

# Structural Bioinformatics

## Lecture 4

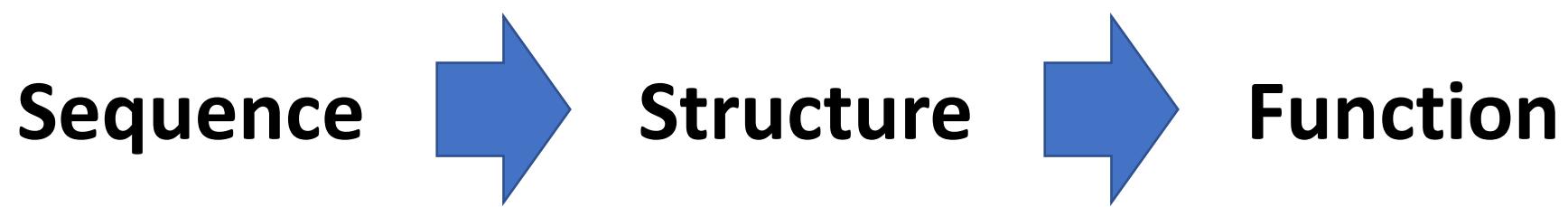
Prediction of structural features from sequence



UNIVERSITÄT  
DES  
SAARLANDES



# Proteins



# Outline

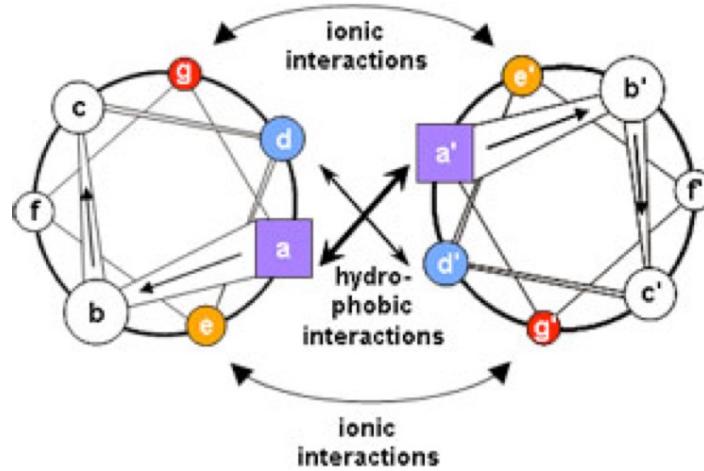
- Coiled-coil regions
- Membrane proteins
- Prediction of membrane-embedded segments
- Prediction of secondary structure

# Predicting coiled-coil regions

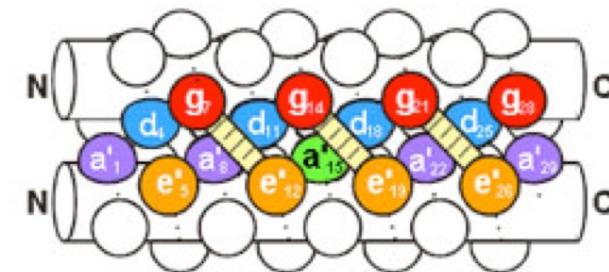
# What are coiled-coils?

- 2 to 7  $\alpha$ -helices coiled together (**super-secondary structure**)

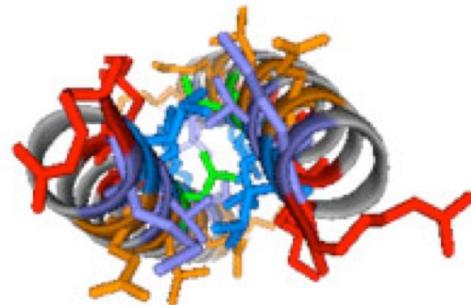
A



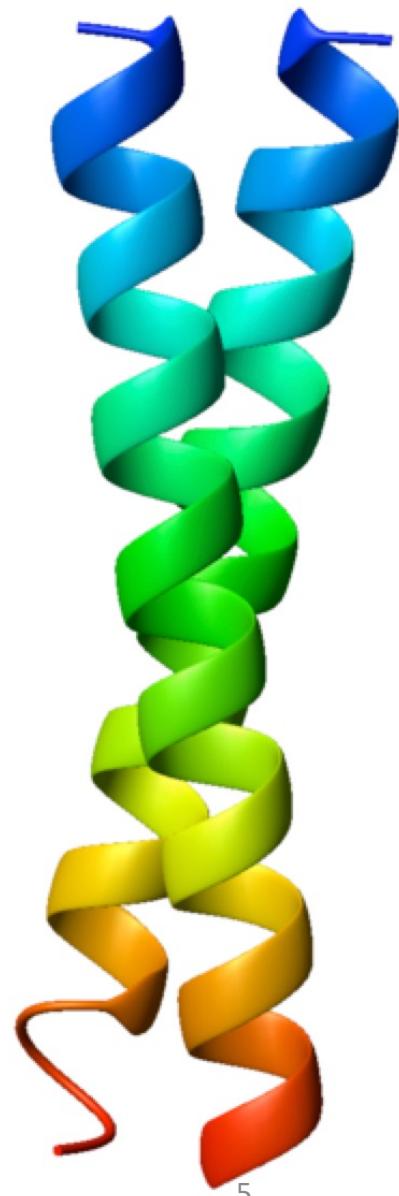
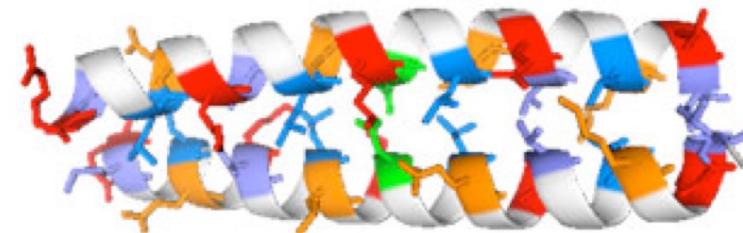
B



C

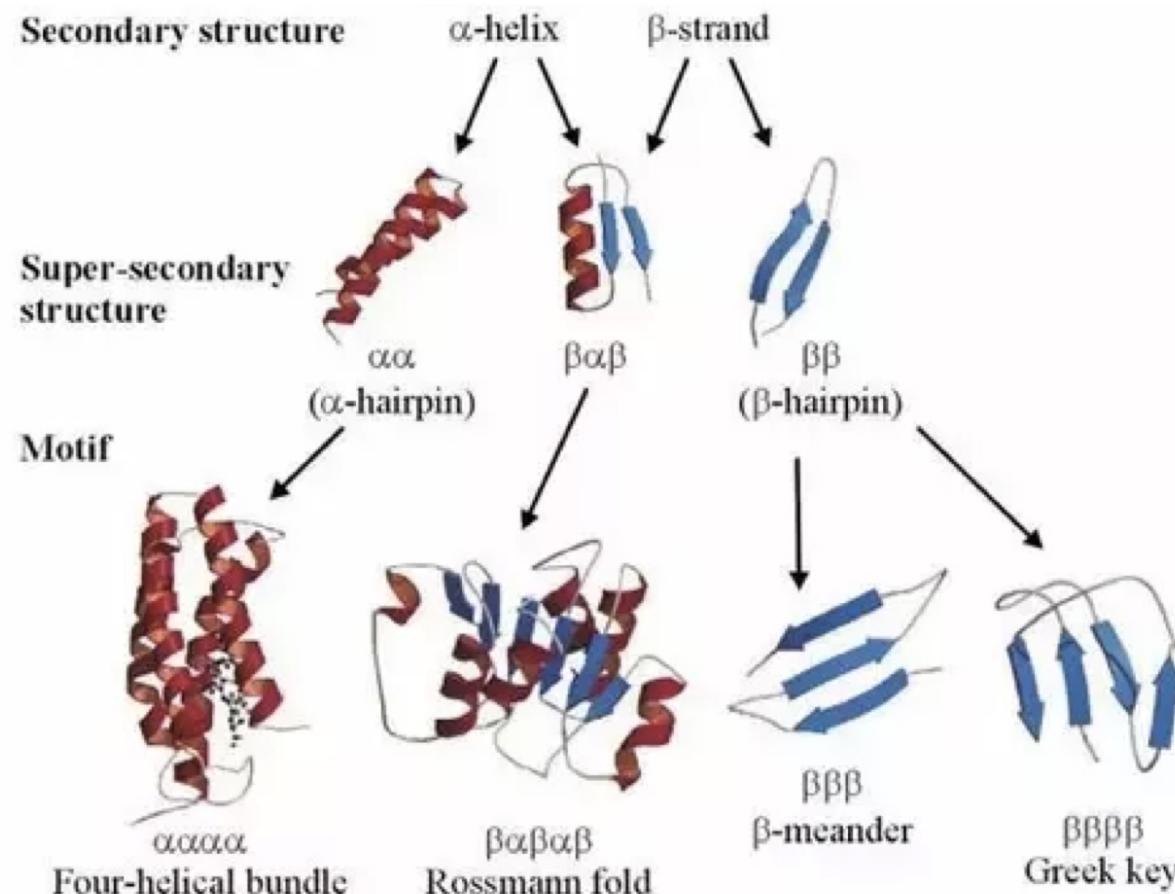


D

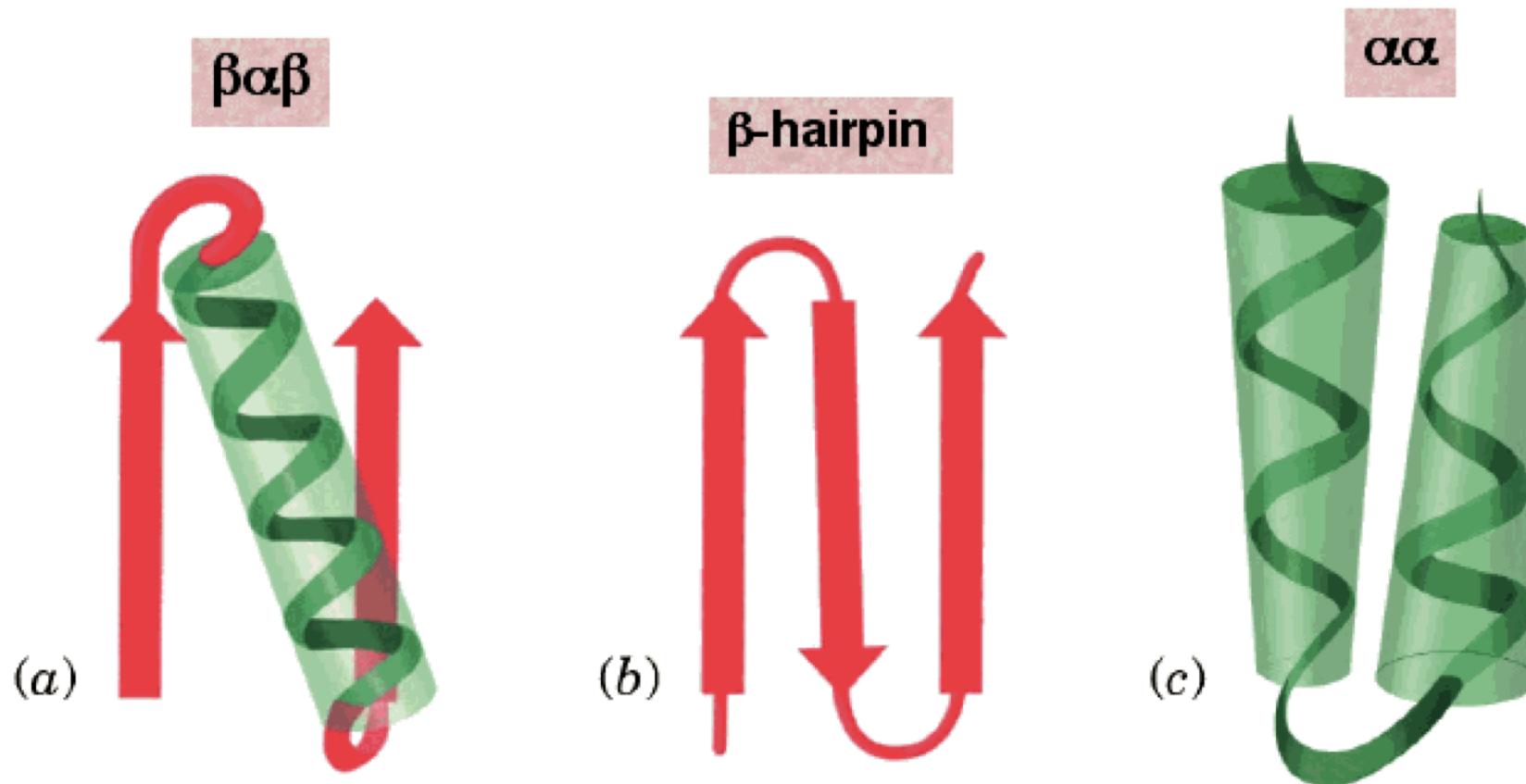


# Detour: super-secondary structure

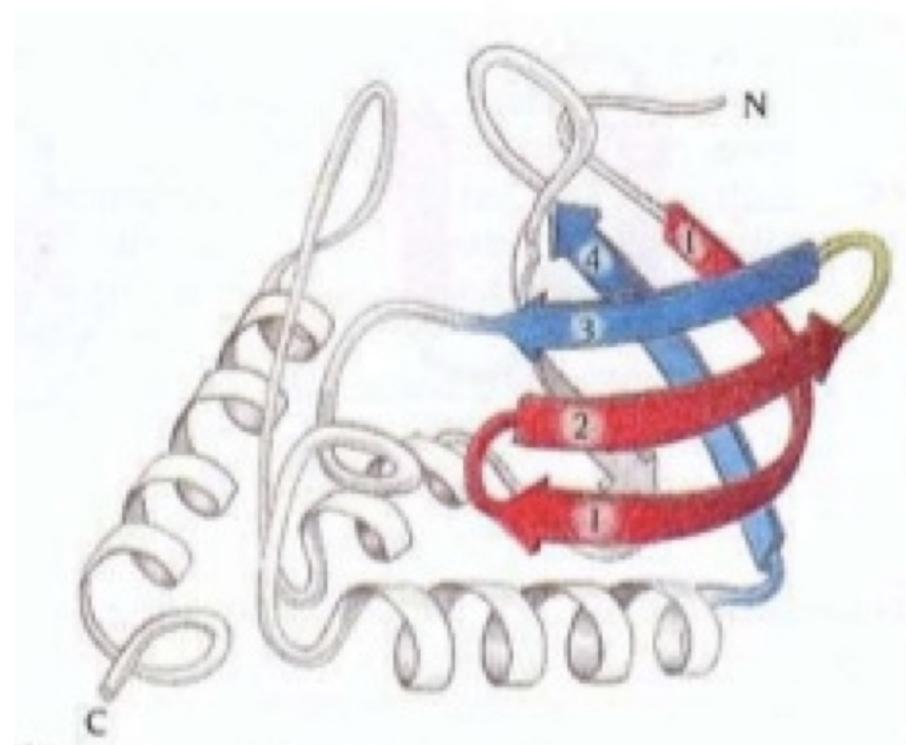
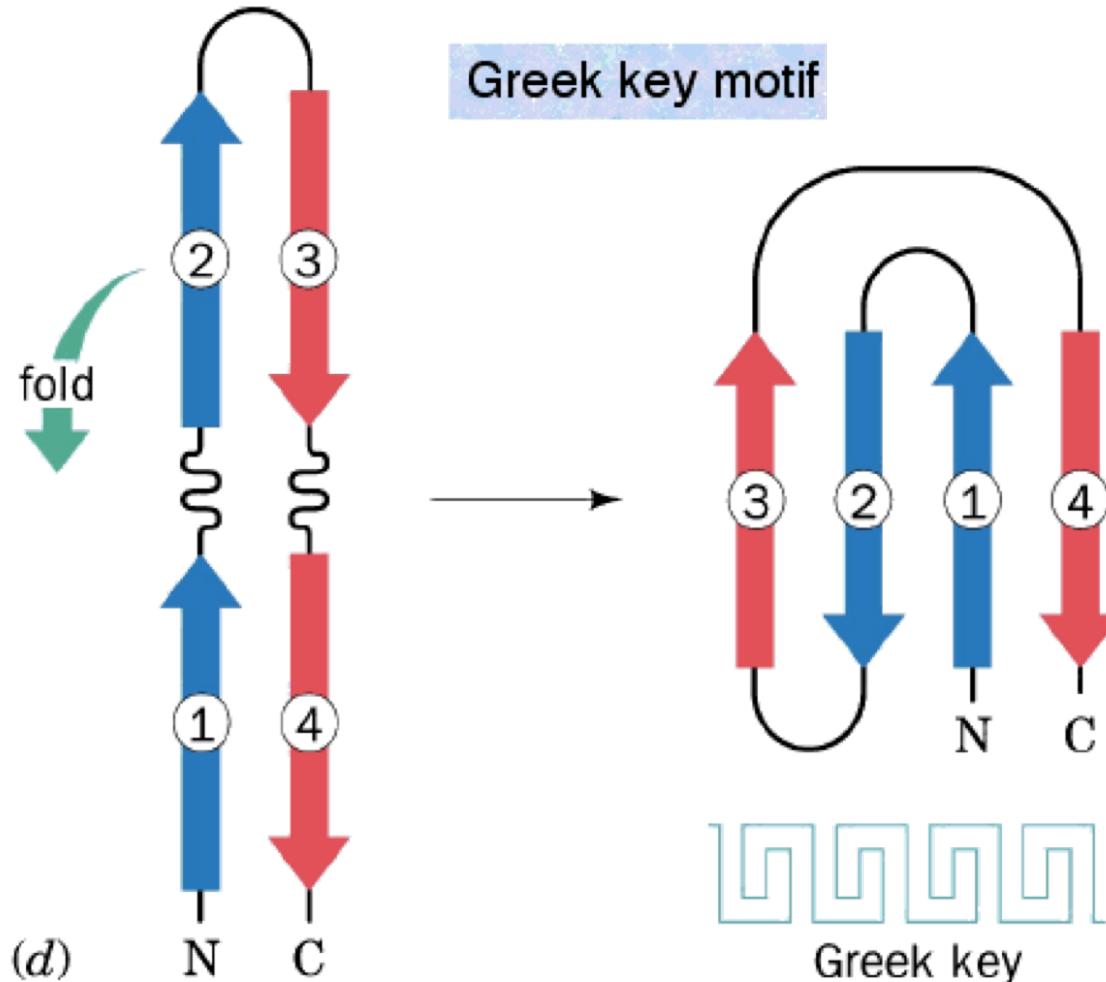
- Halfway between secondary and tertiary structure: spatial arrangement of several secondary structure elements (usually, continuous)



# Super-secondary structure motifs

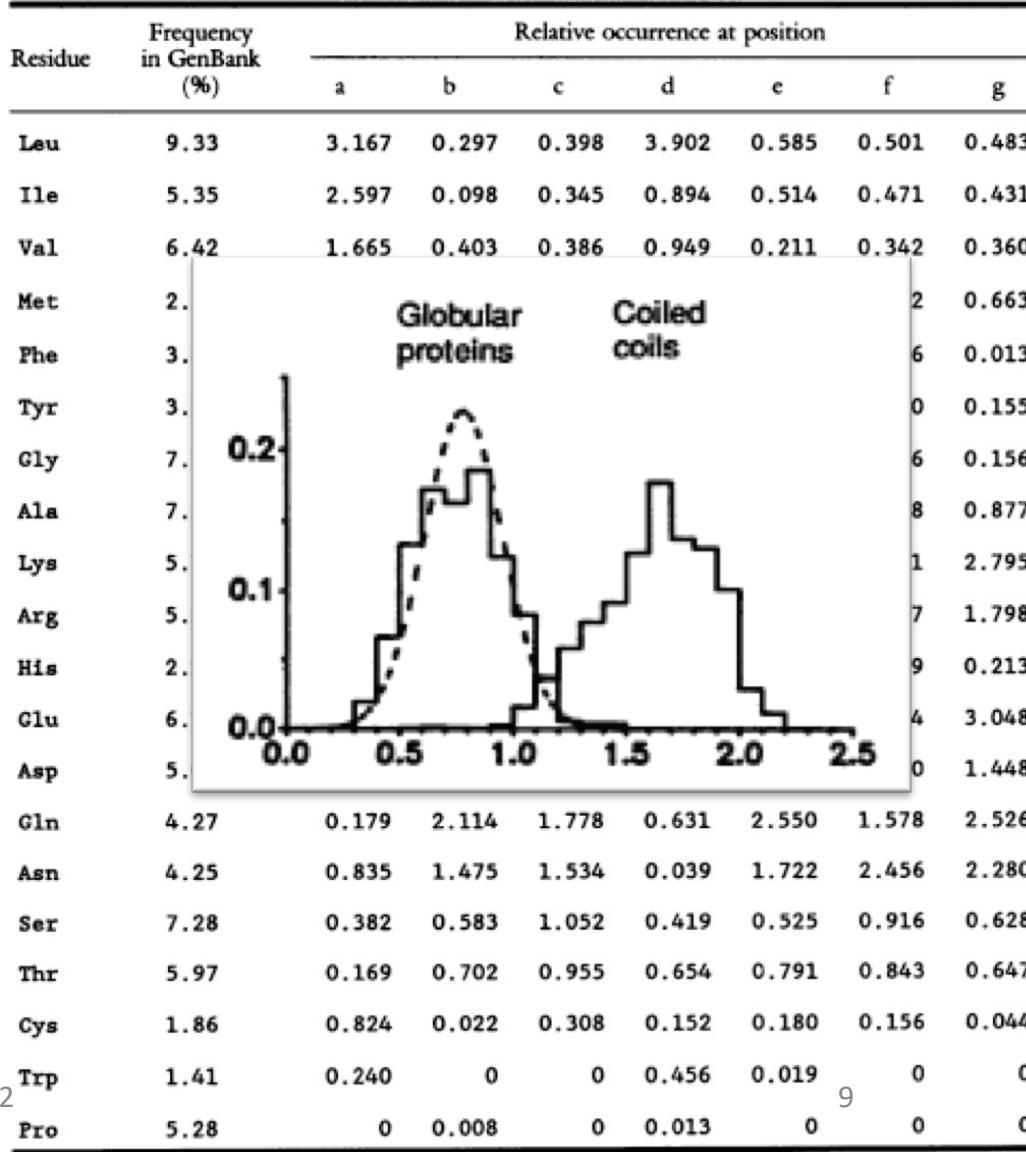
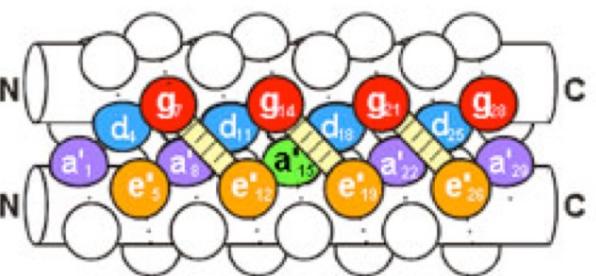
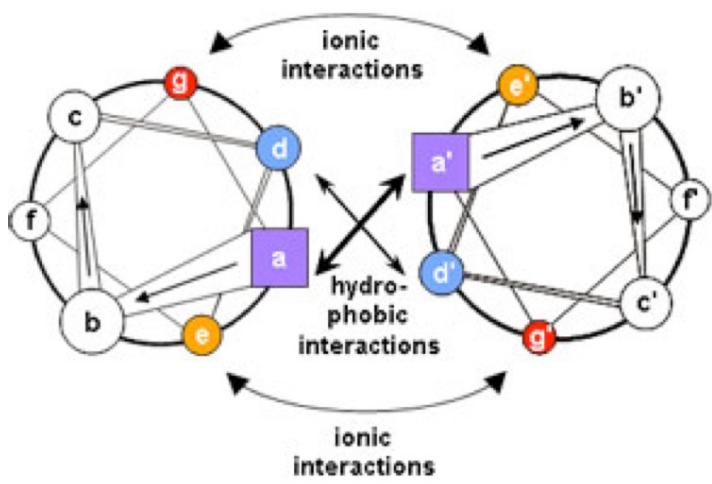


