistical method
- Based on the property (propensity) of a single residue
- Propensities are calculated from a set of proteins whose structures are known

→ Propensity for α -helix, β -sheet, or turn
 $= \log\{\text{Pr}(\text{amino acid } i \text{ in } \alpha\text{-helix, } \beta\text{-strand, or turn}) / \text{Pr}(\text{amino acid } i)\}$

BIOS477/877 L25 - 32

32

Secondary structure prediction

• Conformational propensity table used in Chou-Fasman method

Name	P(a)	P(b)	P(turn)	f(1)	f(1+1)	f(1+2)	f(1+3)
Ala	142	83	66	0.06	0.076	0.035	0.058
Arg	98	93	95	0.070	0.106	0.099	0.085
Asp	101	54	146	0.147	0.110	0.179	0.081
Asn	67	89	156	0.161	0.083	0.191	0.091
Cys	70	119	119	0.149	0.050	0.117	0.128
Glu	151	37	74	0.056	0.060	0.077	0.064
Gln	111	110	98	0.074	0.098	0.037	0.098
Gly	57	75	156	0.102	0.085	0.190	0.152
His	100	87	95	0.140	0.047	0.093	0.054
Ile	108	160	47	0.043	0.034	0.013	0.056
Leu	121	130	59	0.061	0.025	0.036	0.070
Lys	114	74	101	0.055	0.115	0.072	0.095
Met	145	105	60	0.068	0.082	0.014	0.055
Phe	113	138	60	0.059	0.041	0.065	0.065
Pro	57	55	152	0.102	0.301	0.034	0.068
Ser	77	75	143	0.120	0.139	0.125	0.106
Thr	83	119	96	0.086	0.108	0.065	0.079
Trp	108	137	96	0.077	0.013	0.064	0.167
Tyr	69	147	114	0.082	0.065	0.114	0.125
Val	106	170	50	0.062	0.048	0.028	0.053

BIOS477/877 L25 - 33

33

Secondary structure prediction

➤ **Chou-Fasman method** Chou and Fasman (1974)

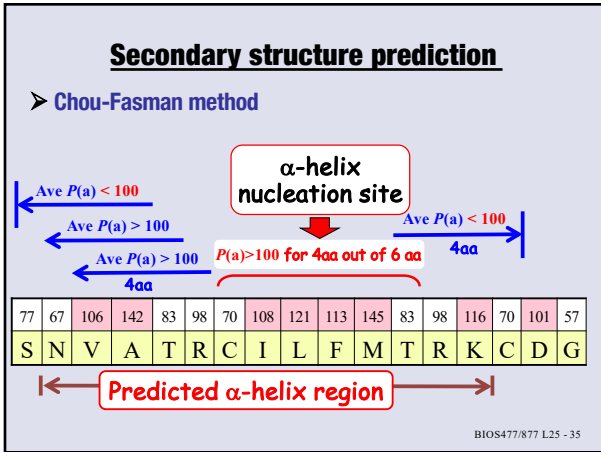
http://fasta.bioch.virginia.edu/fasta_www2/fasta_www.cgi?rm=misc1

1. Identification of helix and sheet nuclei:
 helix: 4 out of 6 aa w/ high helix propensity ($P > 100$)
 strand: 3 out of 5 aa w/ high strand propensity ($P > 100$)
2. Propagation in both directions until 4-aa Average ($P < 100$)
3. Helix if $\Sigma P(a) > \Sigma P(b)$, strand if $\Sigma P(b) > \Sigma P(a)$

→ Turn prediction:
 i) $p(t) > 0.000075$, ii) $P(\text{turn}) > 100$, and
 iii) $P(a) < P(\text{turn}) > P(b)$, where $p(t) = f(j)f(j+1)f(j+2)f(j+3)$

BIOS477/877 L25 - 34

34



35

Secondary structure prediction

➤ **Garnier-Osguthorpe-Robson (GOR)** Garnier *et al.* (1978)

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

- Likelihood of a secondary structure state depends on the neighboring residue
- Chou-Fasman method assumes that each amino acid individually influences secondary structure
- Uses principles of information theory Garnier *et al.* (1996)
- Algorithm
 - For each residue, 8 N-terminal and 8 C-terminal positions are considered (a window size = 17)
 - The frequencies with which each amino acid occupies each of the 17 window positions in helices, strands, and turns are compiled (a 17 x 20 scoring matrix) e.g., $P(\text{Ala at pos 1} | \alpha\text{-helix}) \rightarrow \text{Lod score}$
 - The probabilities that each residue in a target sequence will be involved in a helix, strand, or turn are calculated

BIOS477/877 L25 - 36

36

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

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

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)	Seq 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%	✓	✓	
JPREIQ 2008 (Cole et al., 2007)	81.5%	✓	✓	
SPIDER3 2017 (Jefferson et al., 2017)	84%	✓	✓	✓

Accuracy Achieved by NN Methods

Evolutionary or = MSA profile information

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

BIOS477/877 L25 - 55

55

Improving prediction of protein secondary structure, backbone angles, solvent accessibility and contact numbers by using predicted contact maps and an ensemble of recurrent and residual convolutional neural networks

Jack Hanson^{1,*}, Kuldip Paliwal¹, Thomas Litfin², Yuedong Yang³ and Yaoqi Zhou^{2,4,*}

