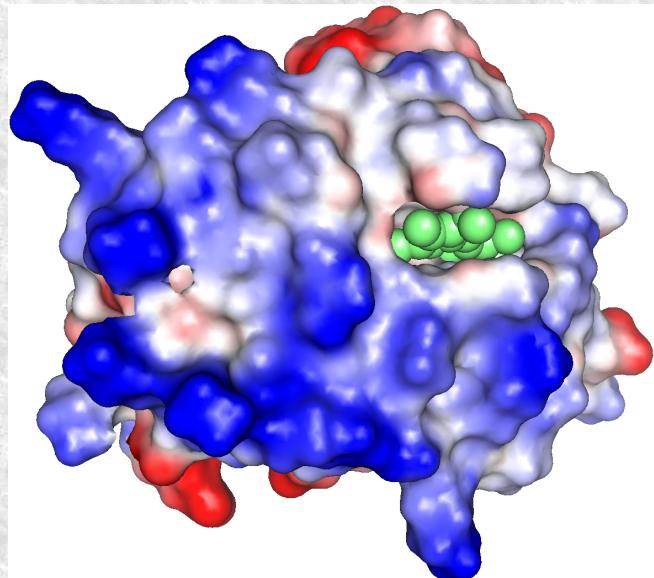


Modeling and predicting protein structure



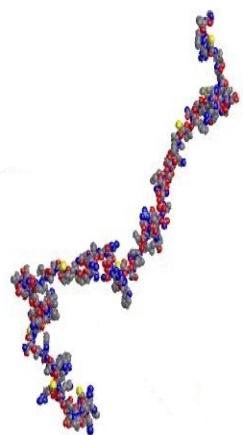
Mastères *Ingénierie et Chimie des BioMolécules et Biologie et Santé*
Module Bioinformatique Structurale
novembre 2022
T. Simonson, Ecole Polytechnique

Importance of modeling

For most known proteins, 3D structure is unknown.

For ~1/3, function is unknown.

A
U
G
C
G
C
U
U
A
U
A
G
C
C
A
A
G
G
:

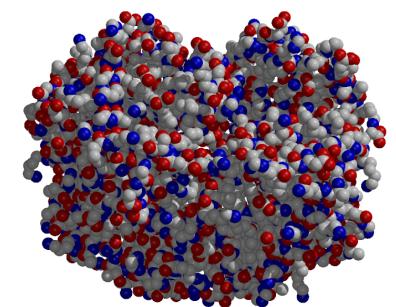


50% of our dry weight...