# Super-secondary structure motifs



# Coiled-coil prediction: COILS (Lupas et al., 1991)

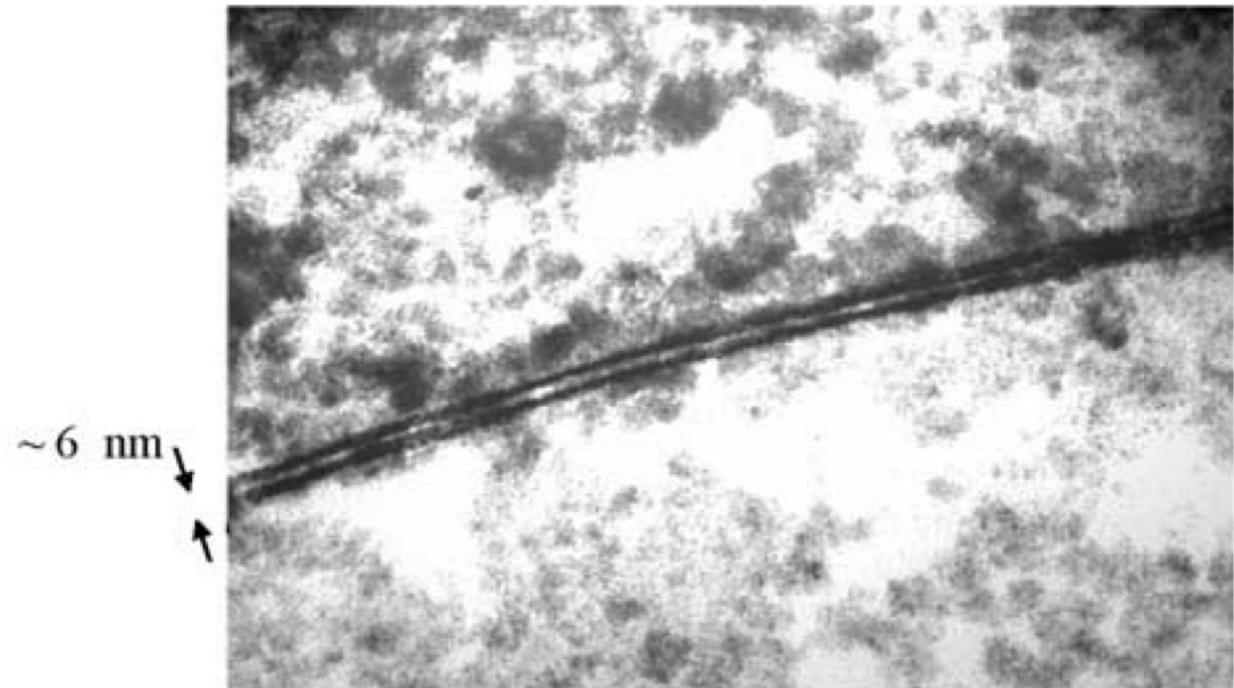
- The 7-residues (**heptad**) pattern is very strong
- A matrix for relative occurrence of all aa. in known coiled-coil regions
- Sliding-window



# Membrane proteins

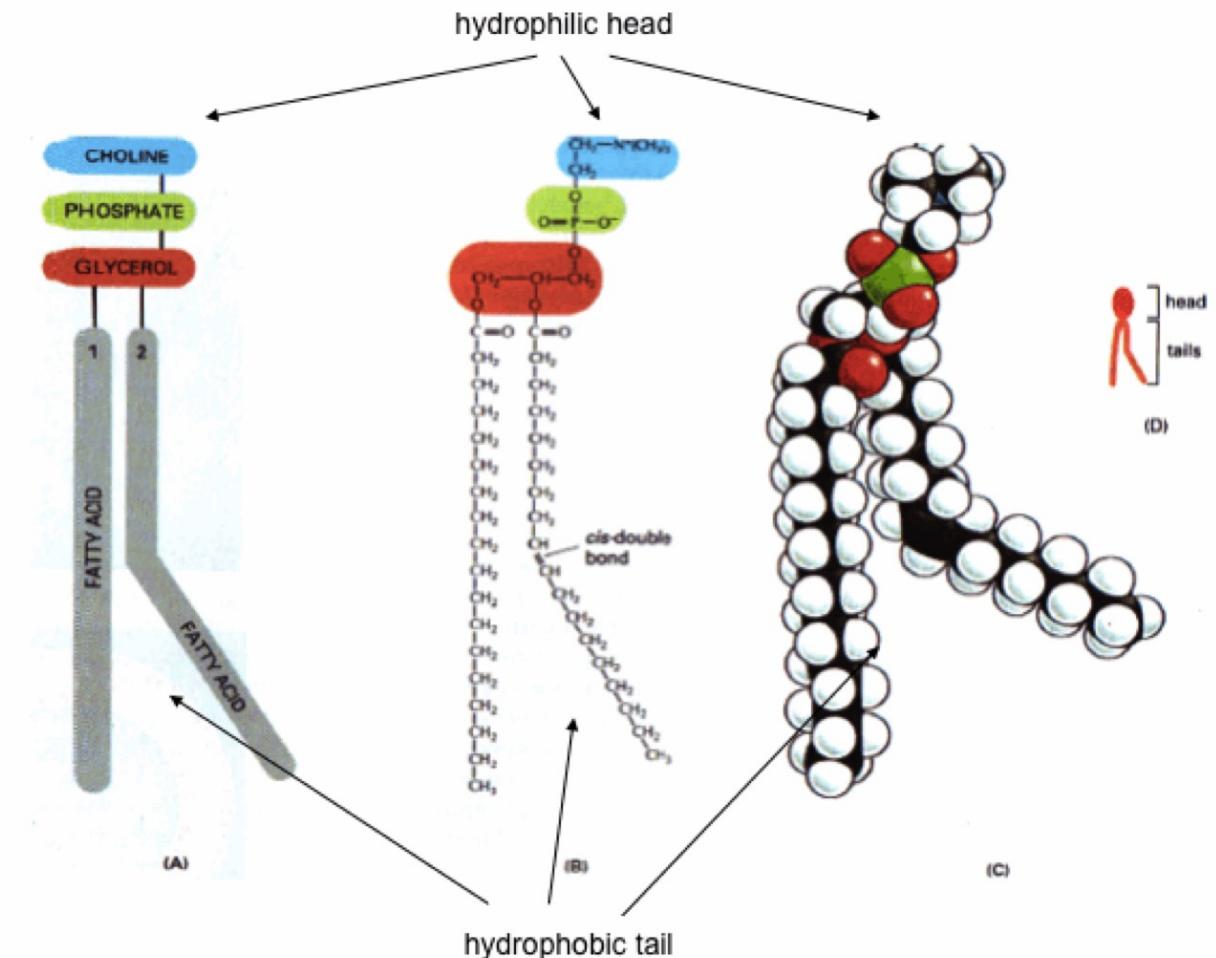
# Biological membranes

- Maintain compartmentalisation
- Cell membranes
  - Plasma membrane: encloses the cell
  - In eukaryotic cells:
    - ER
    - Golgi apparatus
    - Mitochondria
    - Nucleus
    - ...



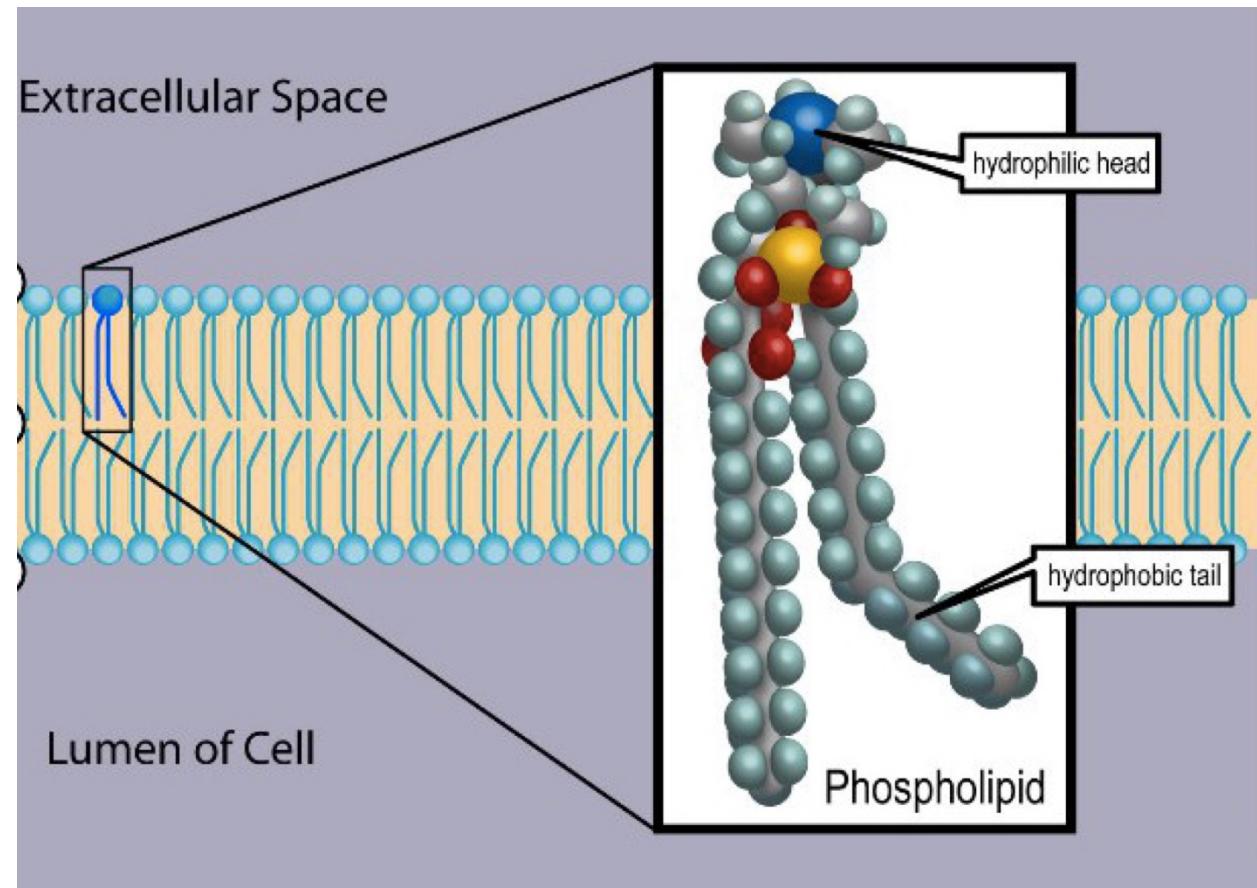
# Lipids

- Lipids are **amphipathic** molecules:  
have a hydrophobic and a  
hydrophilic part

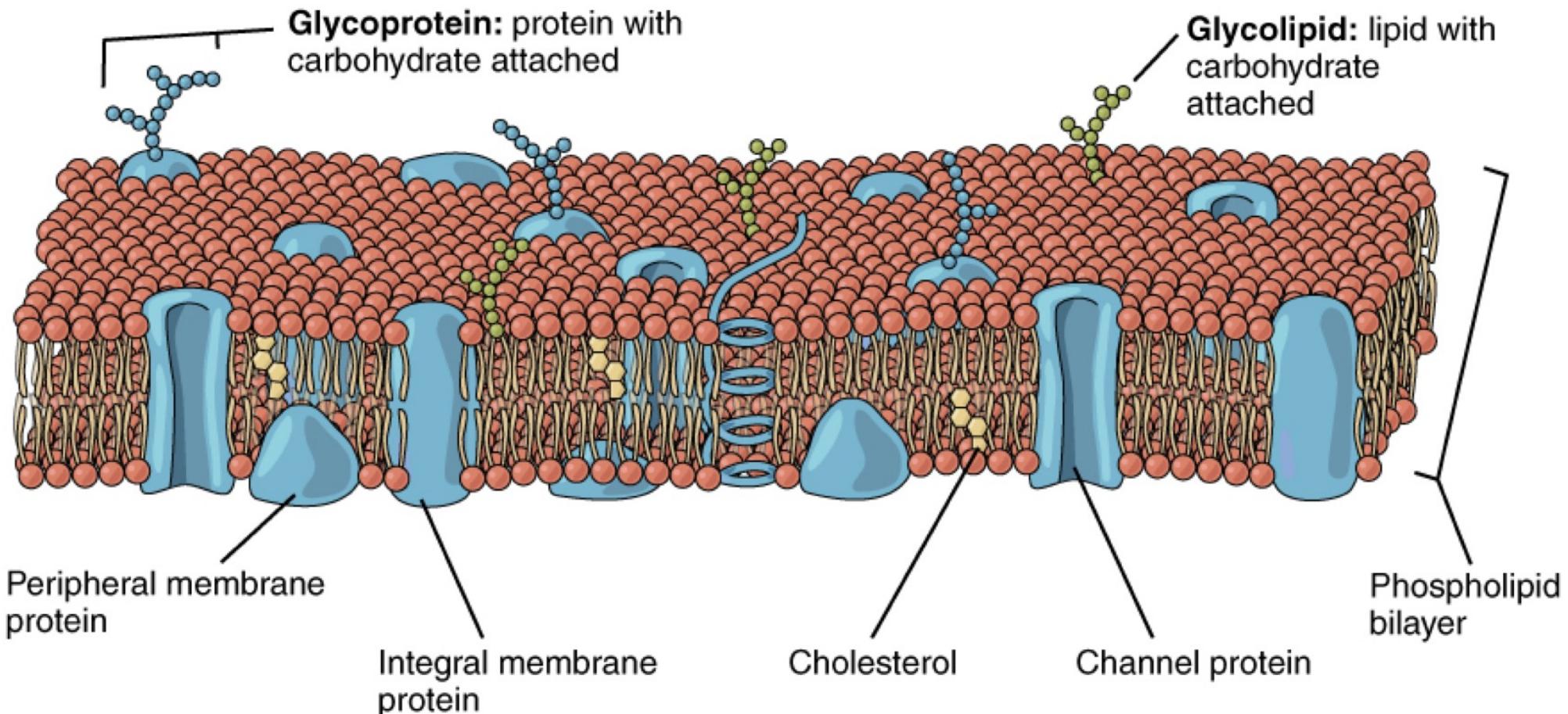


# Lipid bilayer

- Lipids are **amphipathic** molecules:  
have a hydrophobic and a  
hydrophilic part  
=> tend to form bilayers



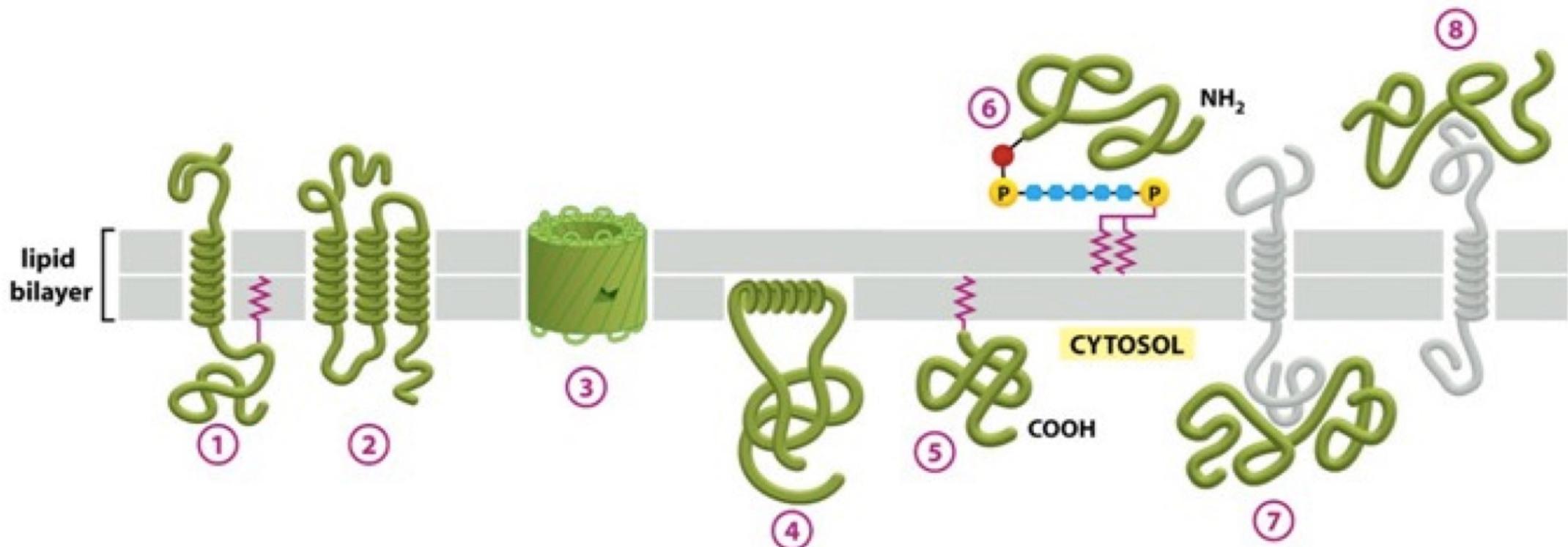
# Lipid bilayer and membrane proteins



**Membrane proteins** can be associated with the lipid bilayer in different ways

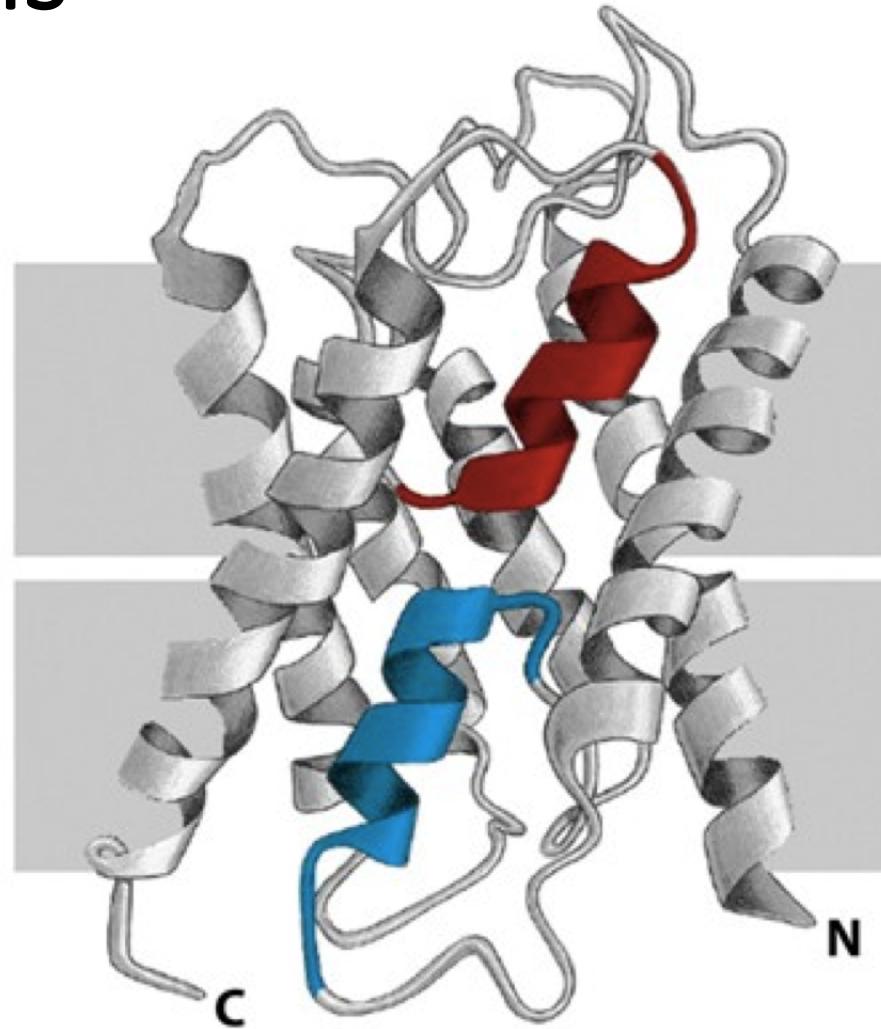
# Membrane proteins

- **Membrane proteins** can be associated with the lipid bilayer in different ways
- Membrane proteins perform most membrane-specific functions

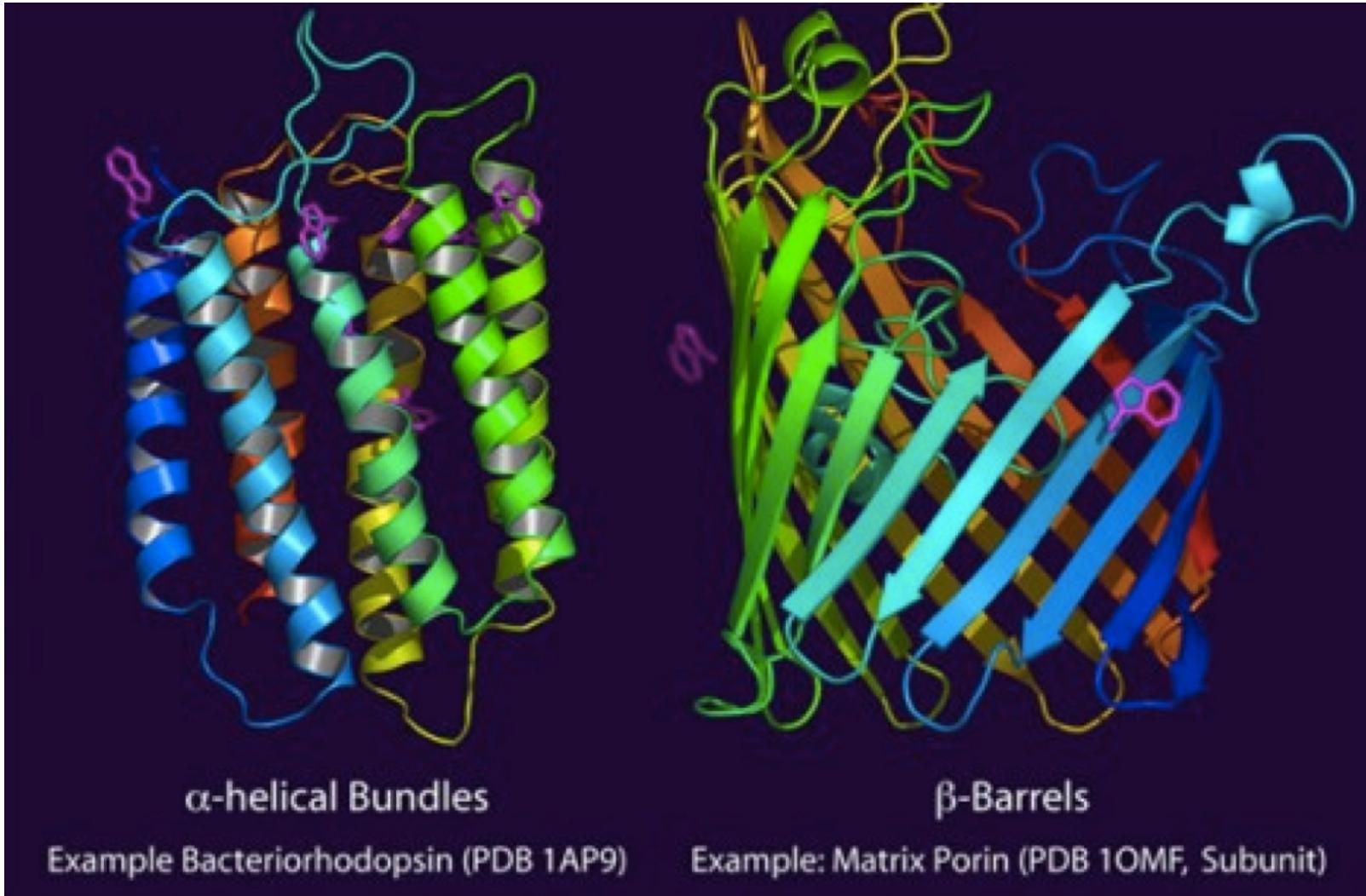


# Transmembrane protein regions

- Some proteins have segments that span the membrane — **transmembrane (TM) regions**
- Proteins can have one (**single-pass, anchored**) or many TM segments (**multi-pass**)
- In multi-pass proteins, TM regions are connected by loops or larger structures
- TM regions are more hydrophobic than extracellular and intracellular loops

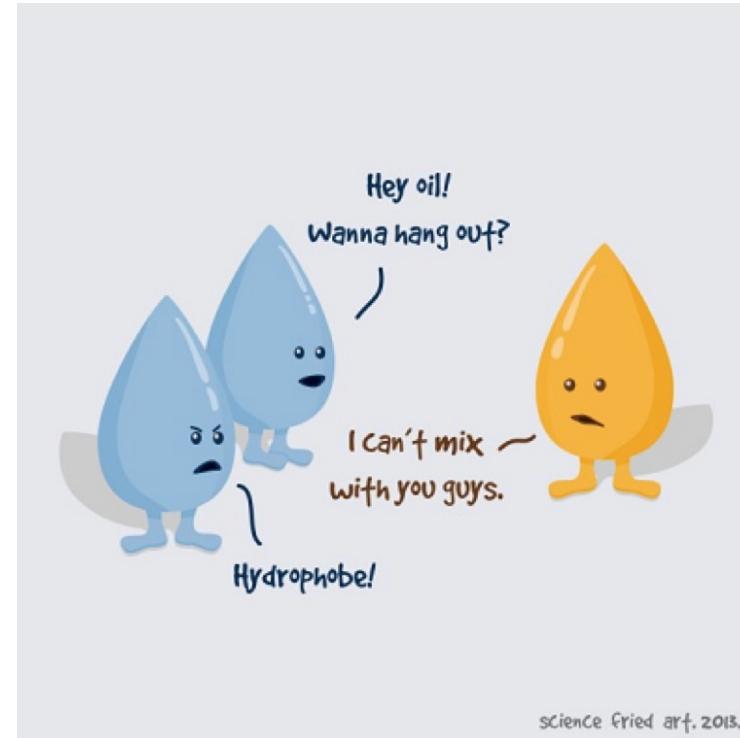


# Two major types of (multi-pass) transmembrane proteins



# Hydrophobicity

- Nonpolar substances tend to aggregate in an aqueous solution and exclude water



# Protein hydrophobic core

- Hydrophobic amino acids tend to be excluded from the protein surface => pack against each other in its core