Bioinformatics, 35, 2019, 2403-10

Table 2. Test performance of several recently developed predictors alongside SPOT-1D on the latest PDB structures (TEST2018) (250 proteins)

Predictor	SS3	SOV3	PSS3	SS8	SOV8	PSS8	ASA	HSE-U	HSE-D	CN	θ	τ	ϕ	ψ
SPIDER-3-Single	72.57	64.08	$<1 \times 10^{-10}$	59.81	57.86	$<1 \times 10^{-10}$	0.570	0.603	0.533	0.619	11.07	45.39	23.77	43.05
RaptorX-DeepCN	61.62	66.58	$<1 \times 10^{-10}$	70.43	65.66	$<1 \times 10^{-10}$	-	-	-	-	-	-	-	21.01 35.95
PSRSM	81.94	74.22	$<1 \times 10^{-10}$	-	-	-	-	-	-	-	-	-	-	-
SPIDER-3	83.84	73.89	$<1 \times 10^{-10}$	-	-	-	0.768	0.764	0.716	-	7.73	29.62	18.38	28.10
PORTER-5	84.10	74.04	$<1 \times 10^{-10}$	73.22	70.27	$<9.89 \times 10^{-9}$	-	-	-	-	-	-	-	-
MUFOLD	84.78	77.56	$<2.73 \times 10^{-8}$	73.66	71.34	$<2.15 \times 10^{-9}$	-	-	-	-	-	-	-	17.78 27.24
NetSurfP-2.0	85.31	78.58	$<2.20 \times 10^{-3}$	73.81	71.14	$<3.64 \times 10^{-4}$	0.801	-	-	-	-	-	-	17.90 26.63
SPOT-1D-base	85.66	78.77	$<1.08 \times 10^{-2}$	74.26	71.45	$<1.33 \times 10^{-4}$	0.799	0.812	0.775	0.837	7.03	26.86	17.15	25.41
SPOT-1D	86.18	79.00	-	75.41	73.30	-	0.803	0.814	0.779	0.841	6.91	25.94	16.89	24.87

SS3=Q3, SS8=Q8
 ASA: Solvent accessible surface area
 CN: Contact numbers
 HSE: Half-sphere exposure

Backbone torsion angles (mean absolute error)

BIOS477/877 L25 - 56

56

Increasing the accuracy of single sequence prediction methods using a deep semi-supervised learning framework

Lewis Moffat^{1,2,*} and David T. Jones^{1,2,*}

Bioinformatics, 37, 2021, 3744-51

Uses bidirectional gated recurrent unit (GRU) with 1024 hidden dimensions in each direction

S4PRED (= next version of PSIPRED-single) achieved 75% accuracy without using profiles

BIOS477/877 L25 - 57

57

Reaching alignment-profile-based accuracy in predicting protein secondary and tertiary structural properties without alignment

Jaspreet Singh¹, Kuldip Paliwal², Thomas Litfin³, Jaswinder Singh⁴ & Yaoqi Zhou⁵

Scientific Reports 12, Article number: 7607 (2022) | Cite this article

Model	SS3	SS8	ASA	HSE-u	HSE-d	CN	ψ	ϕ	θ	τ
SPIDER3-Single	72.57	59.81	0.647	0.523	0.487	0.547	43.05	23.78	11.07	45.38
ProteinUnet	72.57	60.30	0.620	0.537	0.510	0.545	42.93	23.42	10.28	44.94
SPOT-1D-Single	74.28	72.17	0.665	0.573	0.563	0.585	40.58	22.16	9.35	42.32
NetSurfP-2.0(profile)*	85.35	73.48	0.783	-	-	-	26.63	17.90	-	-
SPOT-1D (profile)	86.18	75.41	0.787	0.732	0.737	0.777	24.87	16.88	6.91	25.94
SPOT-1D-LM (This work)	86.74	76.47	0.814	0.759	0.761	0.690	23.74	15.99	6.46	24.60

No MSA

*NetSurfP-3.0: no MSA, LM + deep learning is also now available

SS3=Q3, SS8=Q8
 ASA: Solvent accessible surface area
 HSE: Half-sphere exposure; CN: Contact numbers
 LM: language model (natural language processing)

Backbone torsion angles (mean absolute error)

BIOS477/877 L25 - 58

58

Secondary structure prediction

► Prediction program servers

- PSIPRED Server <http://bioinf.cs.ucl.ac.uk/psipred/>
- MPI Bioinformatics Toolkits (Quick2D) <https://toolkit.tuebingen.mpg.de>
- PredictProtein Server <http://www.predictprotein.org/>

[Recent review papers (on Canvas)]

Meng and Kurgan (2016); Jiang et al. (2017); Oldfield et al. (2019); Torris et al. (2020); Ismi et al. (2022)

BIOS477/877 L25 - 59

59

Secondary structure prediction

MPI Bioinformatics Toolkits (Quick2D) <https://toolkit.tuebingen.mpg.de>

Protein ID: AAN59974.1 histone H2A (Homo sapiens)

SS = α -helix β -strand; τ -helix CC = Coiled Coils TM = Transmembrane DO = Disorder

BIOS477/877 L25 - 60

60