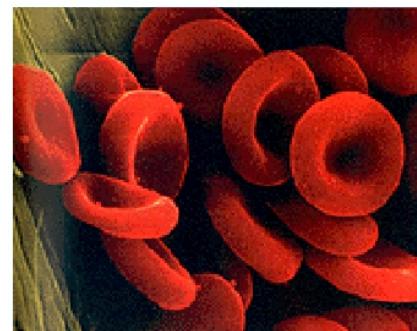
Sequence



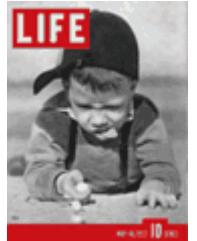
Structure



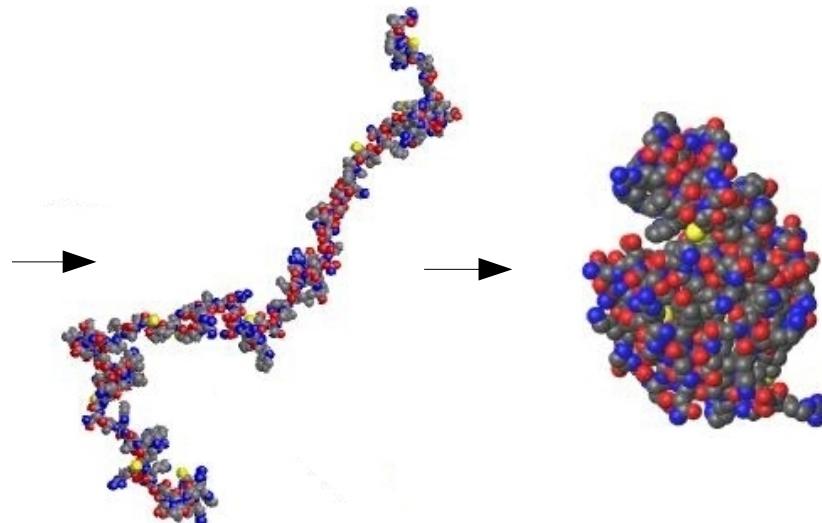
Function



The folding problem



K
L
H
G
G
P
M
L
D
S
D
Q
K
F
W
R
T
P
A
A
L
H
Q
N
E
G
F
T



$$N_{\text{états}} \sim 10^n$$

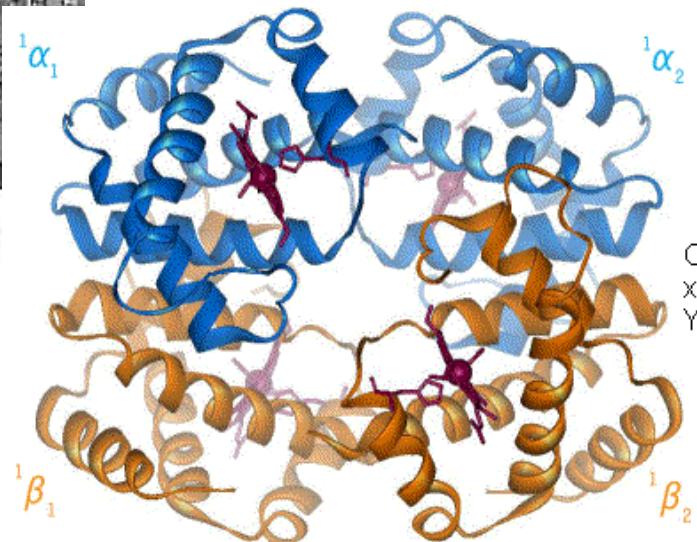
$n = 100-300$

googols of googols....