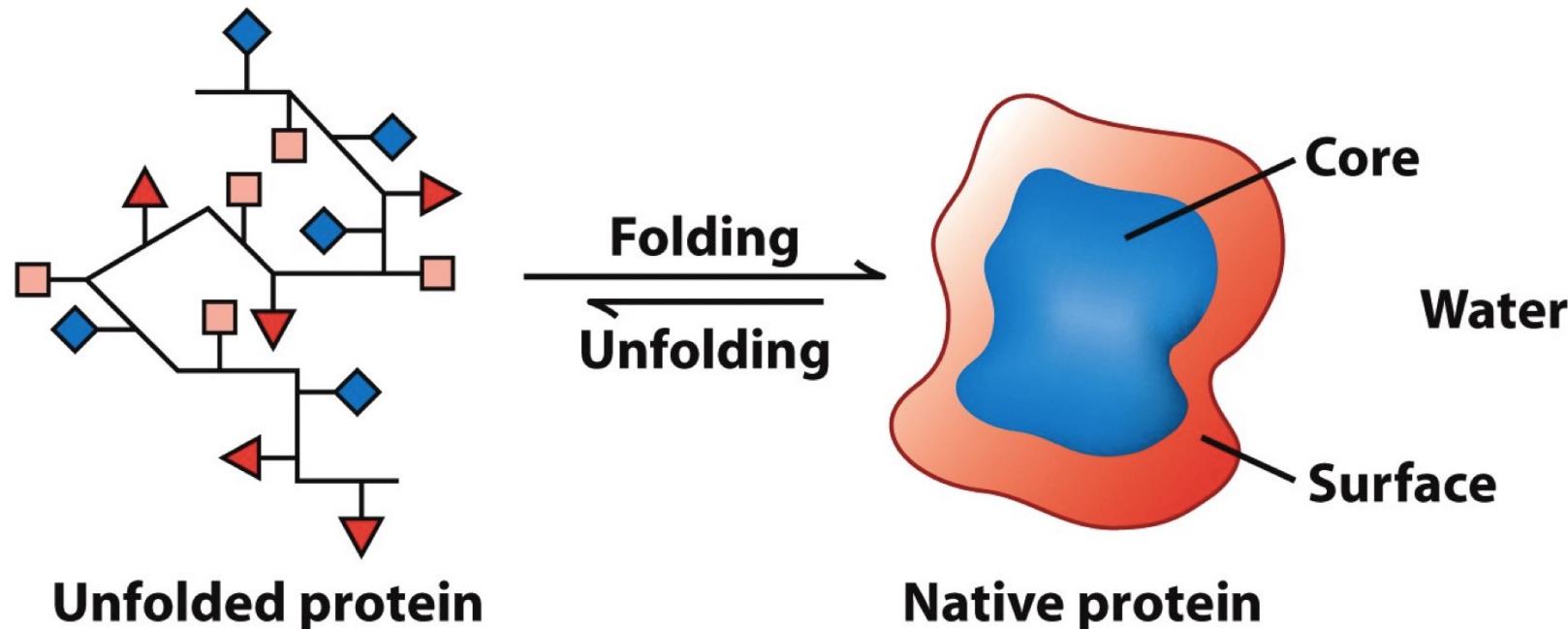


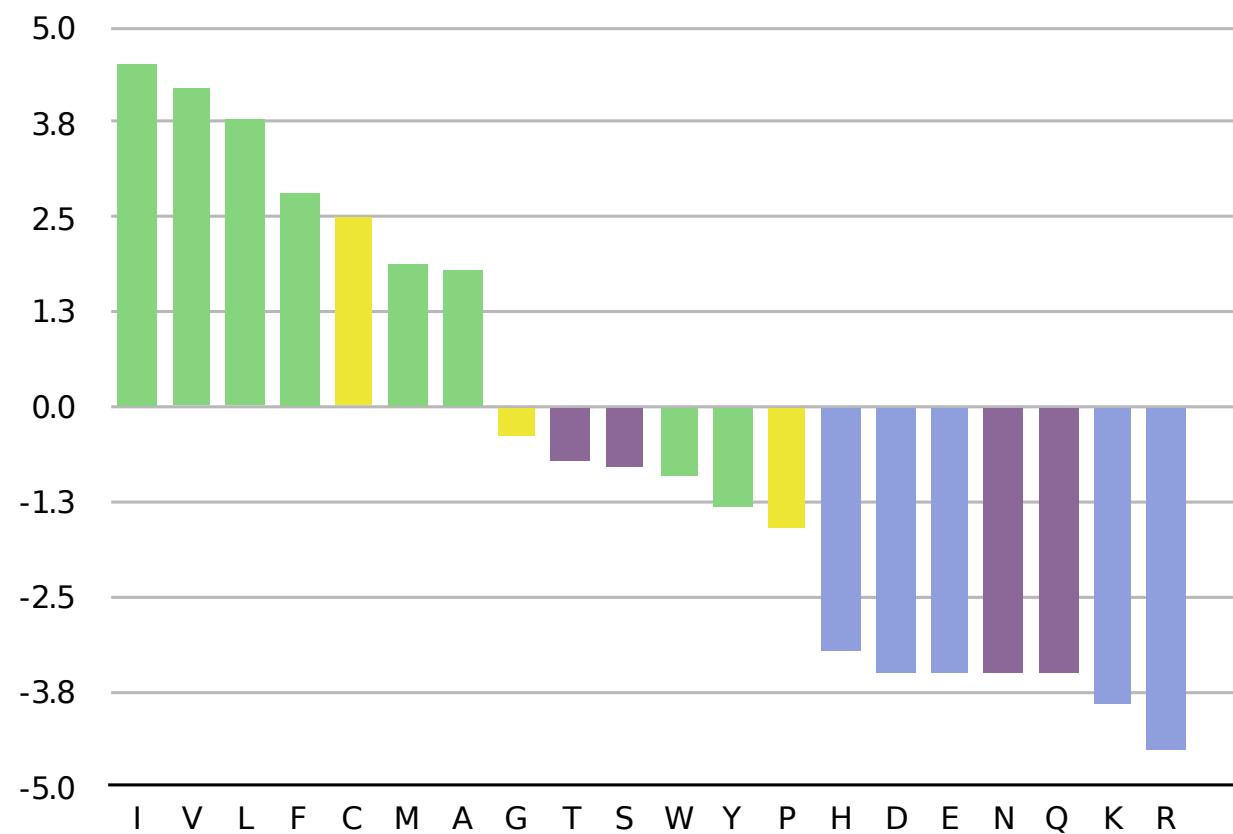
Figure 3-7  
*Molecular Cell Biology, Sixth Edition*  
© 2008 W. H. Freeman and Company

# Hydropathy scale

- Partitioning method:
  - two immiscible liquids (water + something that favours hydrophobic molecules — solvent)
  - measure free energy of transfer of each amino acid between phases
- (Other methods exist)

# Kyte-Doolittle scale (Kyte & Doolittle, 1982)

Side-chain	Hydropathy index	$\Delta G_{\text{transfer}}^{\circ}$ (water-vapor) <sup>a</sup>	Fraction of side-chains 100% buried <sup>b</sup>	Fraction of side-chains 95% buried <sup>c</sup>
Isoleucine	4.5	4.4	4.5	5.2
Valine	4.2	4.2	4.3	4.2
Leucine	3.8	4.5	3.2	2.8
Phenylalanine	2.8	2.5	2.5	3.5
Cysteine/cystine	2.5	1.9	6.0	3.2
Methionine	1.9	1.9	1.0	1.9
Alanine	1.8	3.9	5.3	1.6
Glycine	-0.4	—	4.2	1.3
Threonine	-0.7	-0.6	-0.5	-1.0
Tryptophan	-0.9	-0.9	-2.4	-0.3
Serine	-0.8	-0.8	-0.7	-1.0
Tyrosine	-1.3	-1.1	-3.3	-2.2
Proline	-1.6	—	-2.4	-1.8
Histidine	-3.2	-4.2	-3.6	-1.9
Glutamic acid	-3.5	-3.9	-2.8	-1.7
Glutamine	-3.5	-3.5	-4.0	-3.6
Aspartic acid	-3.5	-4.5	-2.5	-2.3
Asparagine	-3.5	-3.8	-3.1	-2.7
Lysine	-3.9	-3.2	—	-4.2
Arginine	-4.5	—	—	—

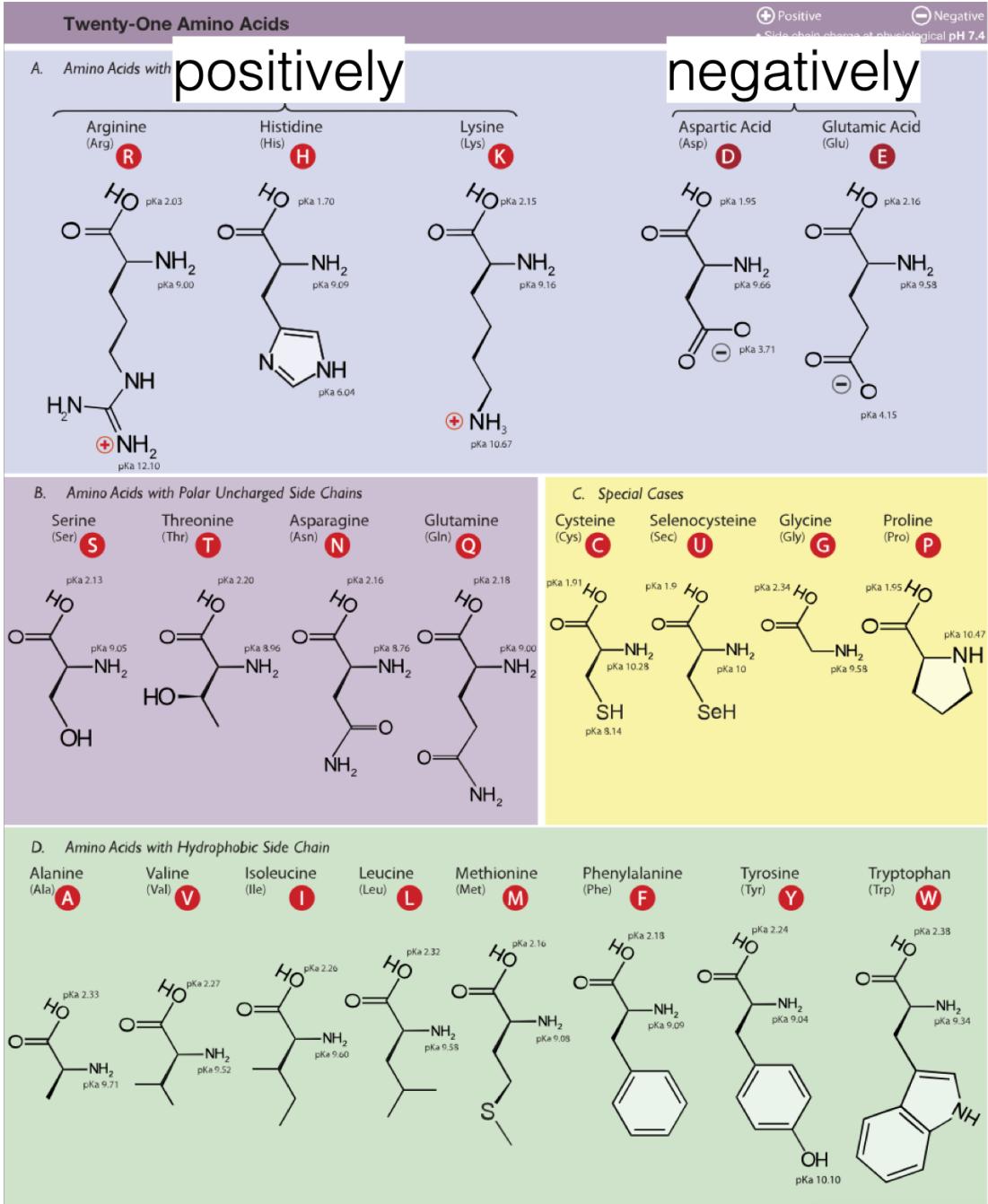


# Amino acids have different chemical properties

charged

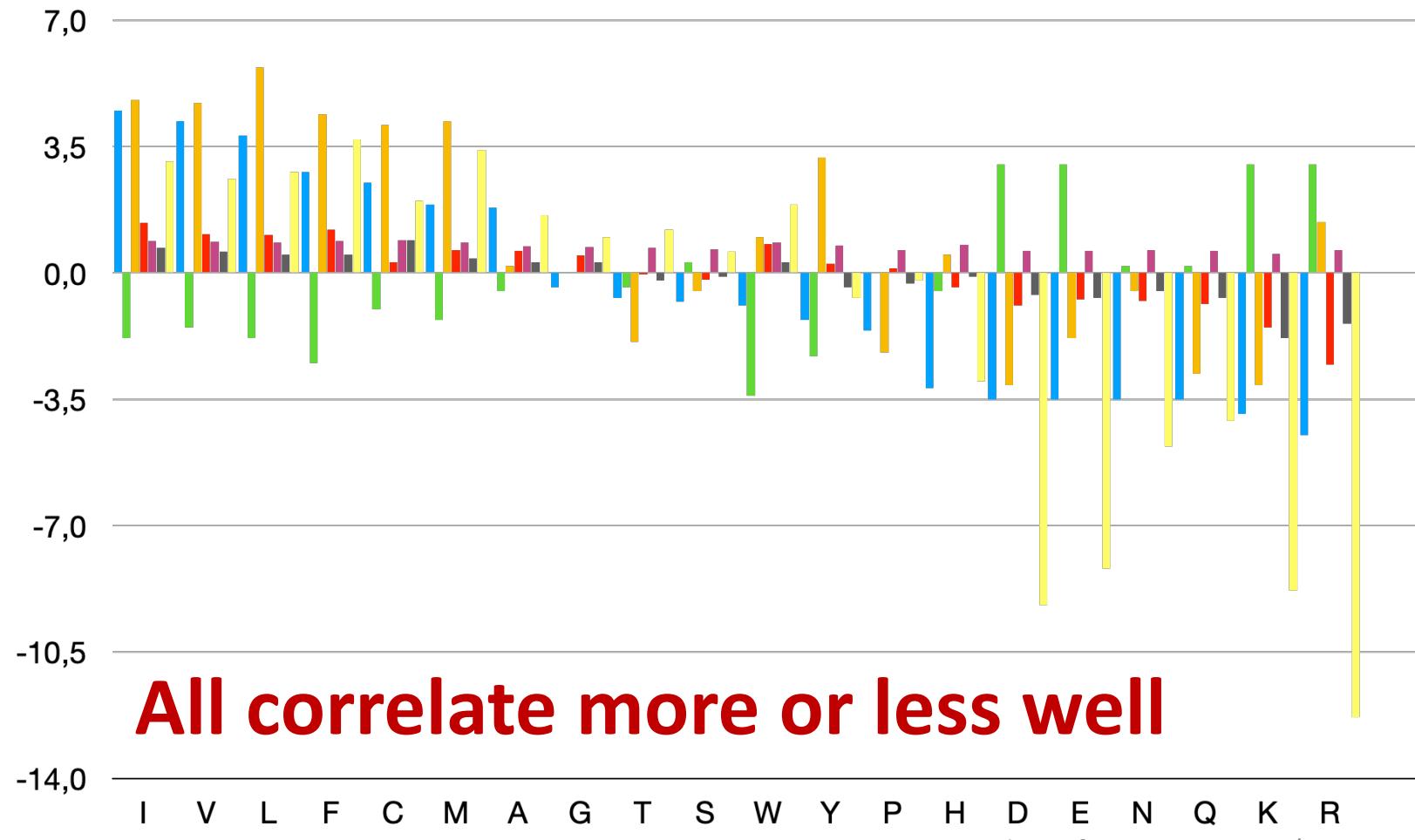
polar

hydrophobic



# Other hydropathy scales

Kyte-Doolittle Hopp-Woods Cornette Eisenberg Rose Janin Engelman GES

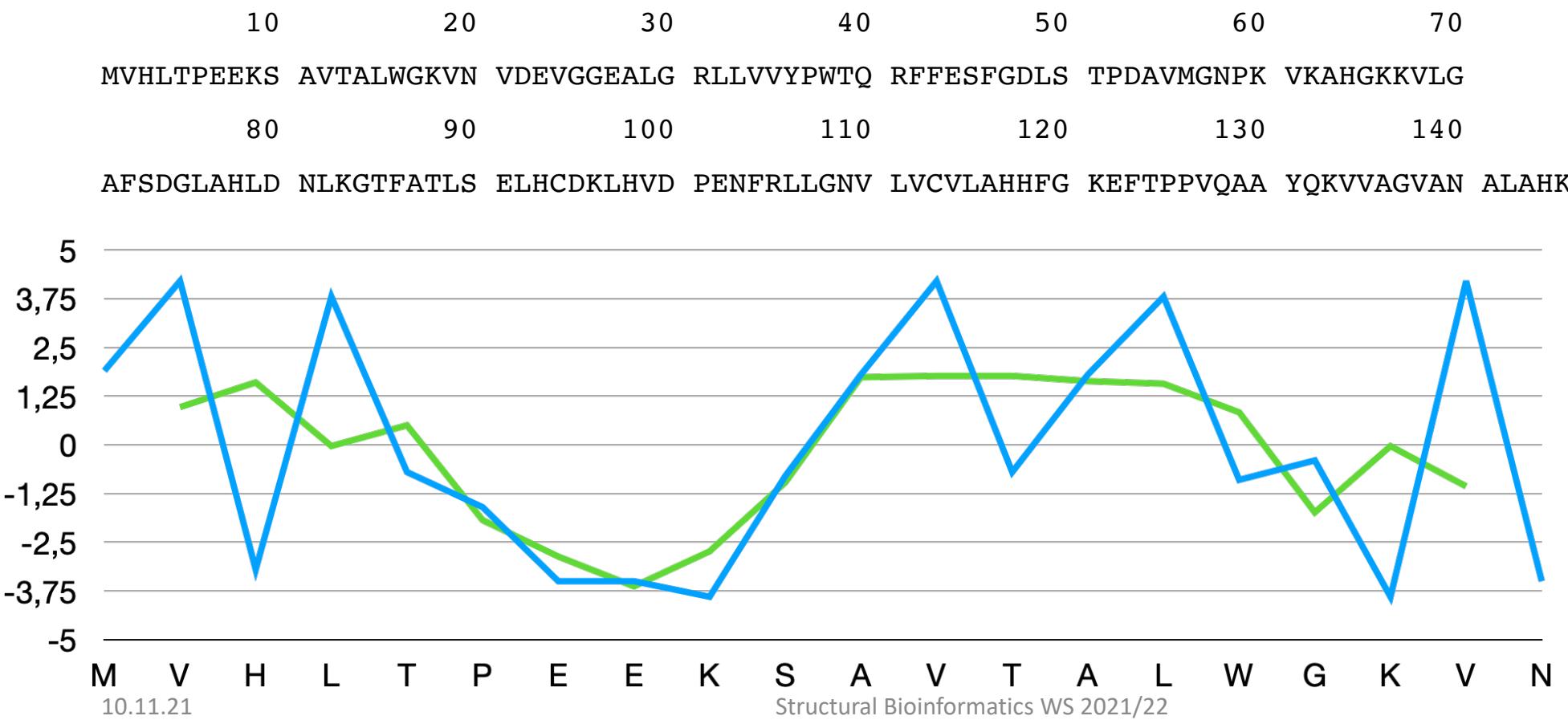


Kyte and Doolittle (1)	Rose, et al (2)	Wolfenden, et al (3)	Janin (1979) (4)
Ile	Cys	Gly,Leu,Ile	Cys
Val	Phe,Ile	Val,ala	Ile
Leu	Val	Phe	Val
Phe	Leu,Met,Trp	Cys	Leu,Phe
Cys		Met	Met
Met,Ala	His	Thr,Ser	Ala,Gly,Trp
Gly	Tyr	Trp,Tyr	
Thr,Ser	Ala		His,Ser
Trp,Tyr	Gly		Thr
Pro	Thr		Pro
His			Tyr
Asn,Gln	Ser		Asn
Asp,Glu	Pro,Arg		Asp
Lys	Asn		Gln,Glu
	Gln,Asp,Glu		
Arg			Arg
	Lys	Arg	Lys

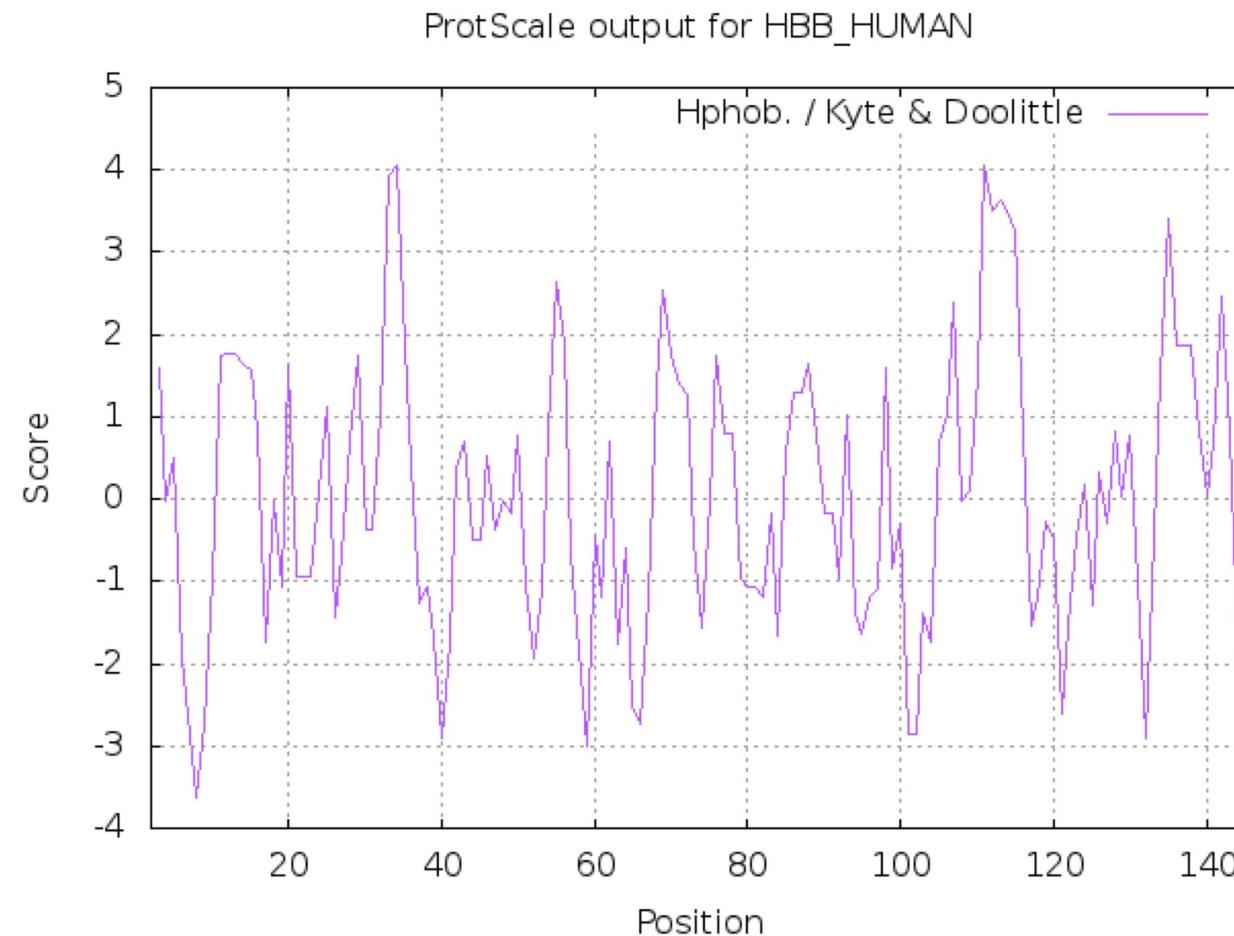
M	1,9	
V	4,2	0,96667
H	-3,2	1,6
L	3,8	-0,03333
T	-0,7	0,5
P	-1,6	-1,93333
E	-3,5	-2,86667
E	-3,5	-3,63333
K	-3,9	-2,73333
S	-0,8	-0,96667
A	1,8	1,73333
V	4,2	1,76667
T	-0,7	1,76667
A	1,8	1,63333
L	3,8	1,56667
W	-0,9	0,83333
G	-0,4	-1,73333
K	-3,9	-0,03333
V	4,2	-1,06667
N	-3,5	24