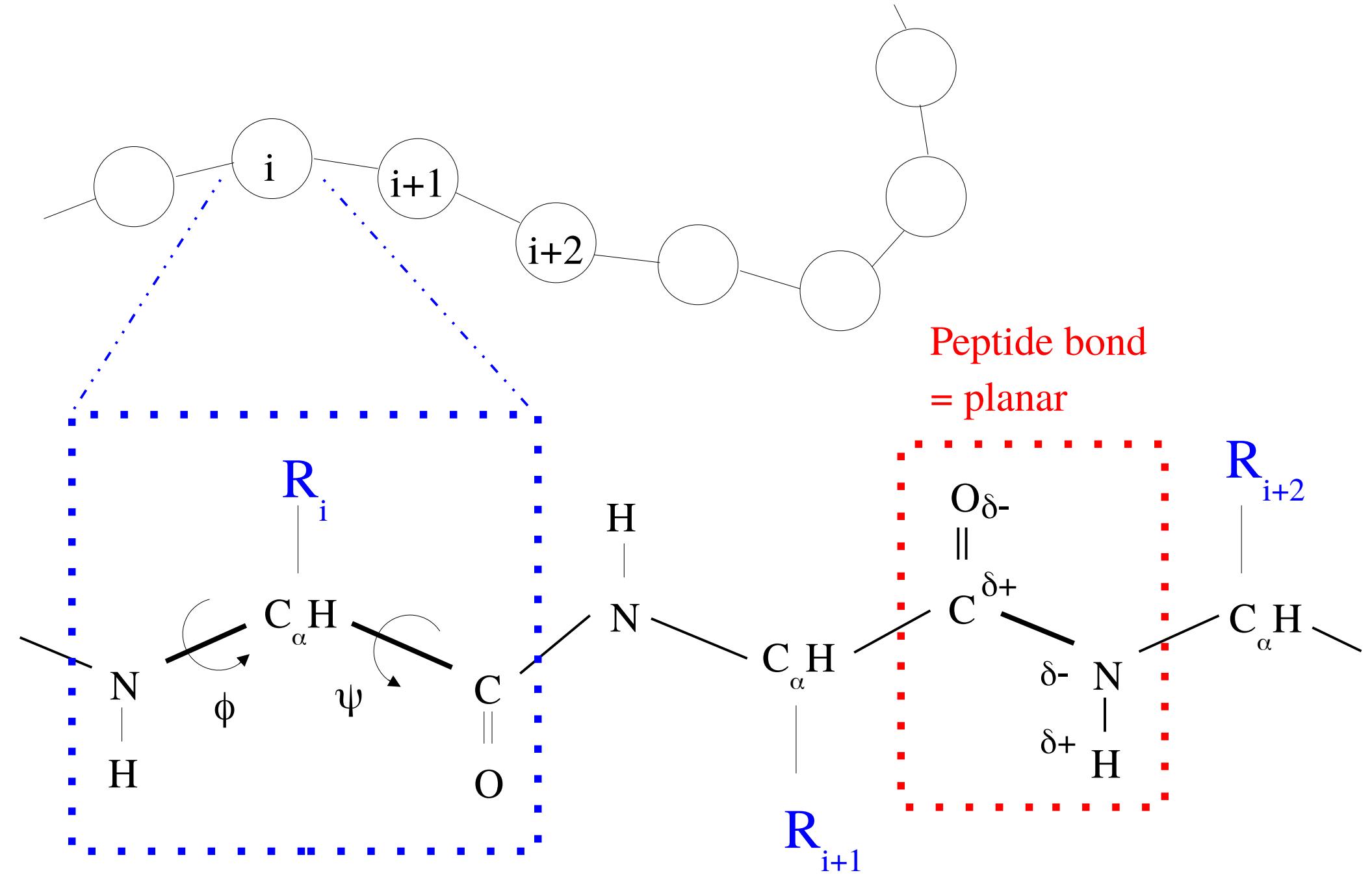
Levinthal paradox

- Protein structure and stability
- Molecular modeling, molecular dynamics
- *Predicting secondary structure*