# Calculating hydrophobicity

- Human haemoglobin

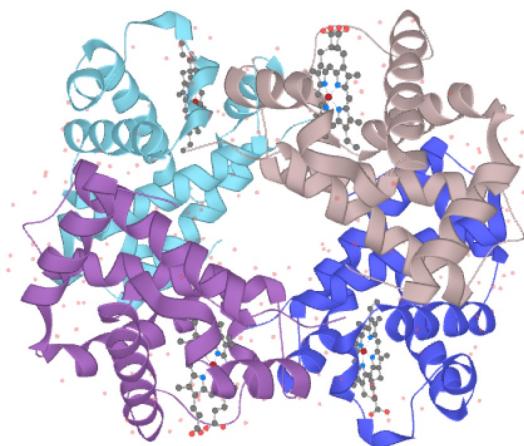
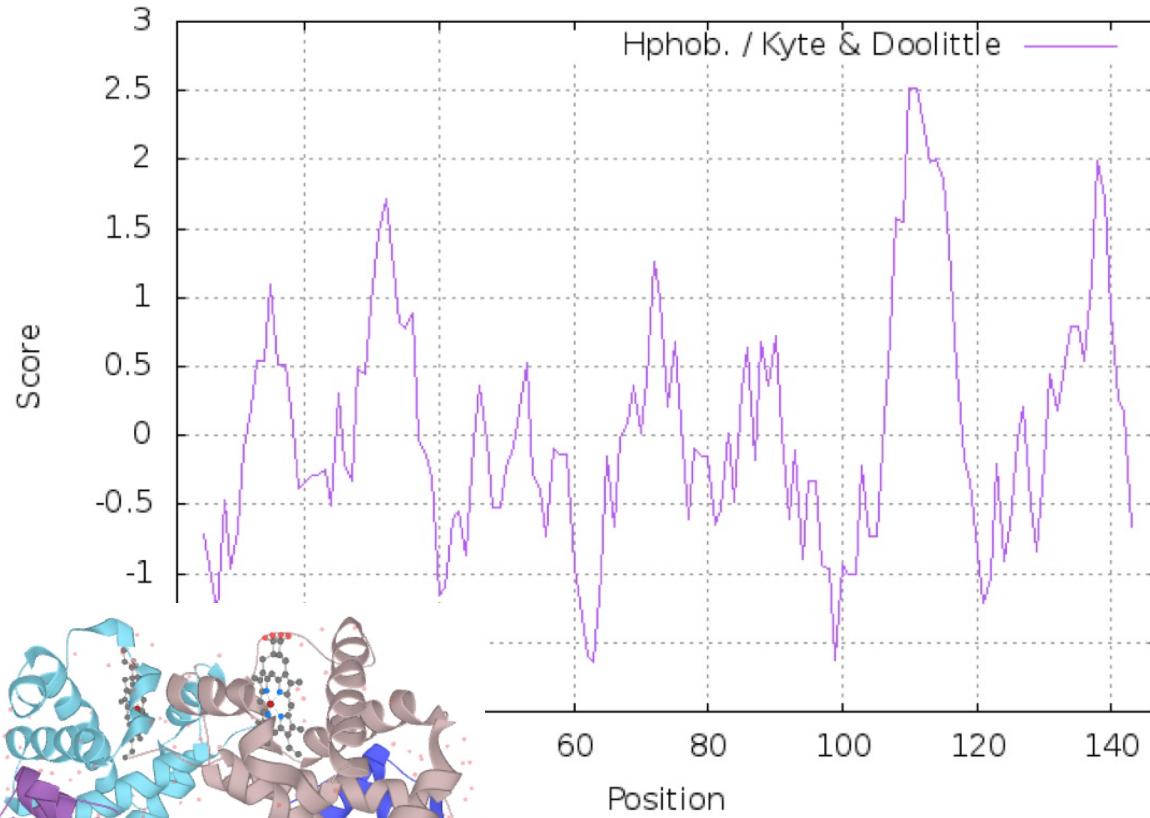


# Calculating hydrophobicity

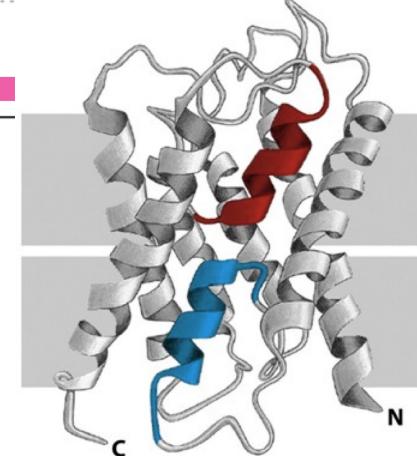
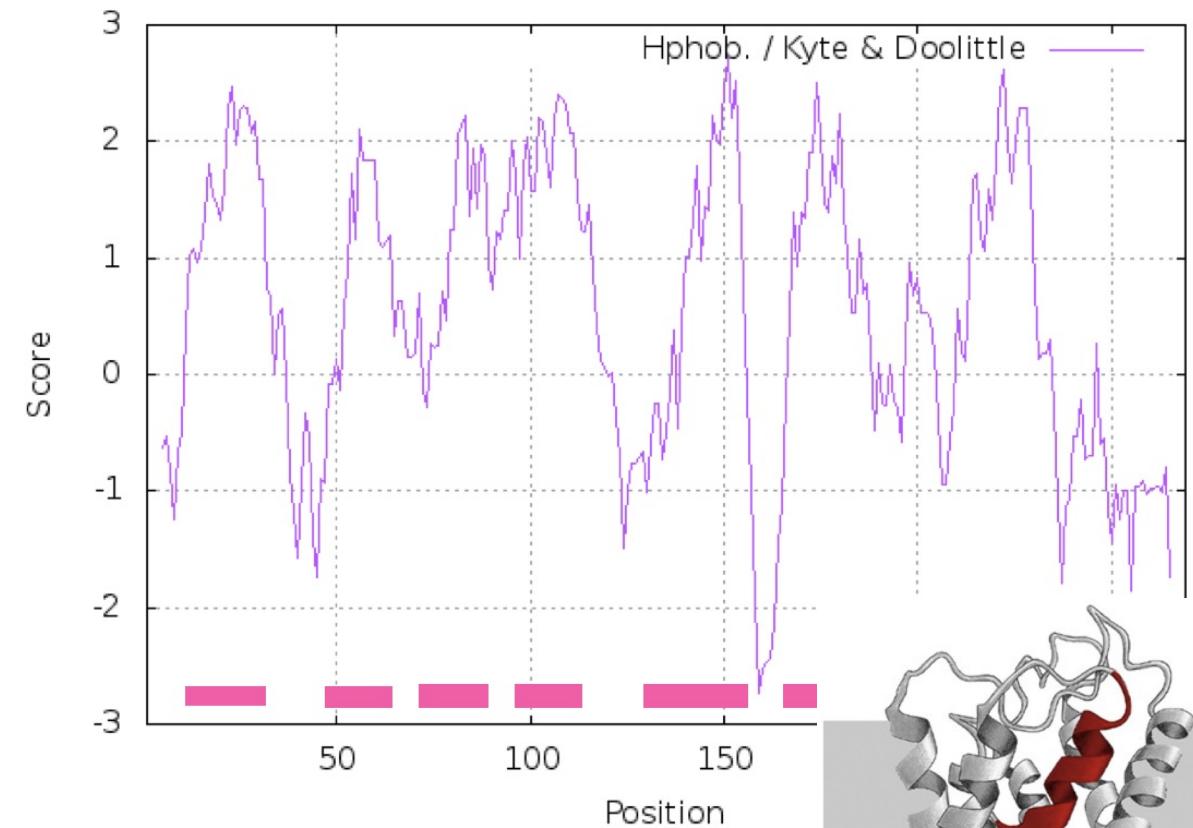


# Globular vs. membrane proteins

Prc Hemoglobin 1AN



ProtSci Aquaporin 1MAN

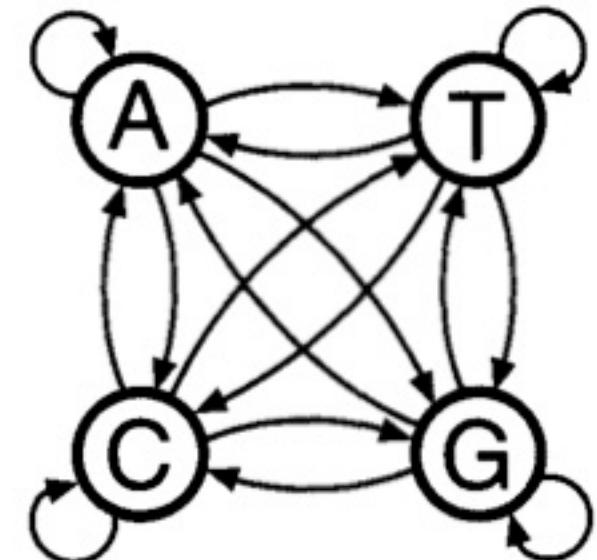


How do we know where hydrophobicity is “high enough”?

# Hidden Markov models

# CS detour: Markov chains

- **Markov property:** Current state of the model depends on its previous state(s)
- Example: in the genome of vertebrate, C's are avoided before G's (CpG context)
  - C's before G's get methylated, which leads to functional consequences
  - in some regions, CpG context is needed for functional reasons
- We can model this sequence with a probabilistic model that is aware of its previous state — **Markov chain**



# Markov chains cont'd: parameters

- In the CpG example, letters: **states**
- We create a sequence as we travel between states
- Parameters: **transition probabilities**

$$a_{st} = P(x_i = s | x_{i-1} = t)$$

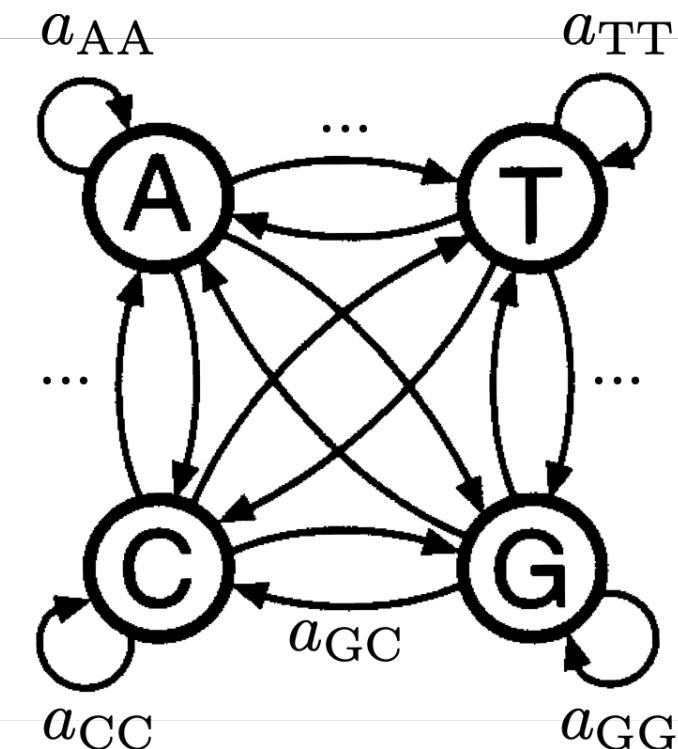
- Probability of a sequence  $X$ :

$$P(X) = P(x_L, x_{L-1}, \dots, x_1) =$$

$$= P(x_L | x_{L-1}, \dots, x_1) \cdot P(x_{L-1} | x_{L-2}, \dots, x_1) \cdot \dots \cdot P(x_2 | x_1) \cdot P(x_1) =$$

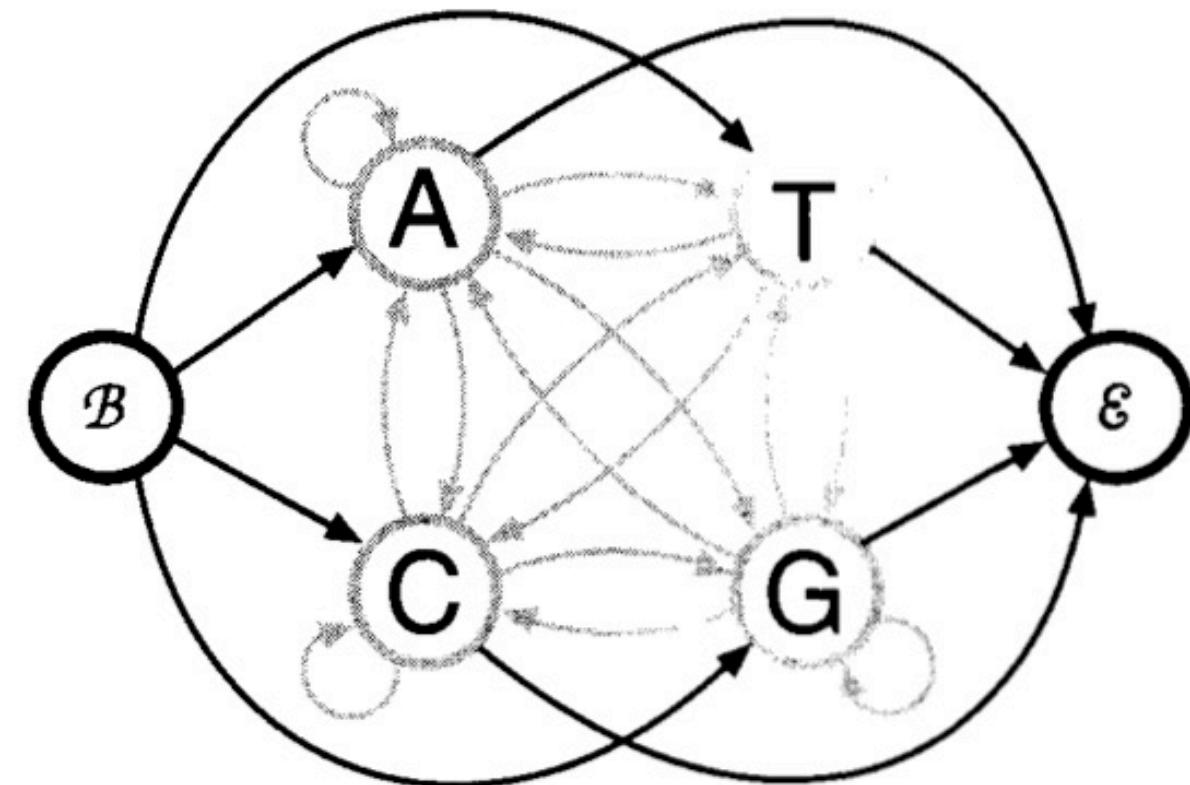
$$= P(x_L | x_{L-1}) \cdot P(x_{L-1} | x_{L-2}) \cdot \dots \cdot P(x_2 | x_1) \cdot P(x_1) =$$

$$= \prod_{i=2}^L a_{x_i x_{i-1}} \cdot P(x_1)$$



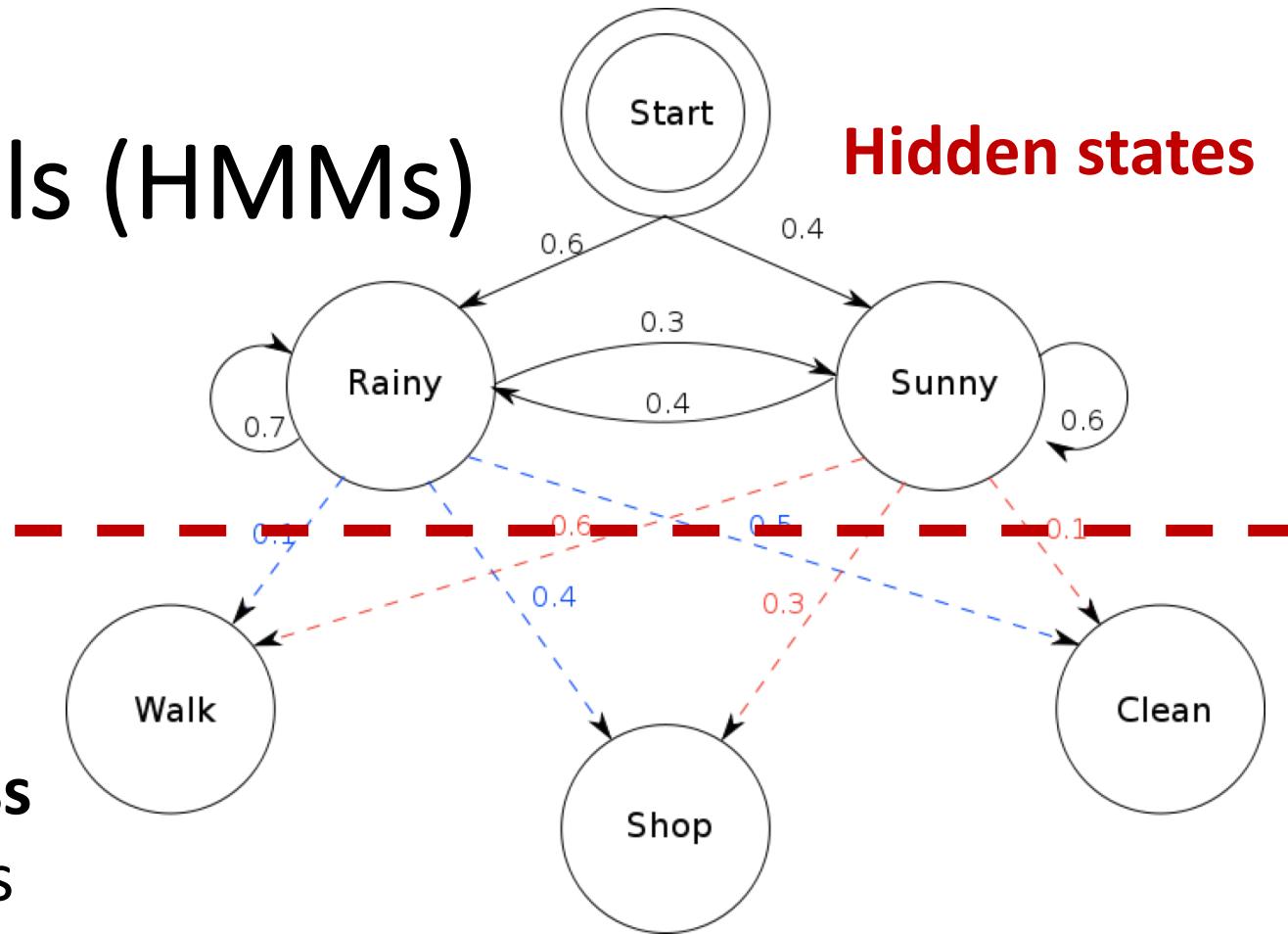
# Markov chains cont'd

- To get rid of  $P(x_1)$ ,  
add technical **start** and **end** states
- Then  $P(X) = \prod_{i=2}^L a_{x_i x_{i-1}}$



# Hidden Markov models (HMMs)

- The modelled system is assumed to be a **Markov process** over unobserved (**hidden**) states
- Hidden states emit observable states
- Task: given the sequence of observable states, **reconstruct the most probable sequence of the hidden states**



# HMMs cont'd: parameters

- **Hidden states ( $n$ ):**

$k$  transitions between the states =>

$k \leq n \times n$  **transition probabilities**

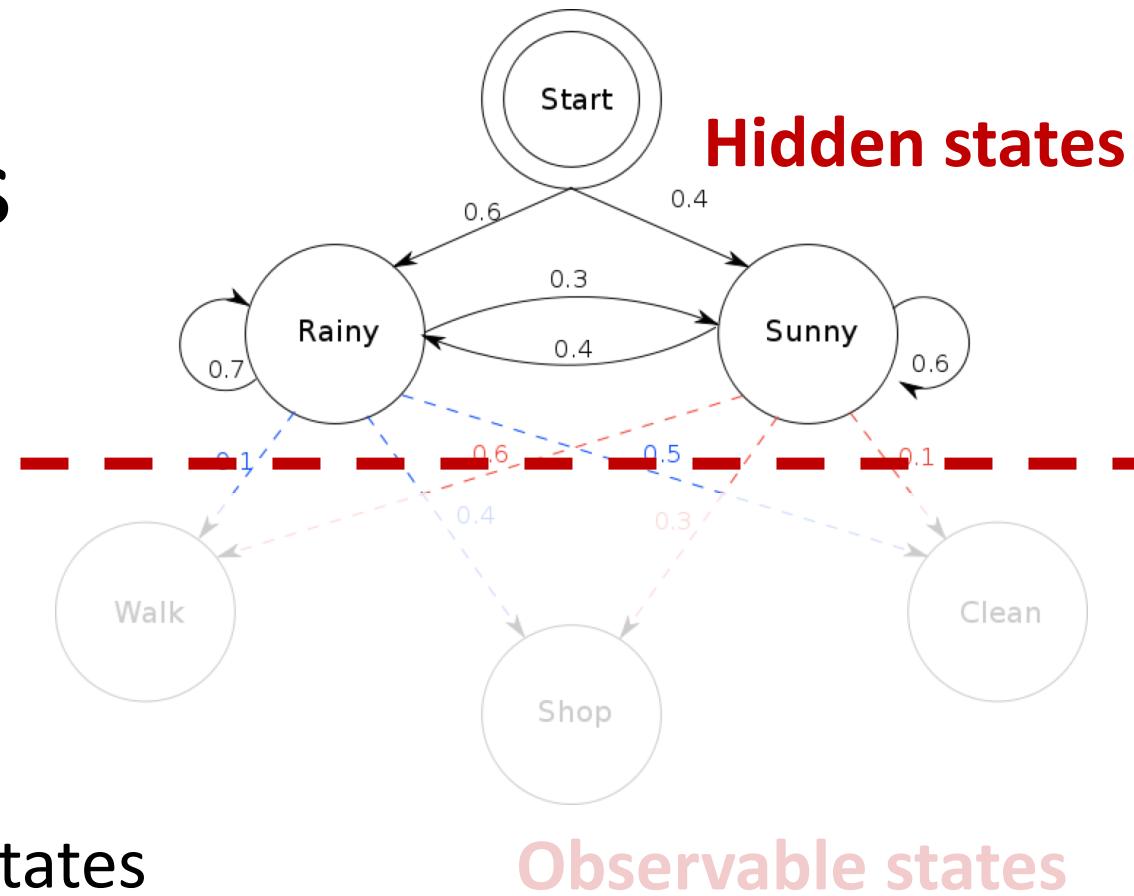
- assume non-existing transitions

have a probability of 0, then  $n^2$

transition probabilities between  $n$  states

- $n \times n$  matrix,  $n \times (n - 1)$  independent parameters

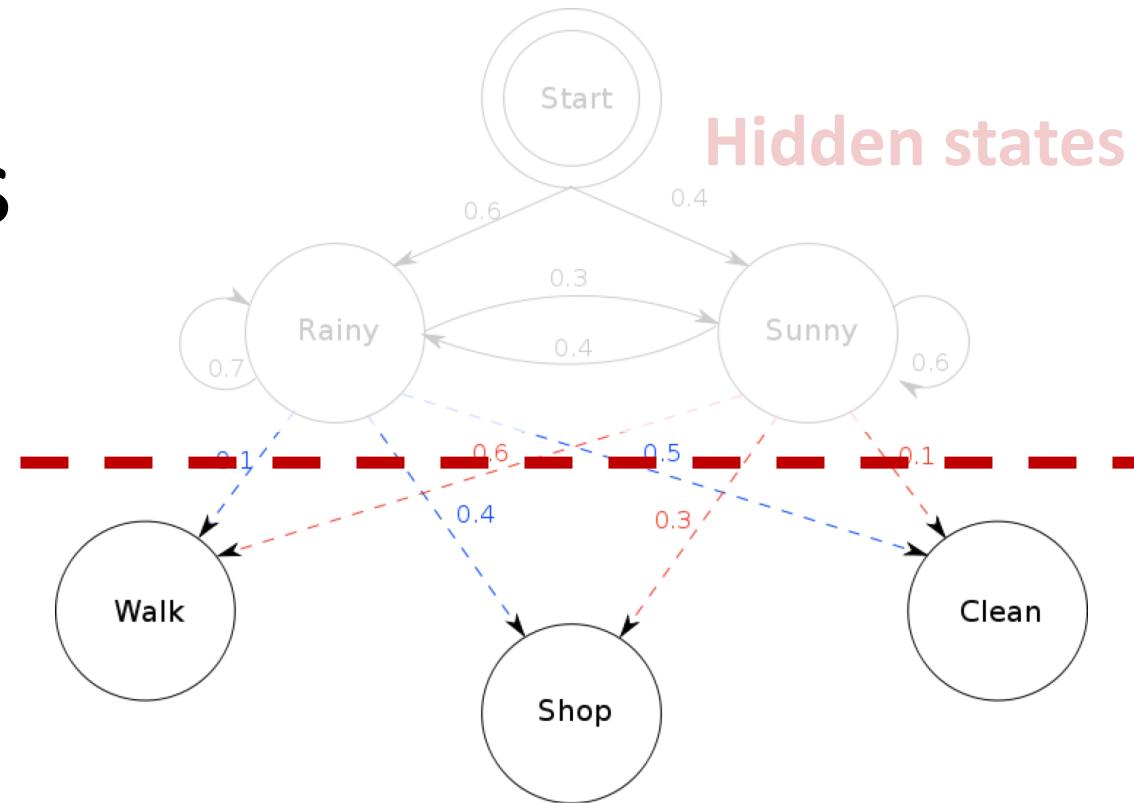
- **Markov matrix:** stochastic, the state at time  $t$  depends only on the state at time  $t - 1$



Observable states

# HMMs cont'd: parameters

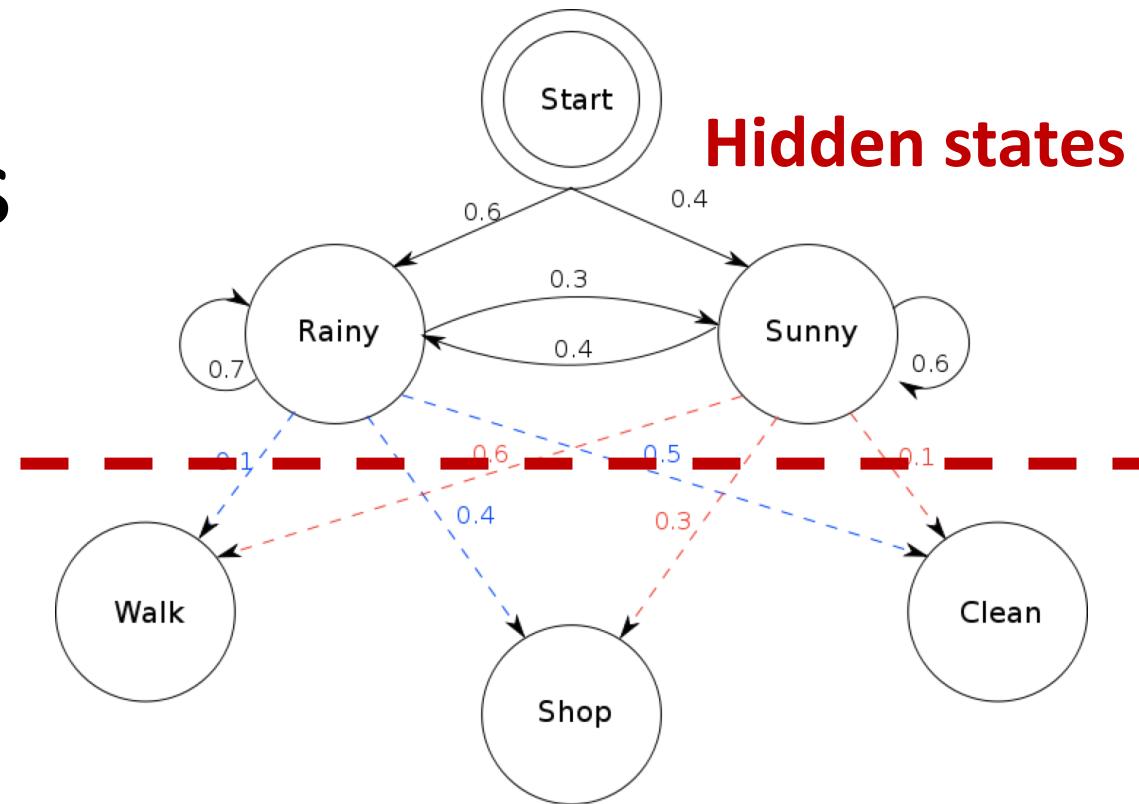
- **Observable states ( $m$ ):**  
 $n$  emitting states;  
 $m$  observable states =>  
 **$n \times m$  emission probabilities**
  - if a state is not emitting,  
assume all emission probabilities of it to be 0
- Emission probabilities:  $n \times m$  matrix,  $n \times (m - 1)$  independent parameters



# HMMs cont'd: parameters

**Hidden states ( $n$ ) and observable states ( $m$ ):**

- Transition probabilities:  
 $n \times n$  matrix,  $n \times (n - 1)$  parameters
- Emission probabilities:  
 $n \times m$  matrix,  $n \times (m - 1)$  parameters
- $\Rightarrow n \times (n - 1) + n \times (m - 1) = \mathbf{n \times (m + n - 2)}$  independent parameters



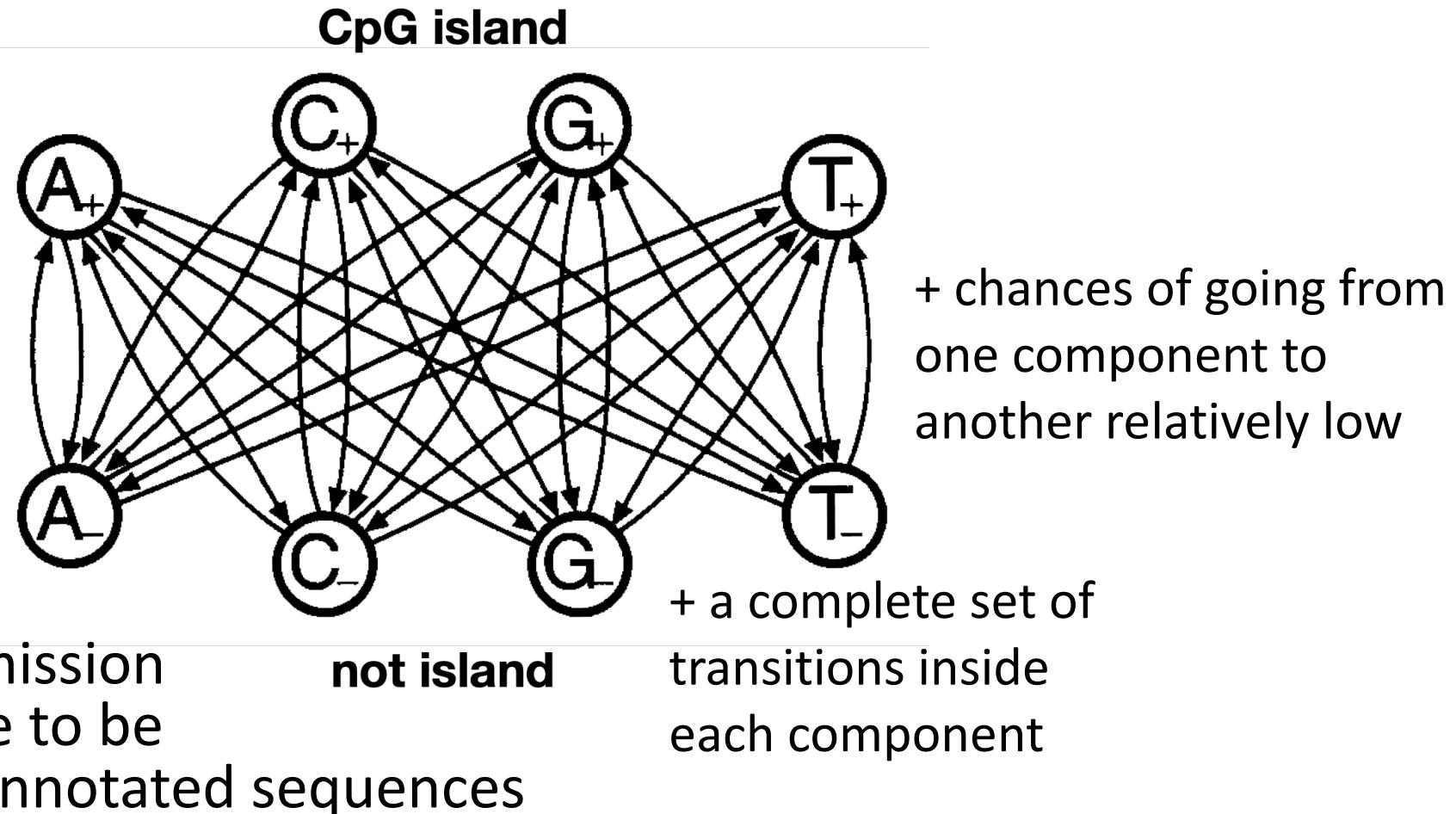
**Observable states**

# Modelling sequences with HMMs

- In a HMM, states **do not correspond** to observed characters
- States can **emit** observed characters (states)
- Example: in the genome of vertebrate, C's are avoided before G's (CpG context)
  - C's before G's get methylated, which leads to functional consequences (beyond the scope of this course)
  - in some regions, CpG context is needed for functional reasons — **CpG islands**
- Find CpG island in an annotated sequence using an HMM
- => **hidden states: island / non-island; observable states: {A, T, G, C}**

# Model topology: piece of creativity

- For CpG islands:



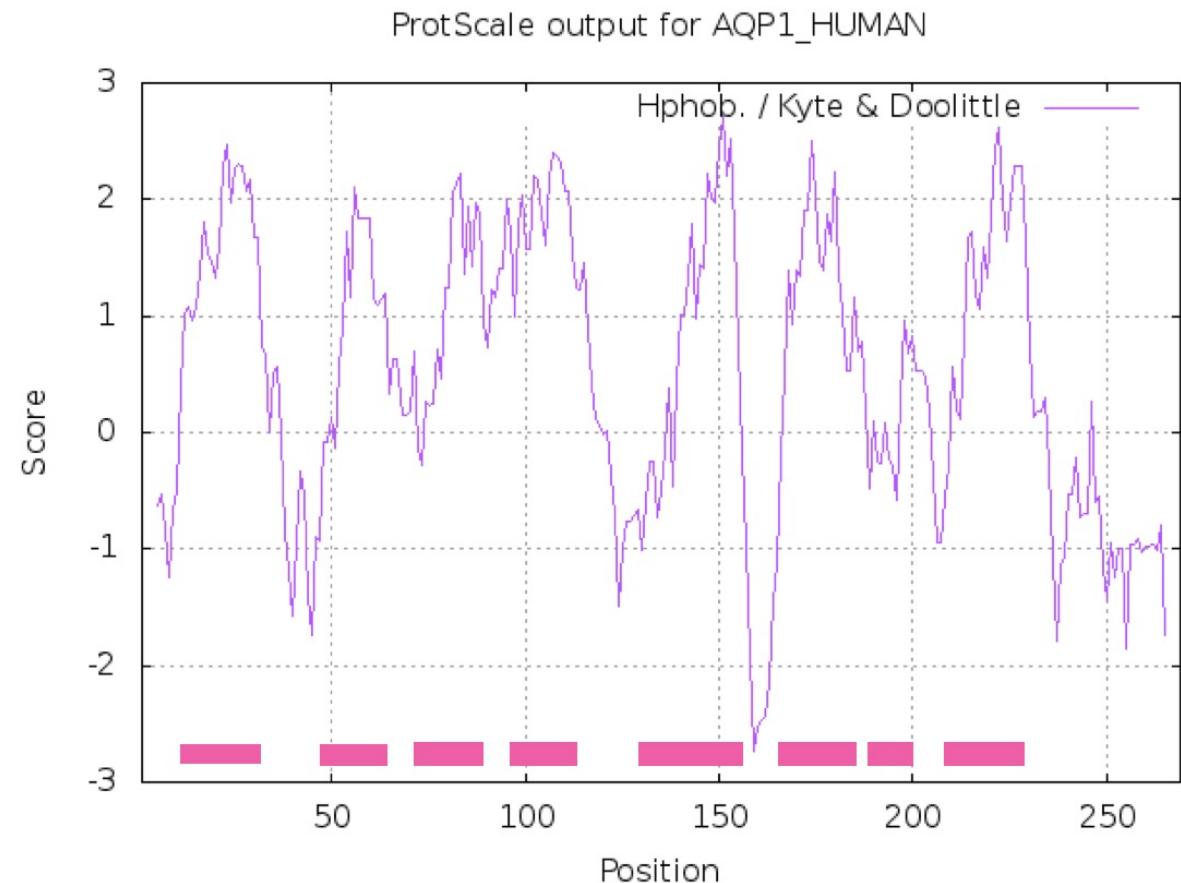
For details: See Durbin et al., "Biological Sequence Analysis", Cambridge University Press, 1998

# Model topology (a.k.a. HMM architecture)

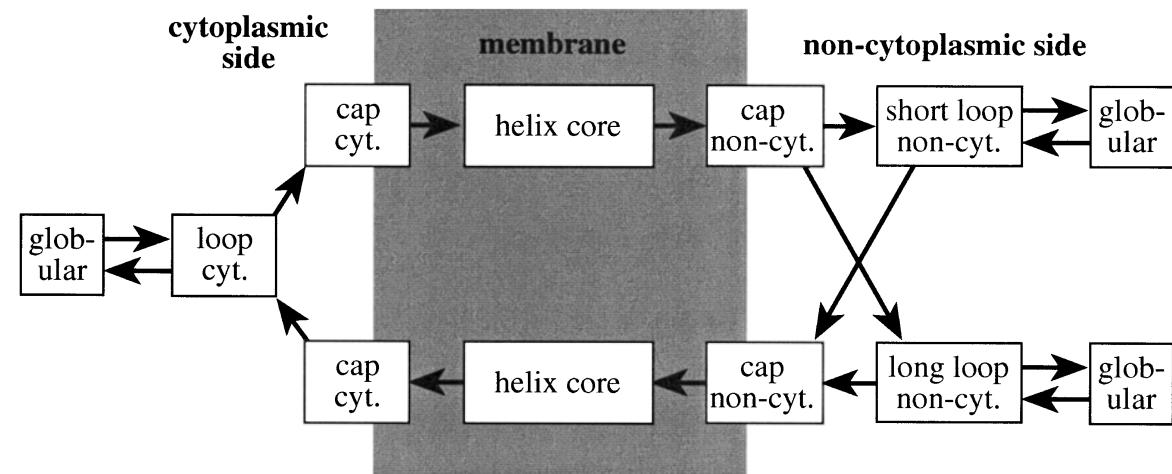
- How to model a sequence of a certain length?
- Suppose first 5 characters of a sequence are always **A**, second five are always **C**, then **A**'s again, etc.?
- Suppose we have two sets of distributions for emission probabilities for different parts of the sequence (membrane vs. non-membrane regions)?

# TMHMM (Krogh, Larsson, von Heijne, Sonnhammer, *J Mol Biol*, 2001)

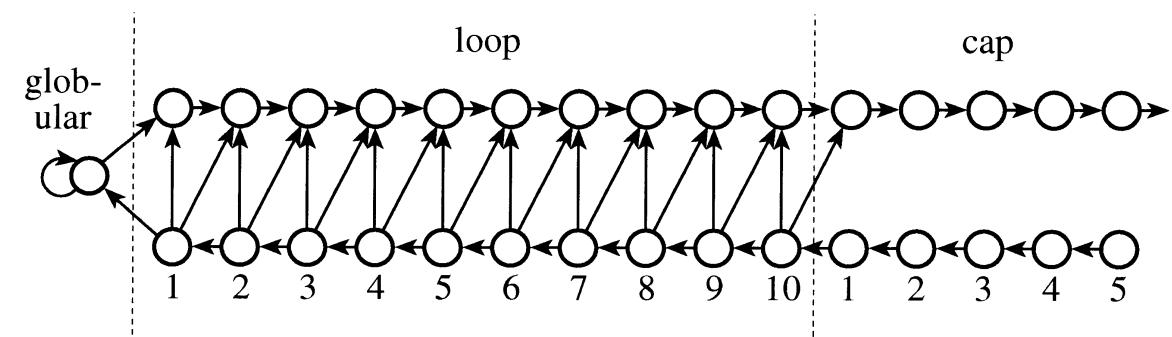
- Predict TM segments from protein sequences
- TM segments correspond to more hydrophobic regions
- Idea: predict TM segments using an HMM



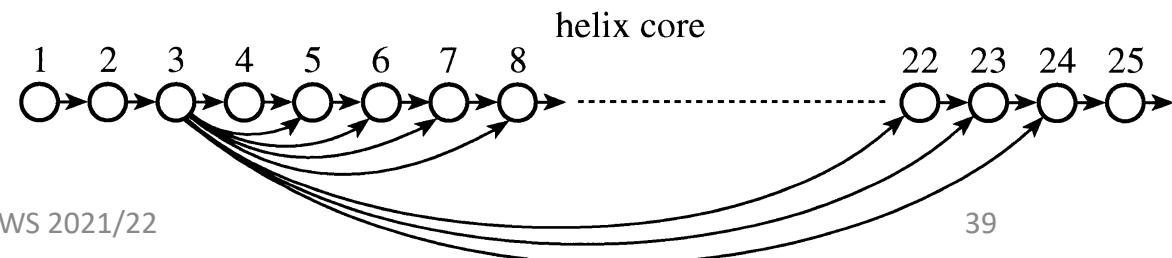
# HMM architecture



(b)



(c)

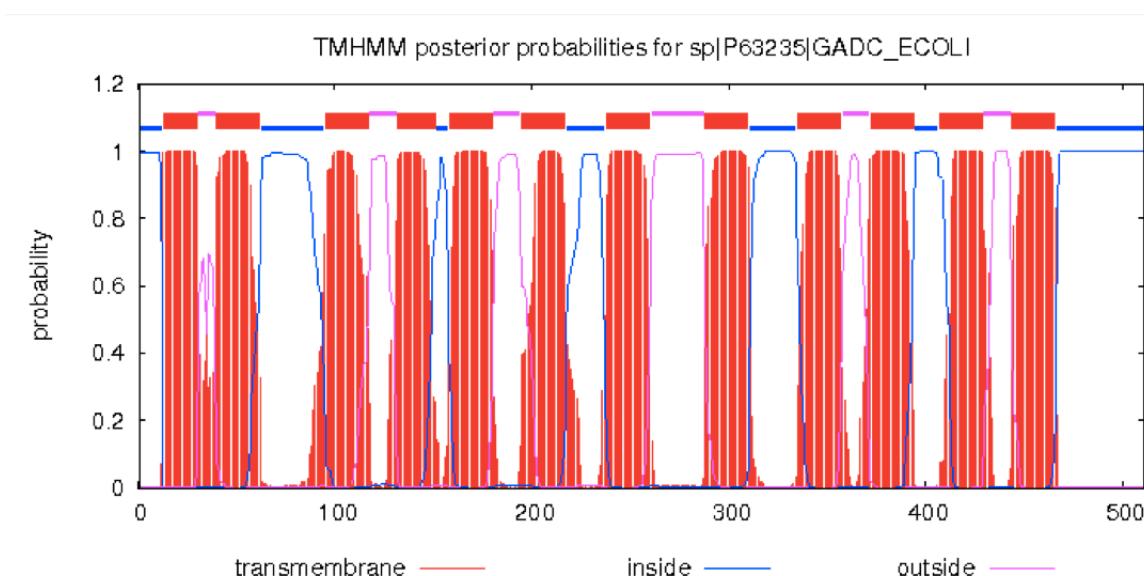


# TMHMM performance

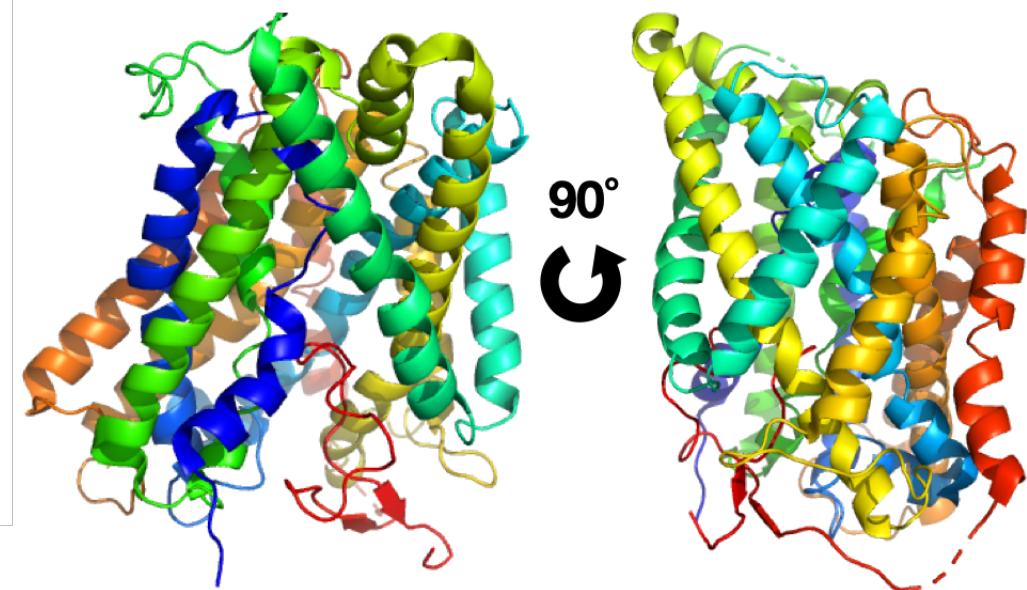
**10-fold  
cross-validation**

	Cross-validation	Mean and std. dev.
Number of proteins	160	
of which single-spanning:	52	32.50 %
Correctly predicted topology:	124	77.50 %
Invertedly predicted topology:	11	6.88 %
Correctly predicted N-terminal:	141	88.12 %
Under-predictions:	16	10.00 %
of which single-spanning:	1	0.62 %
Over-predictions:	12	7.50 %
of which single-spanning:	7	4.38 %
Both over- and under-predictions:	3	1.88 %
of which single-spanning:	1	0.62 %
Total number of real helices:	696	
Number of over-predicted helices:	17	2.44 %
Number of under-predicted helices:	19	2.73 %
Number of shifted helix predictions:	0	0.33
Number of falsely merged helices:	0	0.50
Number of falsely split helices:	0	0

# TMHMM performance: example



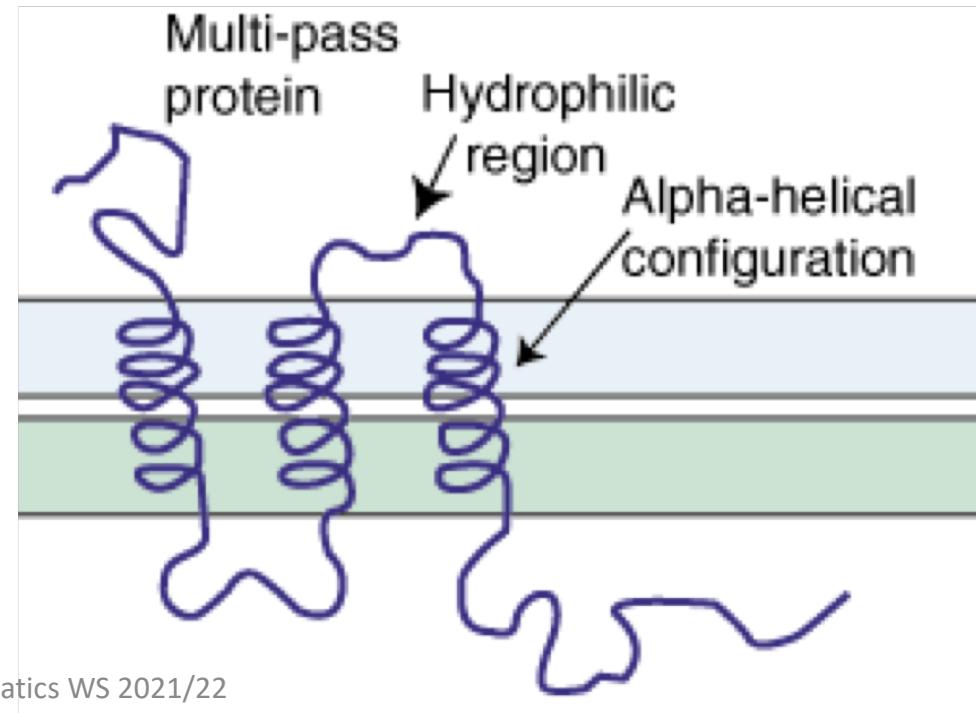
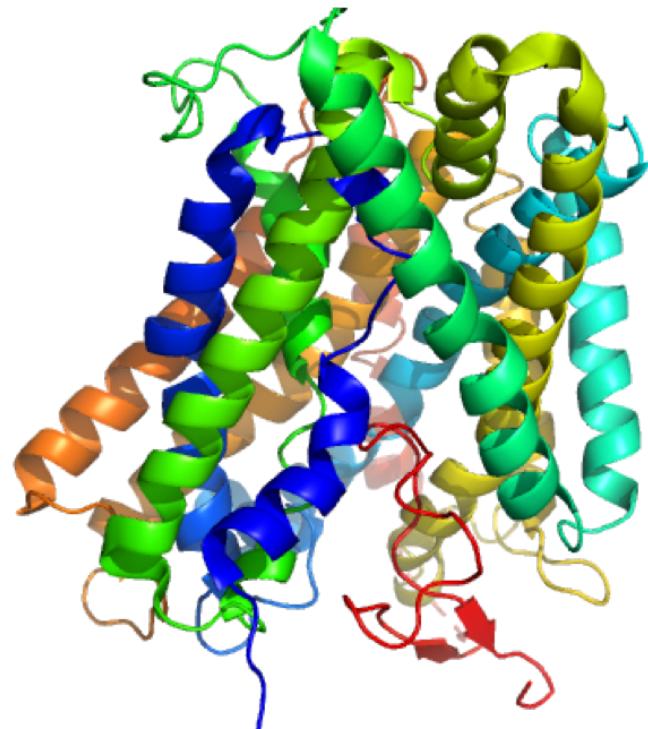
**GadC *E.coli* (glutamate transporter)**