- Predicting 3D structure: homology modeling

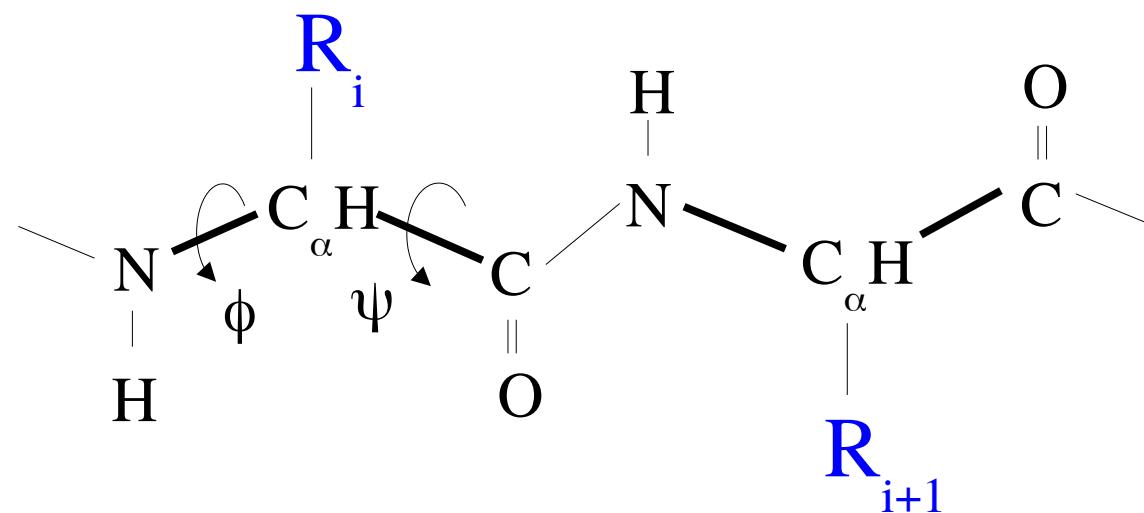
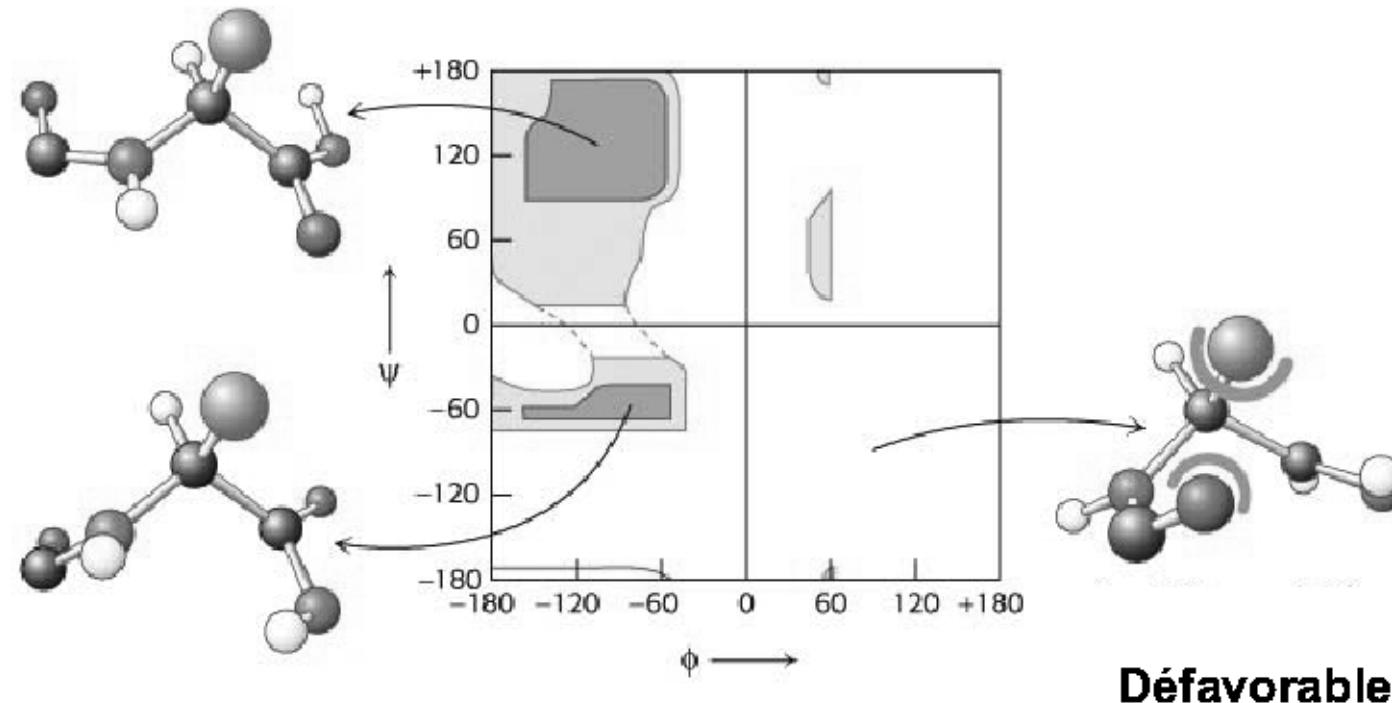
Protein structure