# TMHMM: a caveat

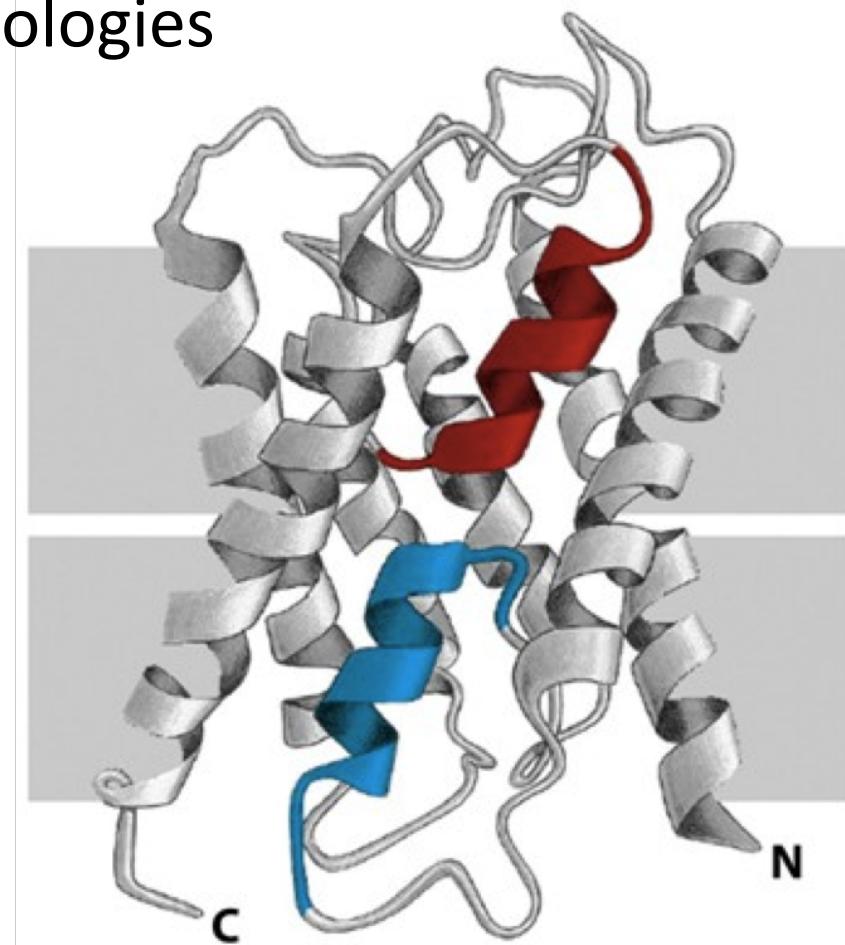
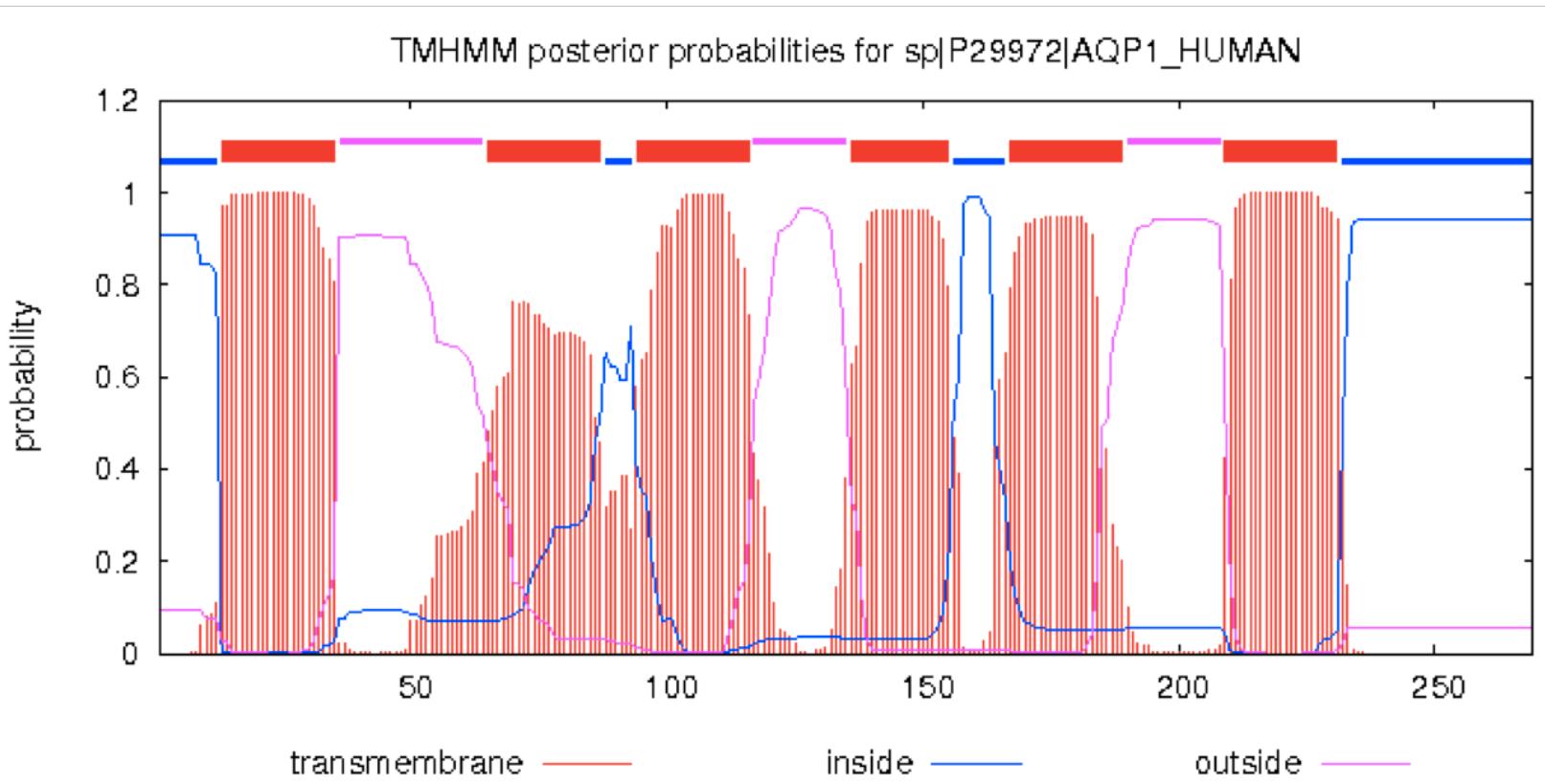
- Can't account for even slightly non-standard topologies

## Standard topology of multi-pass $\alpha$ -helical proteins



# TMHMM: a caveat

- Can't account for even slightly non-standard topologies



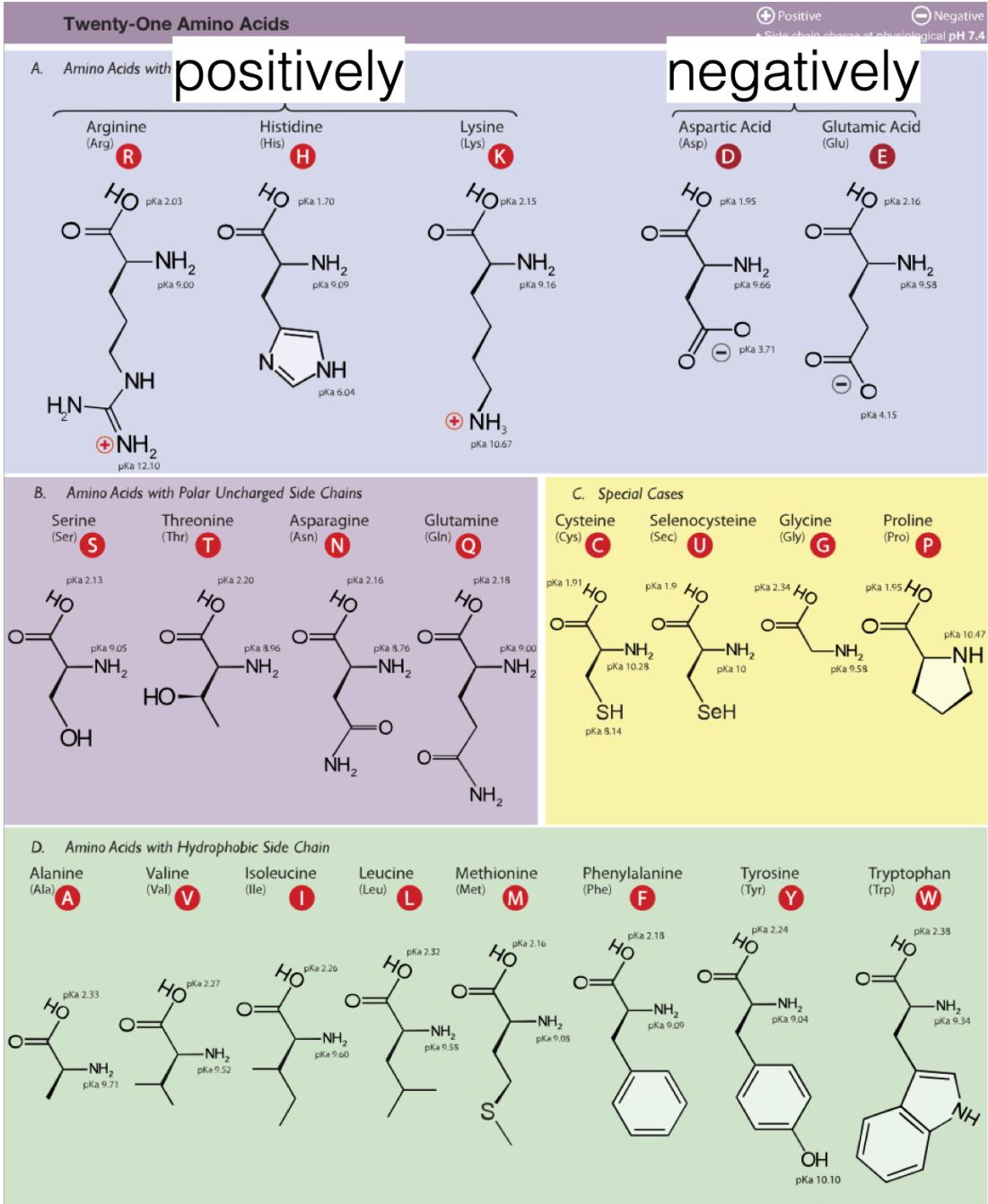
# Prediction of secondary structure

# Amino acids have different chemical properties

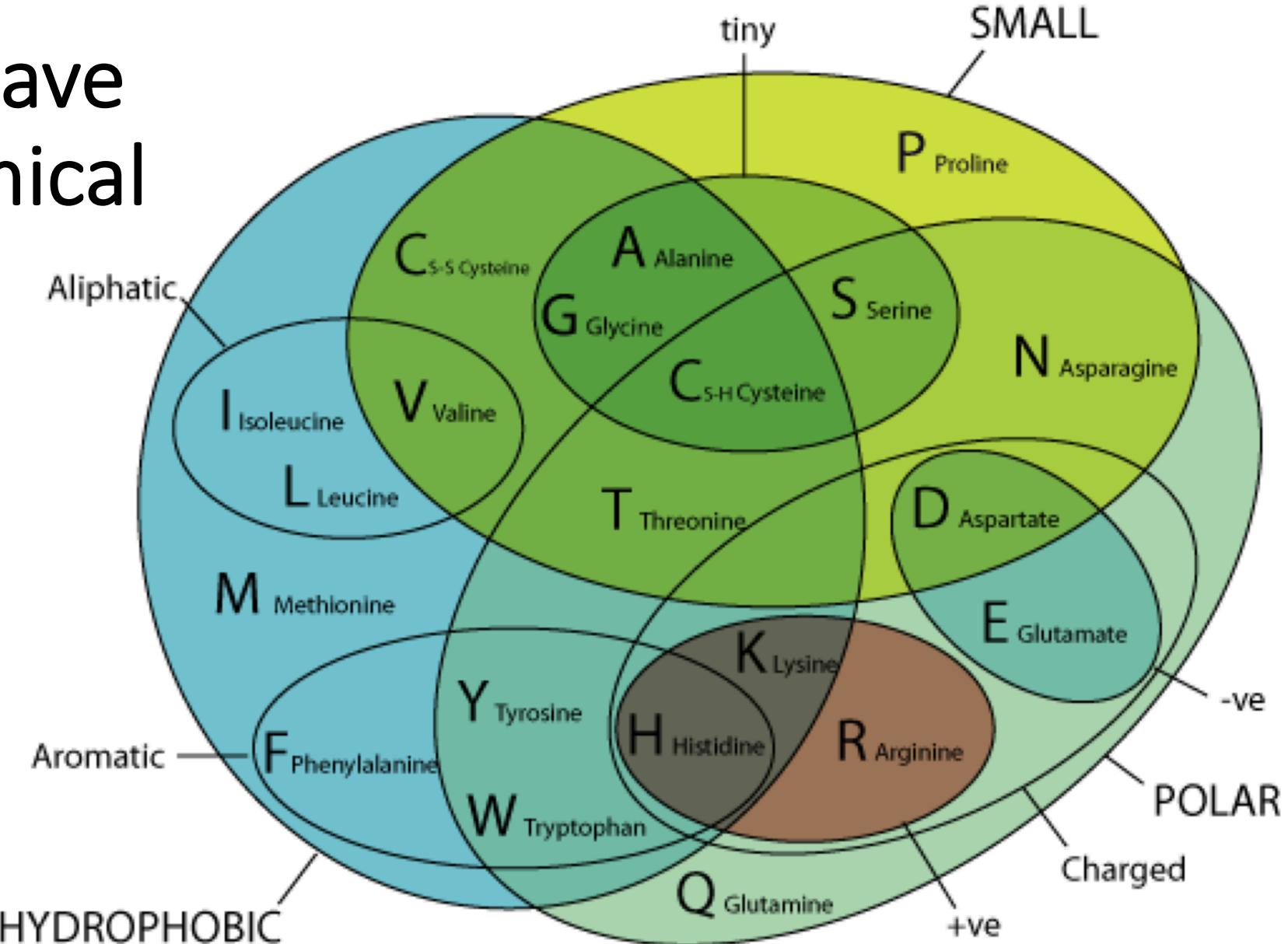
charged

polar

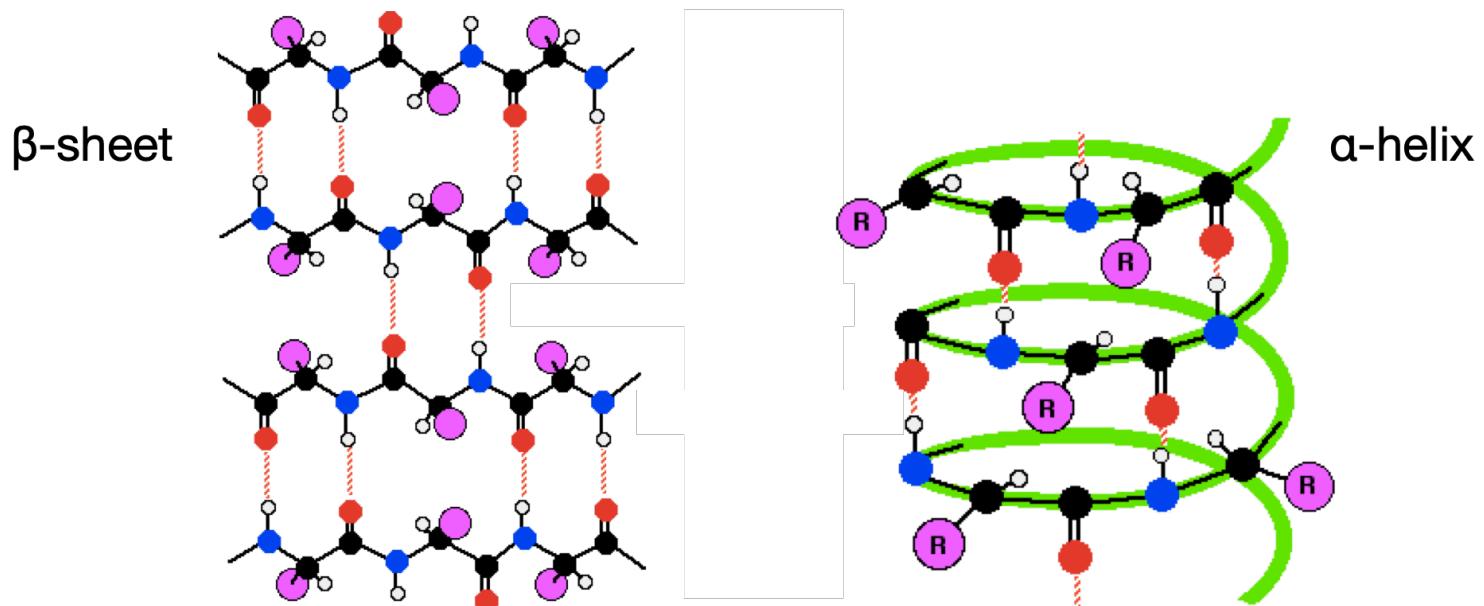
hydrophobic



# Amino acids have different chemical properties

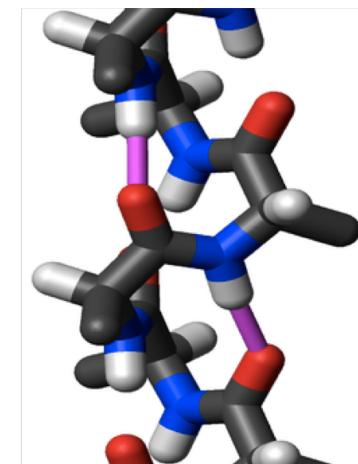
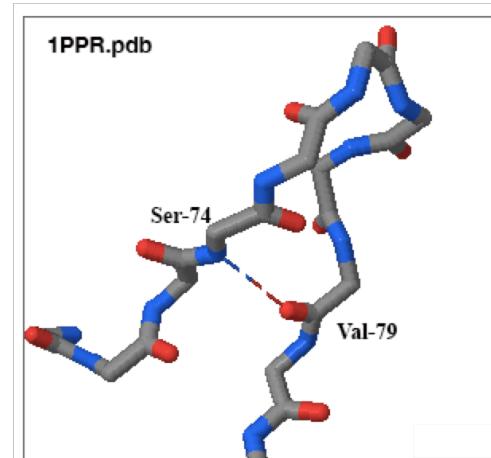


# Recap: most important secondary structure elements

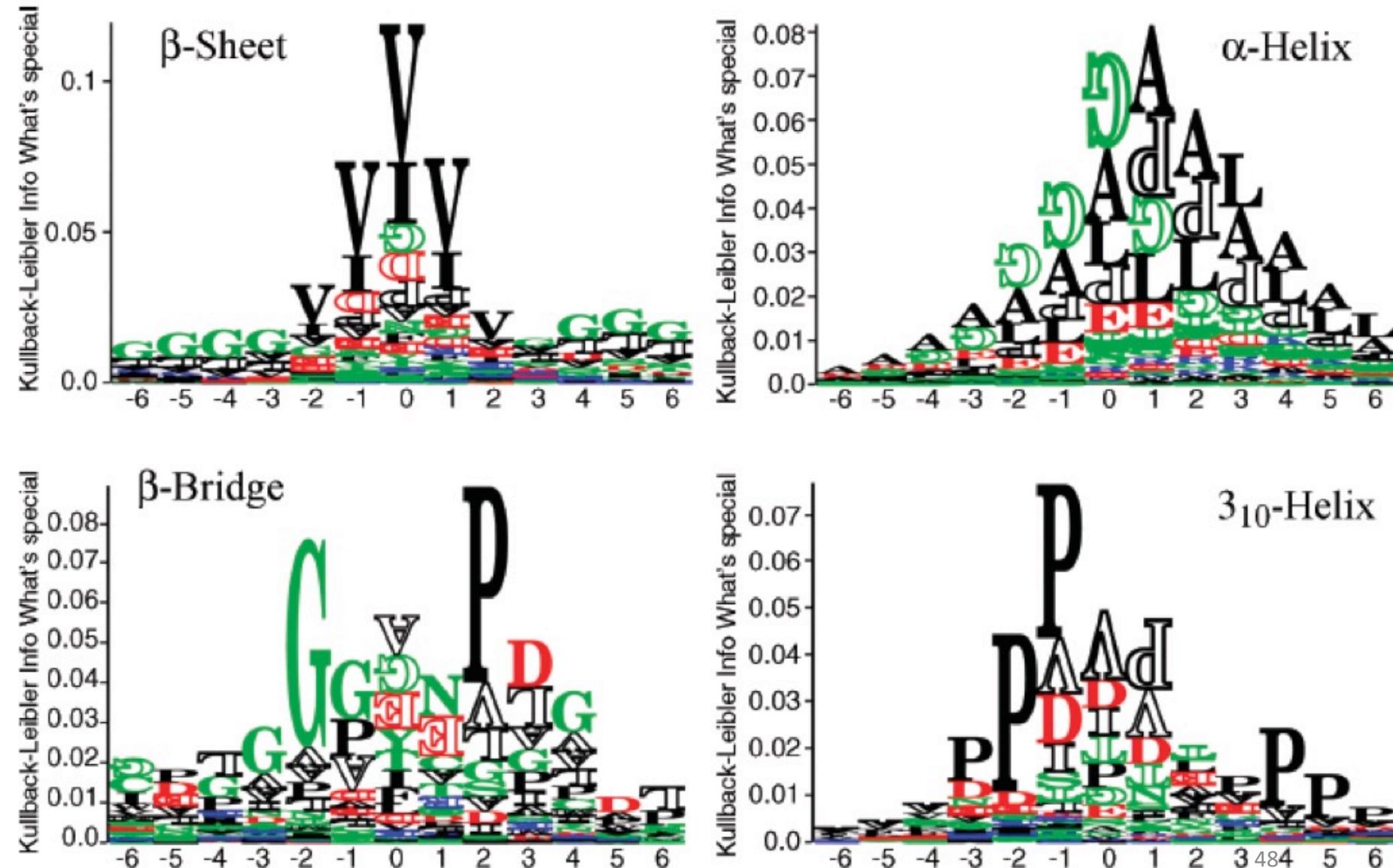


$\beta$ -bridge: two  $\beta$ -strands connected by a single H-bond (rare)

$3_{10}$ -helix:  $i + 3 \rightarrow i$  H-bonding;  
10-15% of all helices

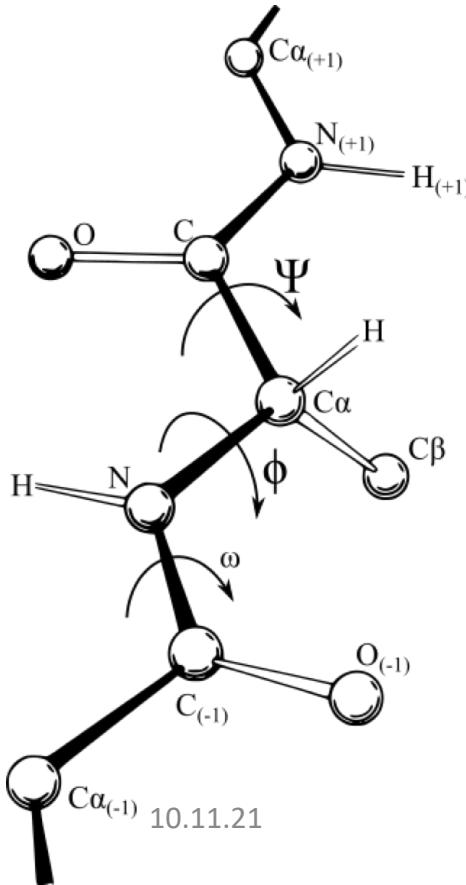


# Amino acid preferences in secondary structure elements

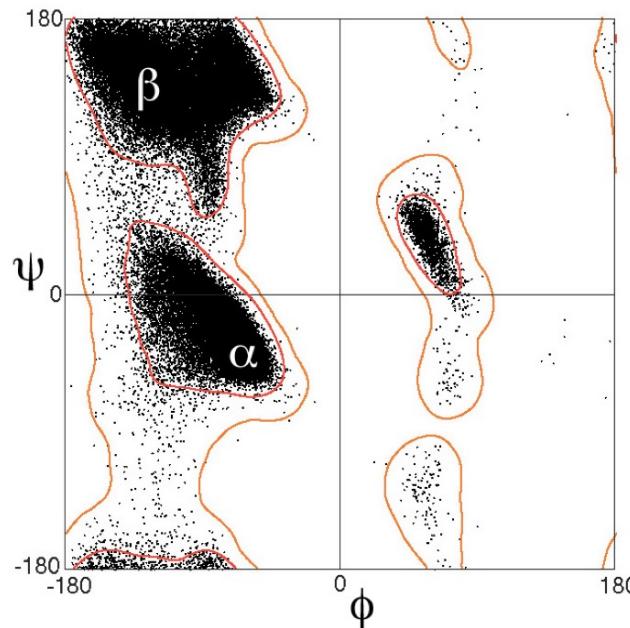


# Ramachandran plot: unusual amino acids

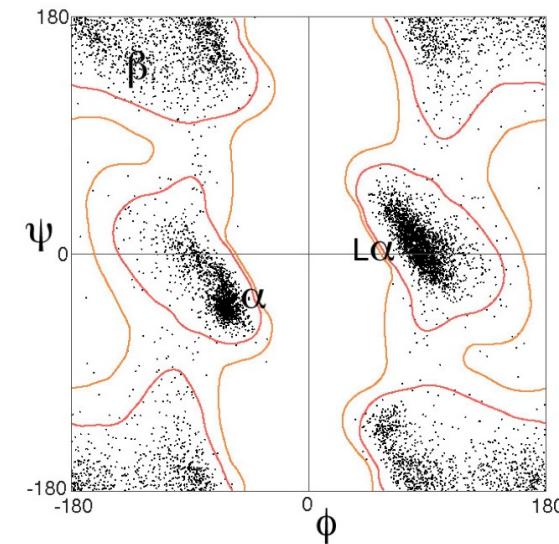
- Steric hindrance: the reason for disallowed regions in the Ramachandran plots



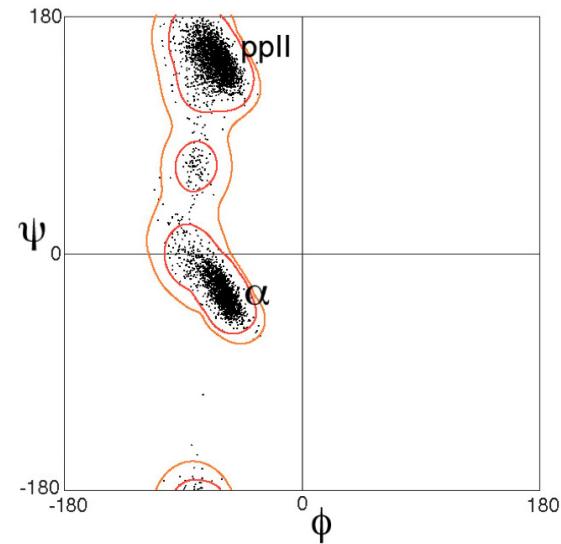
Almost all amino acids



Glycine



Proline



# Secondary structure assignment: DSSP

- **Assignment:** from known coordinates
- DSSP (Kabsch & Sander, 1983): H-bond energy from geometry => H-bonding pattern => secondary structure

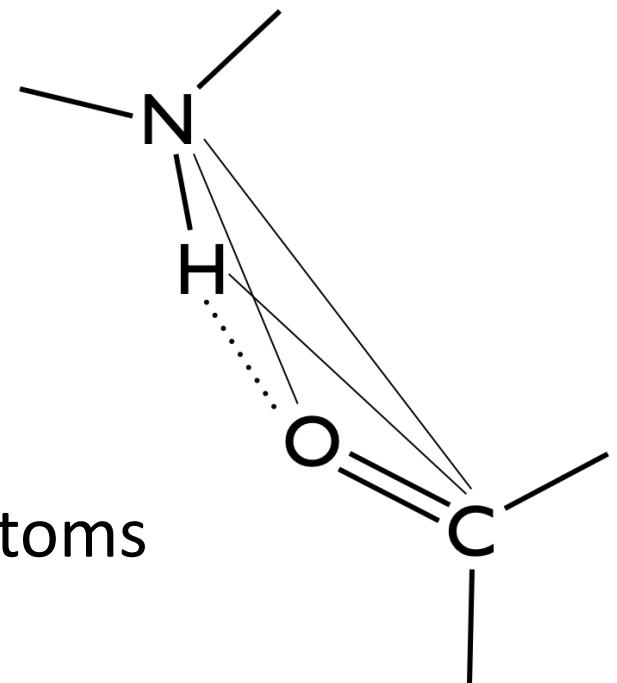
$$E = 0.084 \left( \frac{1}{r_{ON}} + \frac{1}{r_{CH}} - \frac{1}{r_{OH}} - \frac{1}{r_{CN}} \right) \cdot 332 \text{ kcal/mol}$$

0.084: product of partial charges

332: dimensionality factor

- **H-bond exists, if  $E < -0.5 \text{ kcal/mol}$**

- Position of H can be inferred from positions of other atoms (parallel to C=O,  $r_{NH} = 1$ )



# DSSP cont'd

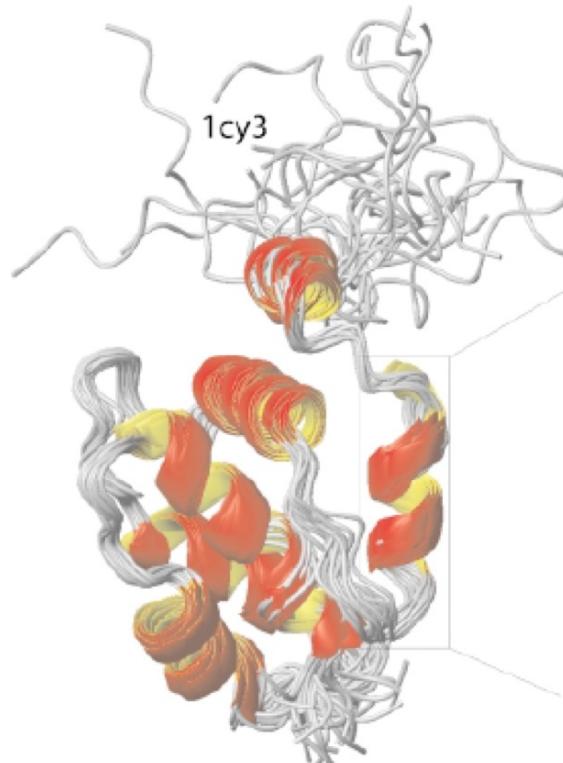
- H-bond energy from geometry => H-bonding pattern => secondary structure
  - Define elementary H-bonding patterns:
    - **n-turn**: H-bond between residues  $(i, i + n)$ ,  $n = 3, 4, 5 \Rightarrow$  **helix**
    - **Bridge** between  $(i - 1, i, i + 1)$  and  $(j - 1, j, j + 1)$ :
      - parallel  $(i, j)$ :  $(\text{H-bond}(i - 1, j) \And \text{H-bond}(i, j + 1)) \parallel (\text{H-bond}(j - 1, i) \And \text{H-bond}(j, i + 1))$
      - anti-parallel  $(i, j)$ :  $(\text{H-bond}(i, j) \And \text{H-bond}(j, i)) \parallel (\text{H-bond}(i - 1, j + 1) \And \text{H-bond}(j - 1, i + 1))$
- => **beta-ladder** => **sheet**

# Reminder: DSSP code

- H = alpha helix
- B = residue in isolated beta-bridge
- E = extended strand, participates in beta ladder
- G = 3-helix (3/10 helix)
- I = 5 helix (pi helix)
- T = hydrogen bonded turn
- S = bend
- C/- = no secondary structure (coil)

# DSSPcont: Continuous secondary structure assignment (Carter et al., 2003)

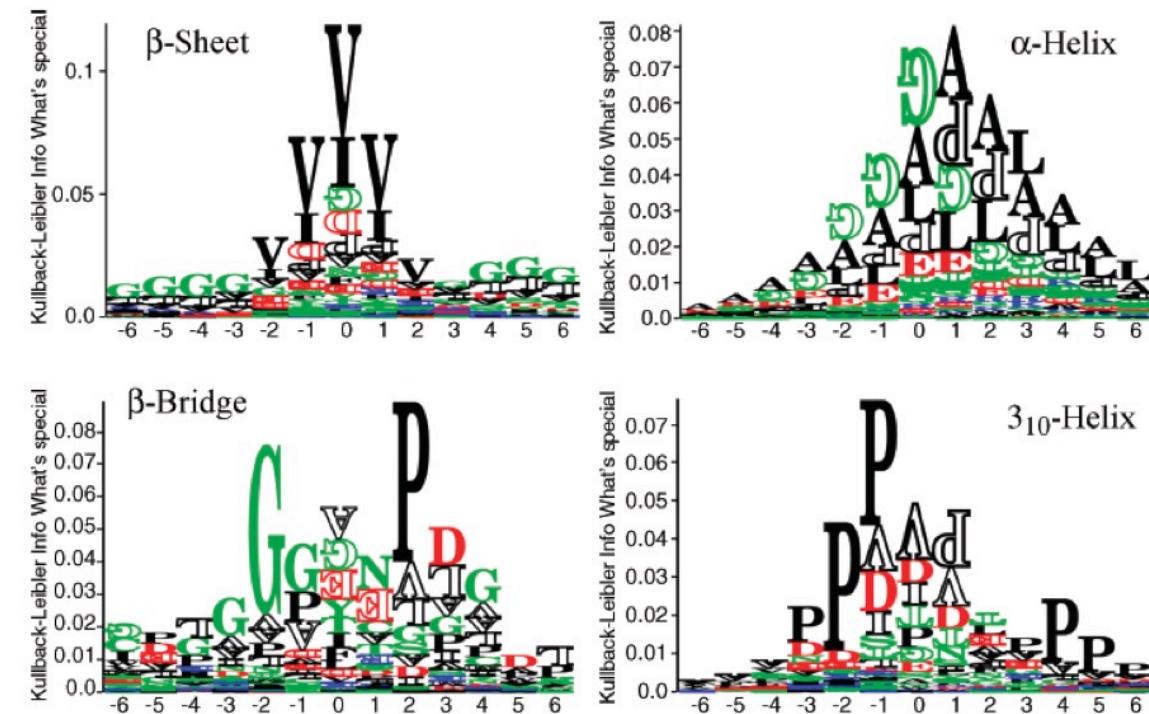
- 9 different H-bond energy thresholds (-1.0 kcal/mol to -0.2 kcal/mol) weighted using NMR structures
- difference between assignment using single models vs. for the average model



	(a) DSSP assignment for models 1-23					(b) DSSPcont assignment for model 1							
	1	5	10	15	20 23	G	H	I	T	E	B	S	L
20 Val	L	L	L	L	L	0	0	0	32	0	0	0	68
21 Ser	T	L	H	T	L	0	0	0	100	0	0	0	0
22 Glu	T	T	H	T	T	0	0	0	100	0	0	0	0
23 Glu	S	T	H	T	H	0	32	0	21	0	0	47	0
24 Ser	H	H	H	H	H	0	99	0	1	0	0	0	0
25 Leu	H	H	H	H	H	0	99	0	1	0	0	0	0
26 Asn	H	H	H	H	H	0	99	0	1	0	0	0	0
27 Lys	H	H	H	H	H	1	99	0	0	0	0	0	0
28 Val	H	T	H	H	T	1	87	0	12	0	0	0	0
29 Arg	T	T	T	T	T	1	53	0	46	0	0	0	0
30 Asn	S	S	S	S	S	0	0	0	10	0	0	90	0
31 Arg	S	S	T	S	S	0	0	0	32	0	0	68	0
32 Glu	S	S	T	L	S	0	0	0	32	0	0	68	0
33 Glu	L	L	L	L	L	0	0	0	0	0	0	0	100

# Prediction of secondary structure

- NO 3D coordinates
- Based on amino acid propensities to be in certain secondary structure elements
- Single amino acids or k-mers
- Information from one sequence or from multiple sequences can be used



# CS detour: Intro to classification problems

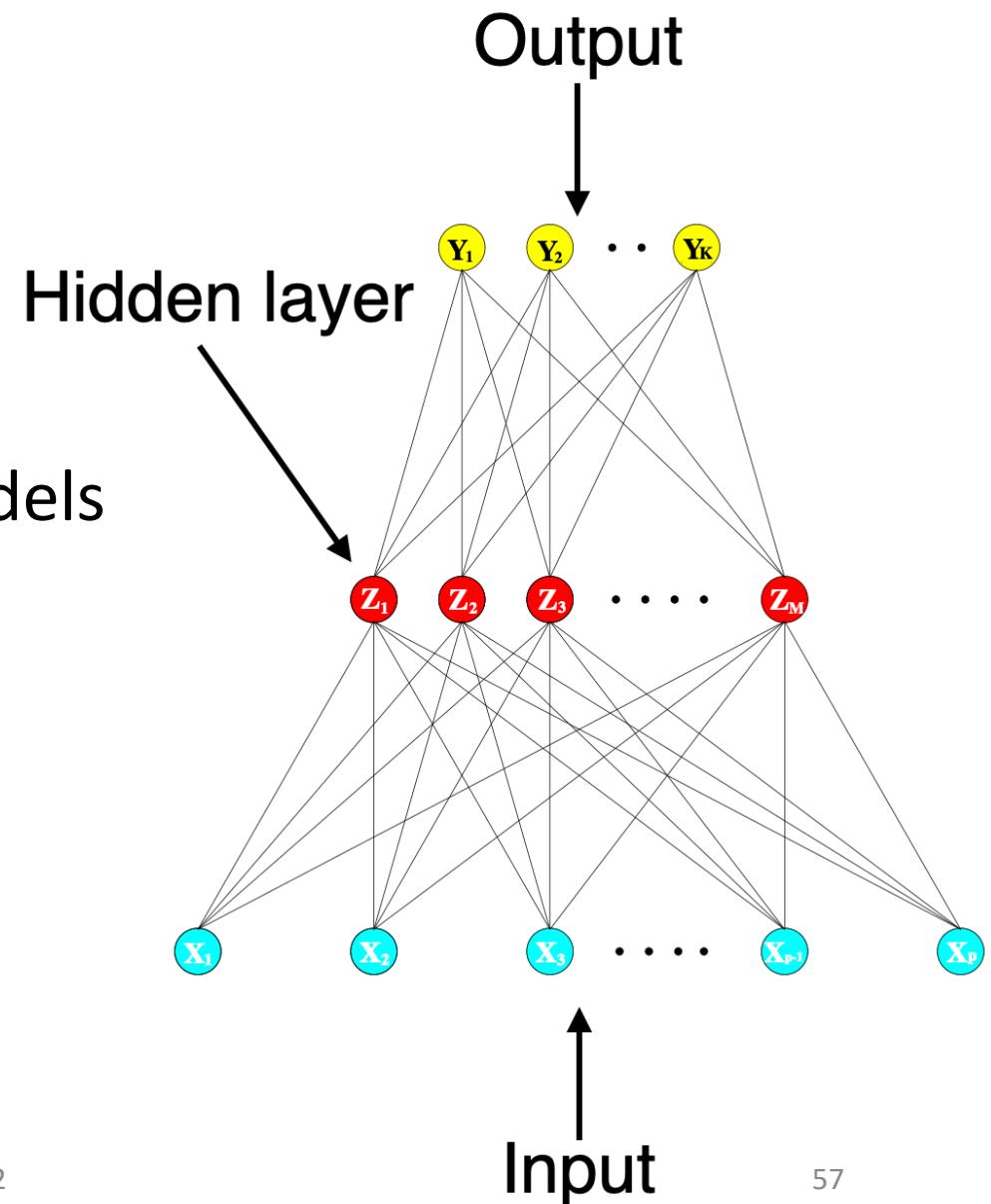
- Given a set of categories (e.g.  $\alpha$ -helices,  $\beta$ -strands, turns, ..., or DSSP categories), identify to which a new observation (e.g. residue) belongs
  - judging based on **features** of the new observation: measurable traits
- **Training set:** a set of observations with known features and categories (sequences with known 3D structures and assigned SSEs)
- **Test set:** new observations with only features known (protein sequences)

# Regression problem (only for completeness)

- Estimate relationship between variables
  - given a set of **independent variables** (~features, inputs), estimate the value of the **dependent variable** (~output)
- **Training set:** a set of observations with known values of independent and dependent variables
- **Test set:** new observations with only values of independent variables known

# Artificial neural networks

- Feed-forward
- Two-stage **regression or classification** models
- Input: numerical features
  - Feature vectors
- Output: class probabilities
  - For regression: only one output unit



# Artificial neural networks

- **Input features:** numerical vectors  $X$
- Derived features (**hidden layer**): linear combination of the input features:

$$Z_m = \sigma(\alpha_{0m} + \alpha_m^T X), m = 1, \dots, M$$

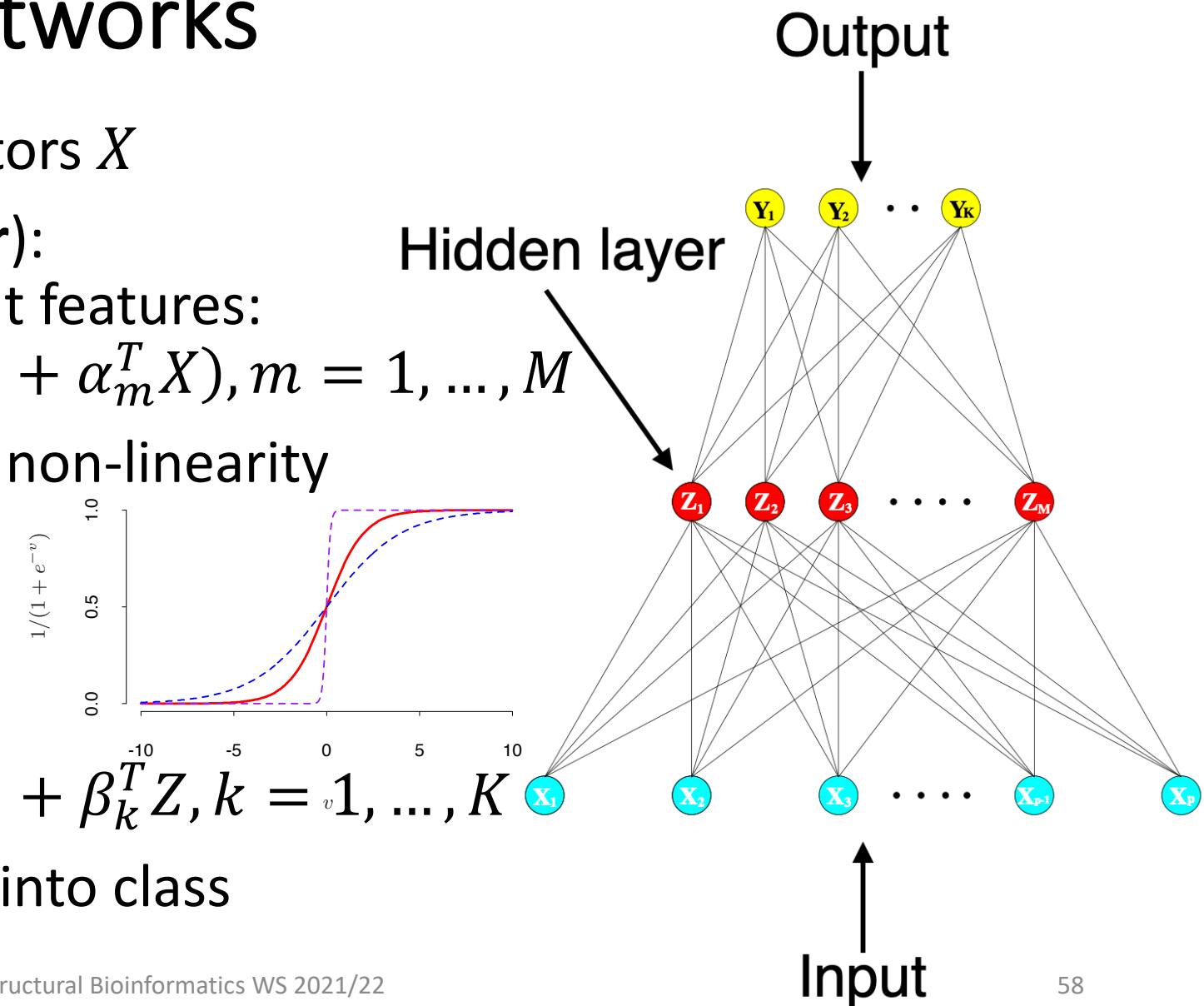
- $\sigma$ : **activation function**, inserts non-linearity

- Sigmoid:  $\sigma(v) = \frac{1}{1+e^{-v}}$

- **Target features ( $K$  classes):**

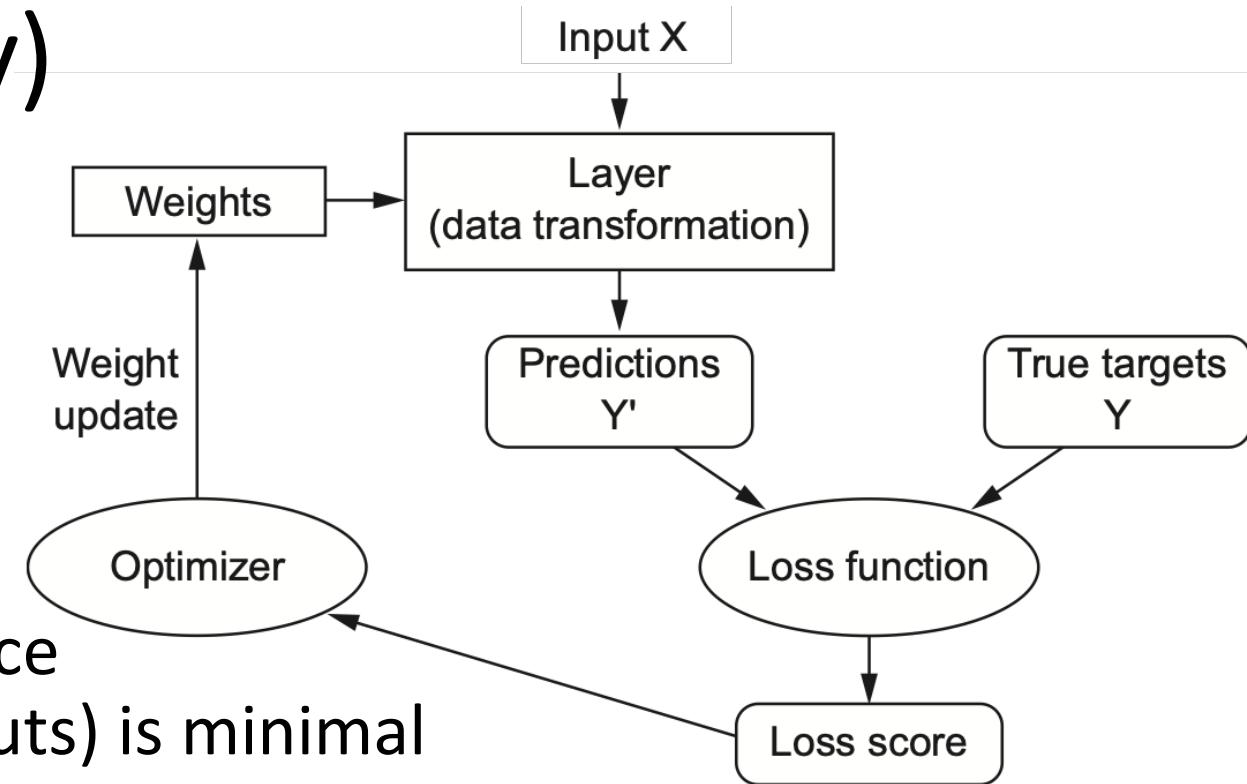
$$T_k = \beta_{0k} + \beta_k^T Z, k = 1, \dots, K$$