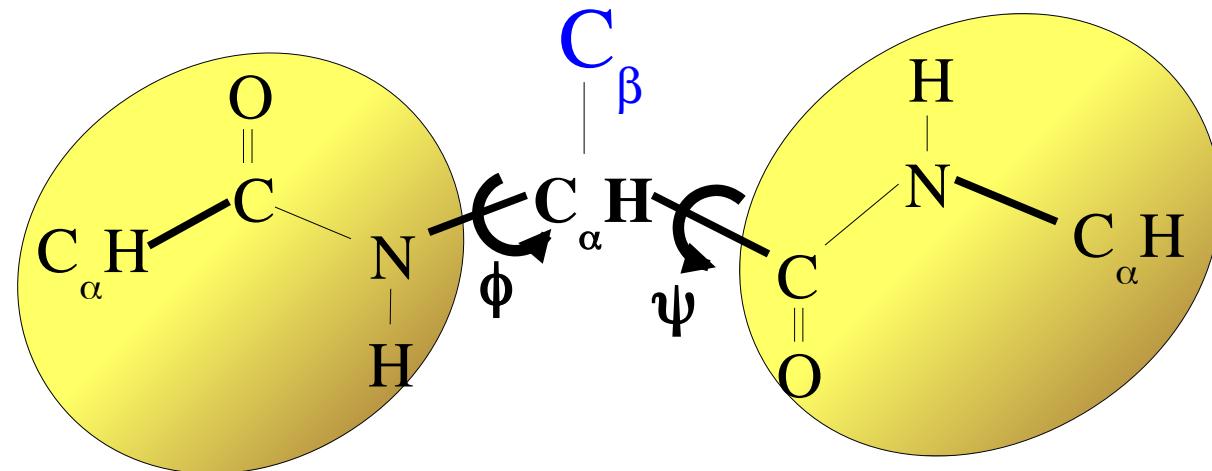
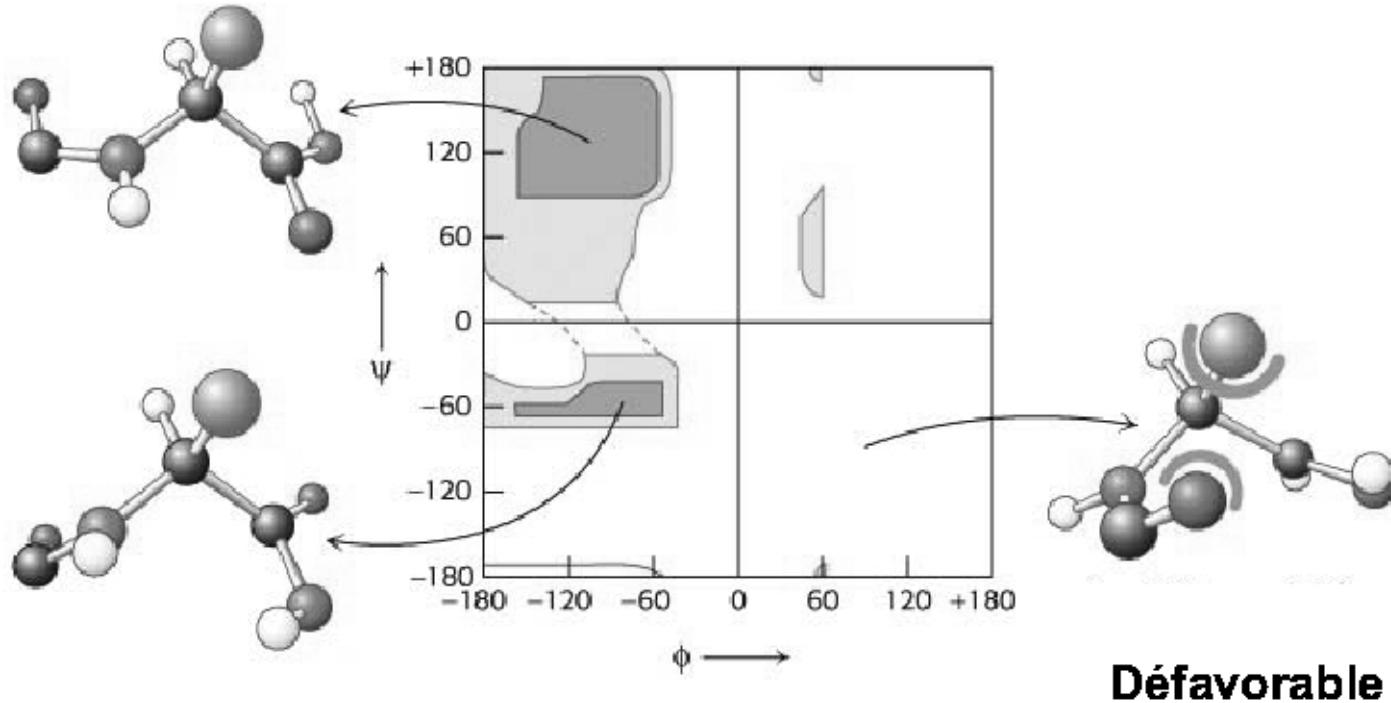
The main chain has two flexible dihedrals per amino acid