- Final normalization to turn  $T_k$  into class probabilities



# Training a network (or any other machine-learning model, actually)

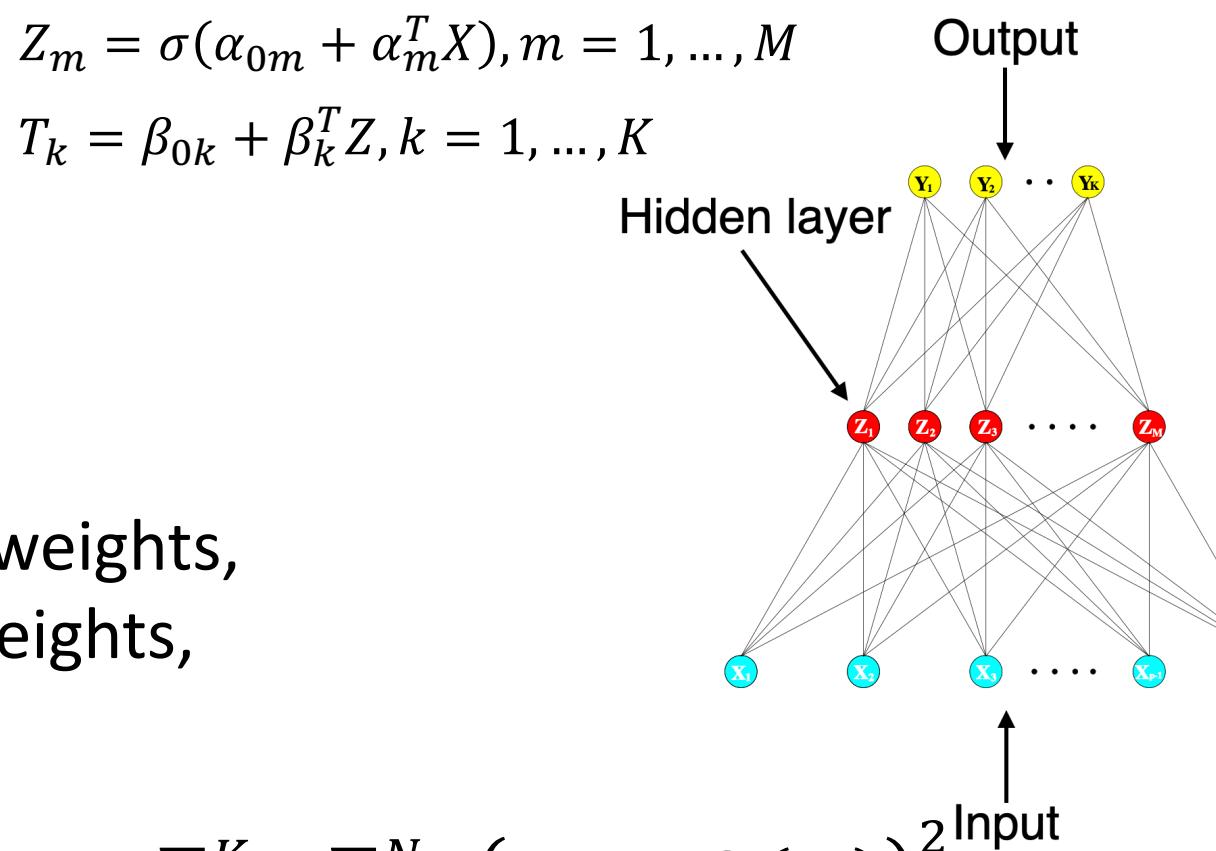
- $N$  samples:
  - Input features:  $x_i$
  - Known outputs:  $y_i$
- **Training:** adjust parameters, such that an loss function (difference between true and computed outputs) is minimal
- **$K$ -fold cross-validation:** split the data in  $K$  roughly equal parts, use  $K - 1$  for training and 1 for testing
- **Leave-one-out cross-validation:** train on all samples but one



# Fitting a network

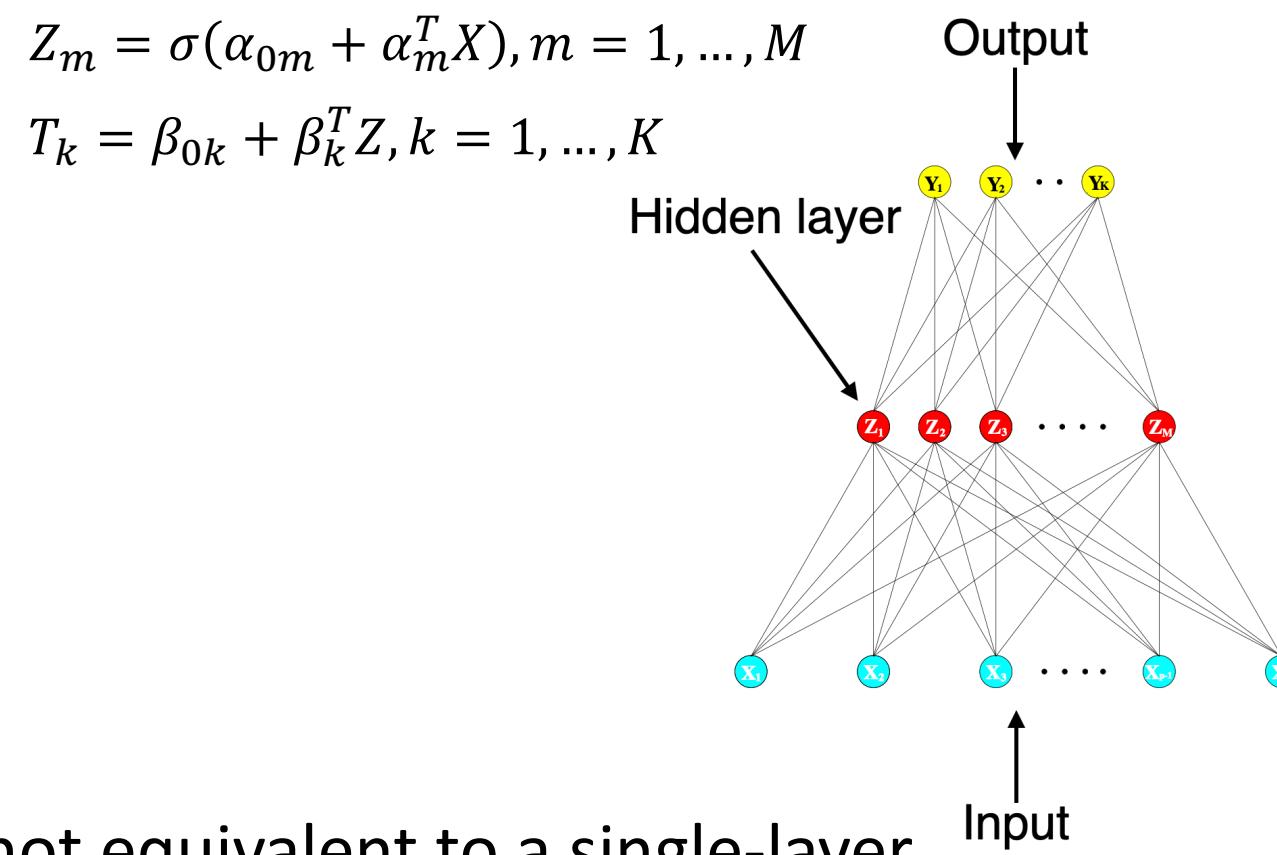
- Unknown parameters (**weights**):  
 $\{\alpha_0, \alpha_m\}, m = 1, \dots, M: M(1 + p)$  weights,  
 $\{\beta_0, \beta_k\}, k = 1, \dots, K: K(1 + M)$  weights,
- **Error function** (to be minimized):
  - For regression **sum-of-squares**:  $R = \sum_{k=1}^K \sum_{i=1}^N (y_{ik} - f_k(x_i))^2$
  - For classification **cross-entropy**:  $R = - \sum_{k=1}^K \sum_{i=1}^N y_{ik} \log f_k(x_i)$
- Start with random values close to 0 of weights, then minimize error by gradient descent

For details: See Hastie et al., “The Elements of Statistical Learning”, Springer, 2009



# Detour: deep learning

- >1 hidden layers
- Non-linear activation function => not equivalent to a single-layer network



# Back to bioinformatics: Protein family profile

- Suppose we have collected a bunch of related proteins (e.g. with BLAST) and built an alignment (multiple sequence alignment methods will not be covered in this course)
  - Can we use it to find other related proteins?
  - Is this better (i.e. more sensitive) than BLASTing individually with the protein sequences?
- YES — **protein family profile**

# Conservation in MSA

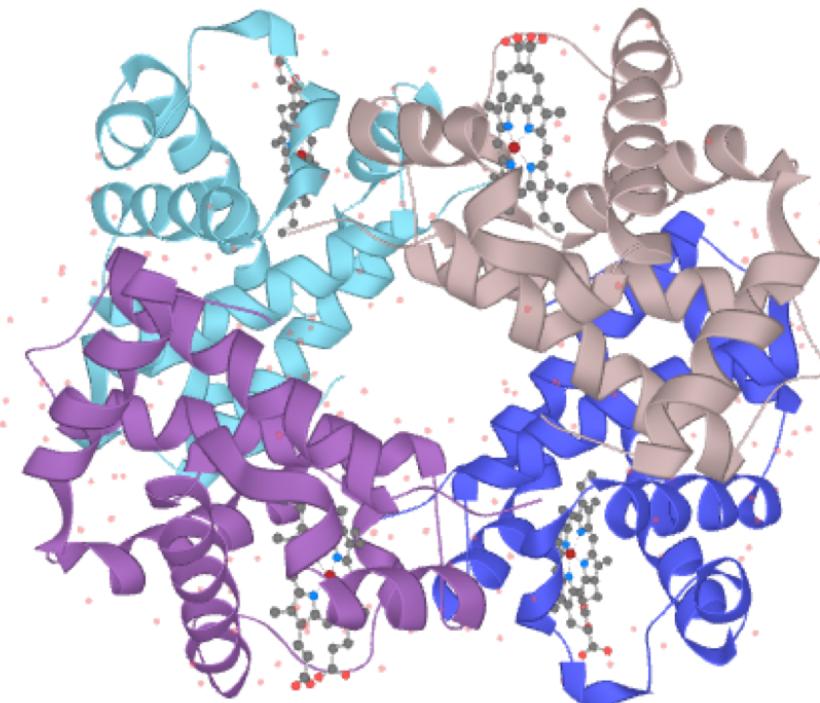
Helix	AAAAAAA HBA_HUMAN	BBBBBBBBB HBB_HUMAN
	-----VLSPADKTNVKAAGKVG-----	HAGEYGAELERMFLSFPTTKTYFPHF
	-----VHITPEEKSAVTALWGKV-----	NVDEVGGEALGRLLVVYPWTQRFFESF
MYG_PHYCA	-----VISEGEWQLVLHWAKVEA-----	DVAGHGQDILIRLFKSHPETLEKFDRF
GLB3_CHITP	-----LSADQ1STVQASFDKVKG-----	DPVGILYAVFKADPSIMAKFTQF
GLB5_PETMA	PIVDTGSVAPLSAAEKTKIRSAWAPVYS-	TYETSGVDILVKFFTSTPAAQEFPKF
LGB2_LUPLU	-----GAI TESQAALVKSSWEFFNA-----	NIPKHTRFFILVLEIAPAAKDLFS-F
GLB1_GLYDI	-----G ISAAQRQVIAATWKDIAGADNGAGVGKDCLIKFLSAHPQMAAVFG-F	
Consensus	I s . . . v a W k v . . .	g . L .. f . P . F F

Helix	DDDDDDDD HBA_HUMAN	EEEEEEEEE HBB_HUMAN	FFFFFFFFF MYG_PHYCA
	-DLS-----HGSAQVKGHGKKVADALTA NAVAHV-----	D--DMPNALSALS DLHAKL-	
	GDLSTPDAMCNP KVKAHGKKVLGAFSDGLAH-----	D--NLKGTFATLSELHCDKL-	
MYG_PHYCA	KHLTEAEMKASEDLKKHGTVLTALGAILKK-----	K-GHHEAHLKPLAQSHATKH-	
GLB3_CHITP	AG-EDLESIKGT TAPFETHANRIVGFFSKIIGE-----	P---NIEADVNTFVASHKPRG-	
GLB5_PETMA	KGLTTADQLKK SADVRWHAERII NAVNDAVASM-----	DDTEKMSM LRLDSGKHAKSF-	
LGB2_LUPLU	LK-CTSEVPQN NPELQAHAGKV FKLVYEAAIQLQVTGV VVTDATLKNLGSV HVSKG-----		
GLB1_GLYDI	SG- AS-----DPGVA ALGAKVLA QICGVAVSH-----	GDEGK MVAGMKAVGV RHKGYGN-----	
Consensus	. t . . . v .. Hg k v . a . . .	d . a 1. 1 H .	

Helix	FFGGGGGGGGGGGGGGGGGG HBA_HUMAN	HHHHHHHHHHHHHHHHHHHHHHHHHHHH HBB_HUMAN
	-RVDPVNFK LSSHCLL VTAAHLP AEETPAV HASLDKFL ASVSTV LTSKYR-----	
	-HVDOPENF RLLGNV LVCV LAHHFG KETPPV QAAYQK VVAGV ANALAH KYH-----	
MYG_PHYCA	-KIPIK YLEF ISEAII HVLHSR HPGD GADA AQGM NKALE LFRKD IAKYKE IGYQG-----	
GLB3_CHITP	-VTHD DQLNN FRAGF VSYMAHT-----	DHA-G AEEAWG ATLDTFF GMIFSKM-----
GLB5_PETMA	-QVDP QYFK VLA AVIA DTVA AG-----	DAGFE KLM SMIC ILLRSAY-----
LGB2_LUPLU	-VADAH FPVV KEA ILKT IKEV VGA KSEEL NSAW TIAY DEL AI VI K EM NDAA-----	
GLB1_GLYDI	KHI KAQY FEPL GASLL SAME HRIG GKIN NAA AKDA WA AA YAD IS GAL ISGLQS-----	
Consensus	v f l . . . . .	f . aa. k . . 1 sky



**Helices are better conserved than loops between them!**

# Protein family profile: idea

- Some protein regions are better conserved than other
  - Conserved sequence regions often correspond to conserved / structured 3D regions
  - Some residues are particularly conserved
    - They often correspond to functionally important residues
  - **Conserved regions and residues are more important**

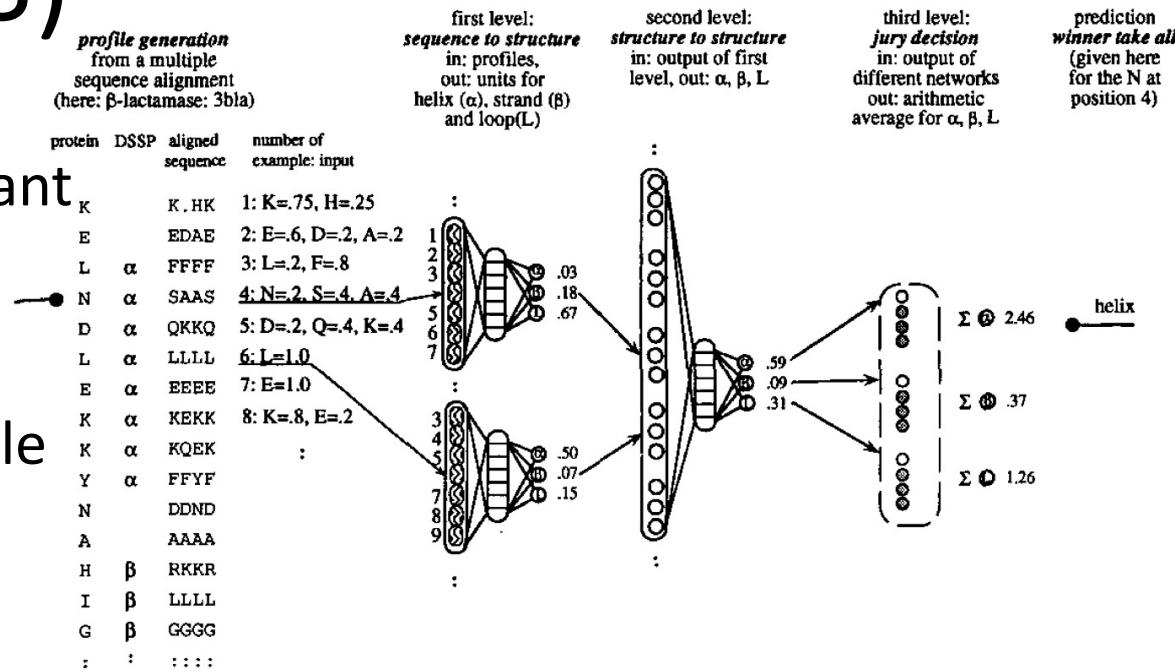
# Calculating protein family profile

- For each amino acid and each alignment column, count the number of occurrences of this amino acid in this column
- Divide by the number of sequences in the alignment
- **How to use it?**
  - For **sequence similarity search** (Lecture 6): can be naturally turned into an HMM
    - Calculate the probability of observing a certain amino acid at a certain position: emission probabilities  $e(x_i)$  = observed frequencies (almost), normalized by background frequencies  $q_{x_i}$
    - Probability of a sequence given a model:  $\prod_{i=1,\dots,L} e_i(x_i)$
    - Transition probabilities: slightly more tricky
    - For **predicting secondary structure** (now)

Helix	AAAAAAAAAAAAAA
HBA_HUMAN	-VLSPADKTNVKAAGKVGA
HBB_HUMAN	-VHLTPEEKSAVTALWGKV--
MYG_PHYCA	-VLSEGEWQLVLHVWAKVEA
GLB3_CHITP	-LSADQ1STVQASFDKVKG
GLB5_PETMA	PIVDTGSVAPLSAAEKTKIRSAWAPVYS
LGB2_LUPLU	-GALTESQAALVKSSWEFNA
GLB1_GLYDI	-GLSAAQRQVIAATWKDIAG
Consensus	Ls.... v a W kv .

# PHD (Rost & Sander, 1993)

- Predicts 3 states: **helix, strand, loop**
- Neural network** trained of **130** non-redundant proteins with known 3D structure
- 3 simple ANNs (1 hidden layer)
- Input:** multiple sequence alignment -> profile
  - Amino acid frequencies
- 7-fold cross-validation
- Overlapping window moves down the sequence:
  - 1st level: independent prediction for central residue
  - 2nd level: combines output from overlapping windows
  - 3rd level: average over several networks
- Overall  $Q_3$  accuracy: **70.8%**



# Quality of prediction

- **$Q_3$  accuracy** = # residues correctly predicted / total # residues
  - Expected to be less than 100%, since SSE can only be identified with a certain uncertainty
  - Does not perfectly represent the quality of prediction

**VLHQASGNSVILFGSDVTVPGATNAEQAR**

HHHHHCCCCEEEECCCIEEECCCCCHHHHHH — actual SSEs

CHHHHCCCEEEECCCCCEEECCCHHHHHH — prediction 1,  $Q_3 = 76\%$ , useful

HHHHHCCCCHHHHCCCHHHCCCCCHHHHHH — prediction 2,  $Q_3 = 76\%$ , bad

- **Sov** (segment overlap):  $\frac{1}{N} \sum_s \frac{\text{minov}(s_1, s_2) + \delta}{\text{maxov}(s_1, s_2)} \cdot \text{length}(s_1)$

**VLHQASGNSVILFGSDVTVPGATNAEQAR**

HHHHHCCCCEEEECCCIEEECCCCCHHHHHH — actual SSEs

CHHHHCCCEEEECCCCCEEECCCHHHHHH — prediction

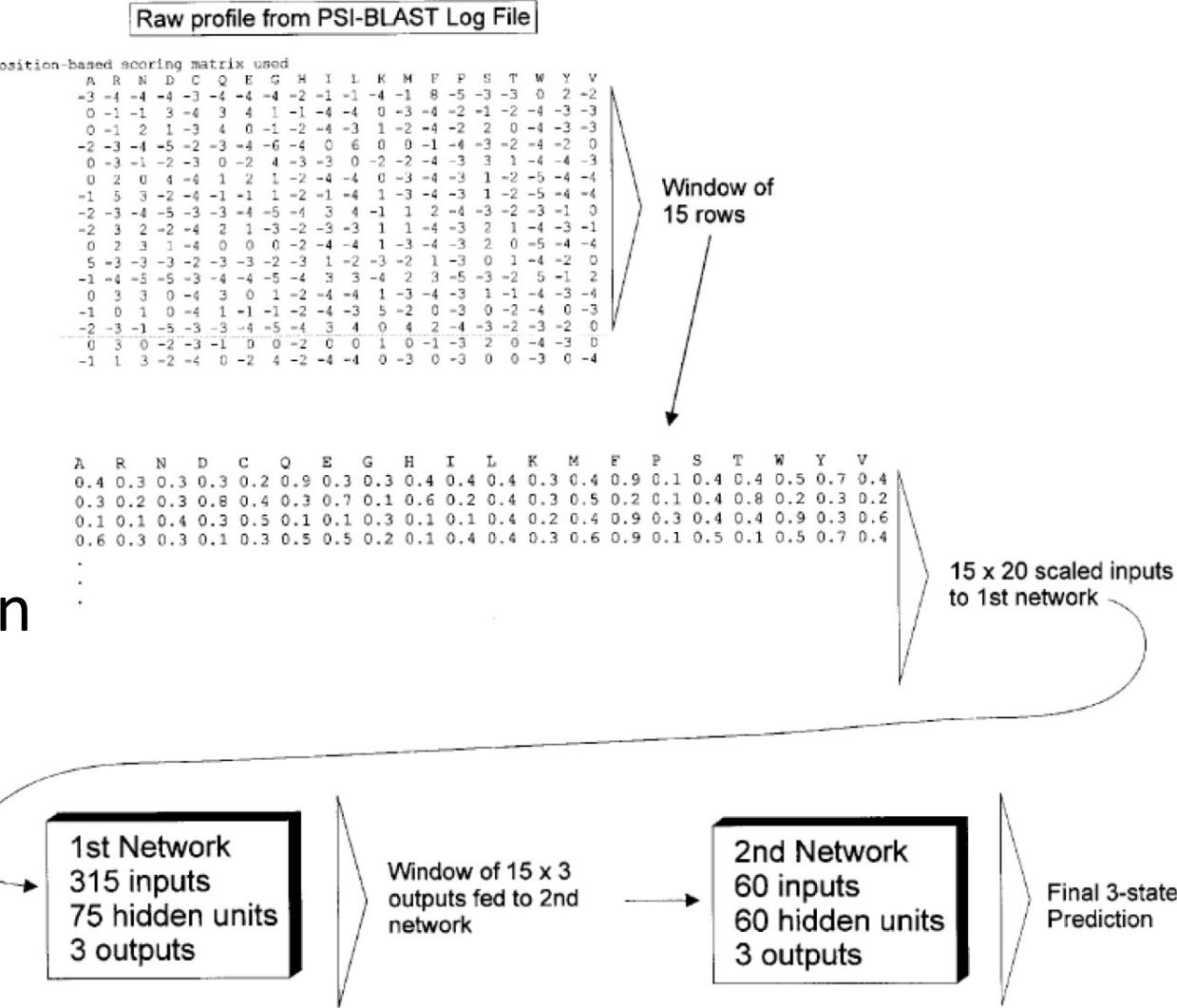
>|||<  
>||||<

— minov

— maxov

# PSIPRED (Jones, 1999)

- PSI-BLAST profiles on a non-redundant set of sequences of globular proteins
- Neural network with a single hidden layer
- Second neural network to average over positions in the window
- 3-fold cross-validation
- Q<sub>3</sub> accuracy: **76.5%**; Sov<sub>3</sub>: **73.5%**



# Nowadays: deep learning

- Deep learning models can be trained to predict secondary structure
  - Usually the same 3 outputs: helix, strand, loop (coil)
- Not usually used on their own, but as a part of models to predict protein 3D structure *ab initio* (Lecture 11)

# Summary and possible exam questions

- What is the residue pattern in the coiled-coil regions?
- What residues are preferred in transmembrane regions? Why?
- What is a Markov chain?
- What is a hidden Markov model?
- What are hidden and observable states, transmission and emission probabilities?
- How hidden Markov models are used to predict transmembrane regions in proteins?
- What amino acids are avoided in regular secondary structure elements, such as  $\alpha$ -helices and  $\beta$ -sheets? Why?
- What is a protein family profile? How is it calculated?
- How is the quality of prediction of secondary structure measured?