Sterically allowed ϕ , ψ values are those of the Ramachandran plot



Sterically allowed values are those of the Ramachandran plot



Lovell et al, 2003,
Proteins, 50:437

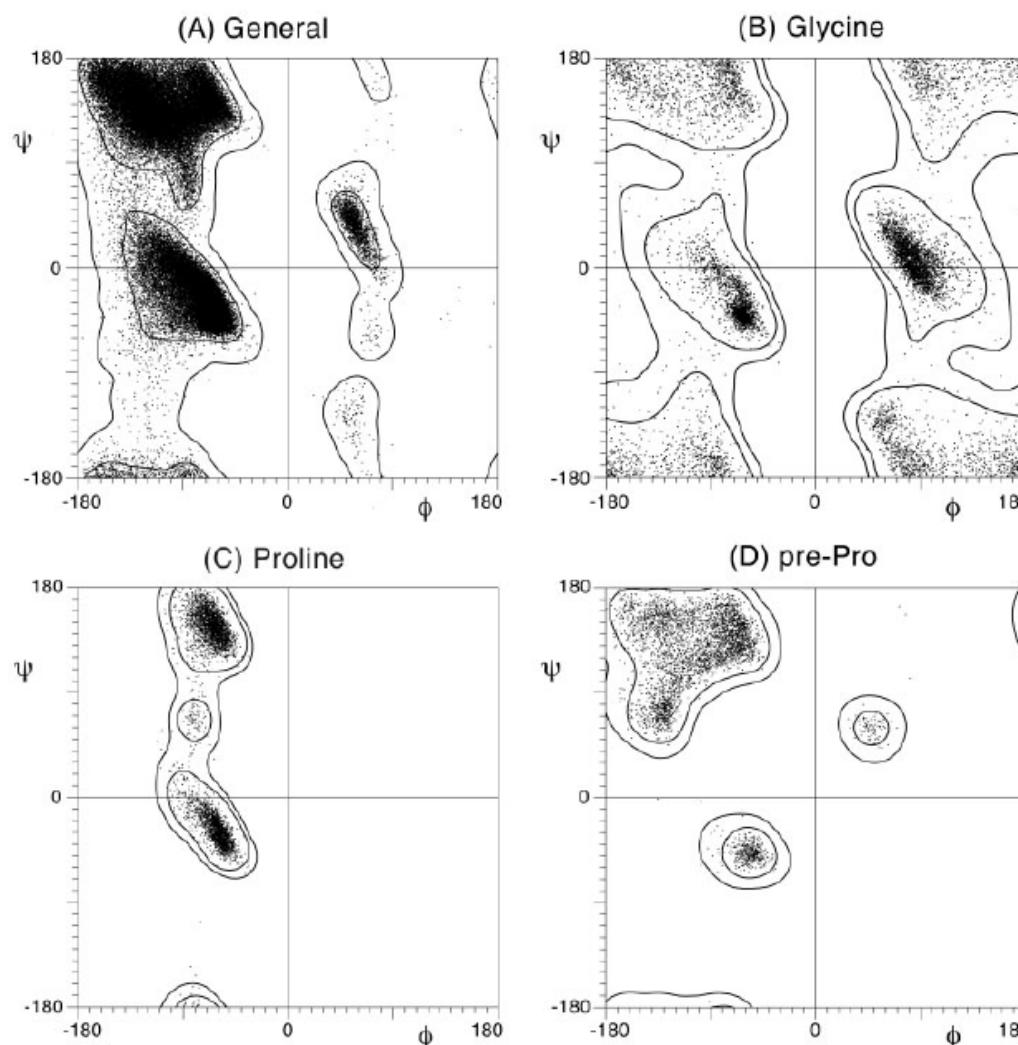
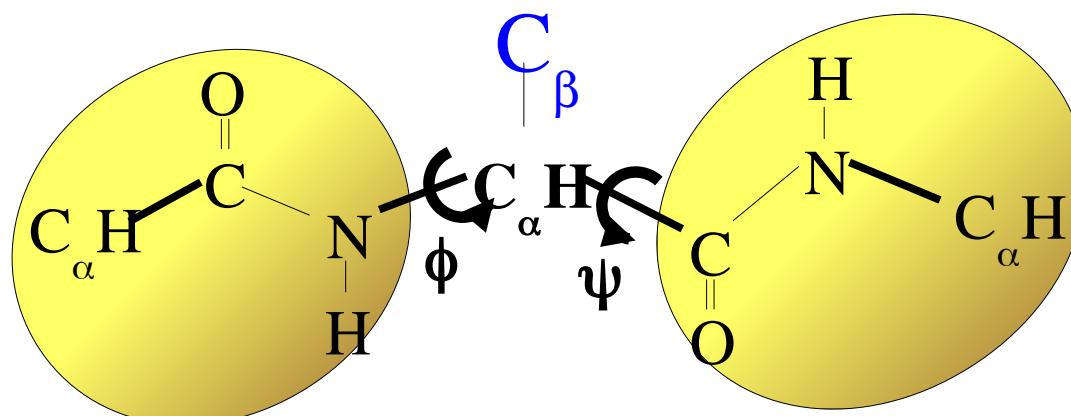
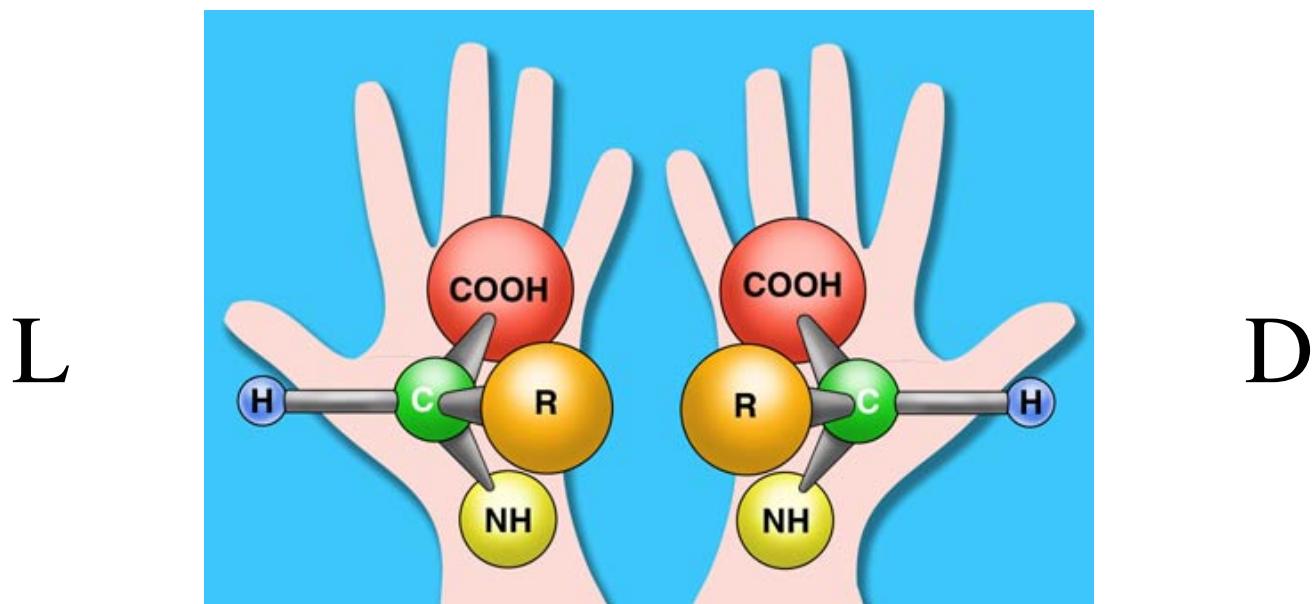


Fig. 4. ϕ, ψ angle distributions for 97,368 residues with backbone B -factor < 30 from the 500-structure high-resolution database, along with validation contours for favored and allowed regions. **a:** The general case of 81,234 non-Gly, non-Pro, non-prePro residues. **b:** The 7705 Gly residues, shown with twofold symmetrized contours. **c:** The 4415 Pro residues with contours. **d:** The 4014 pre-Pro residues (excluding those that are Gly or Pro) with contours.



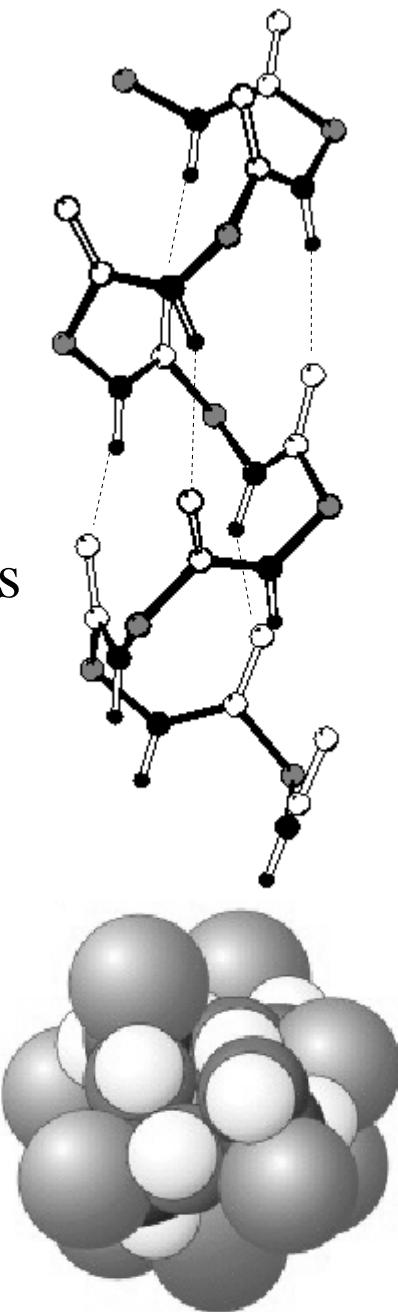
Nature chose L amino acids



“L” not the same as “levorotary”

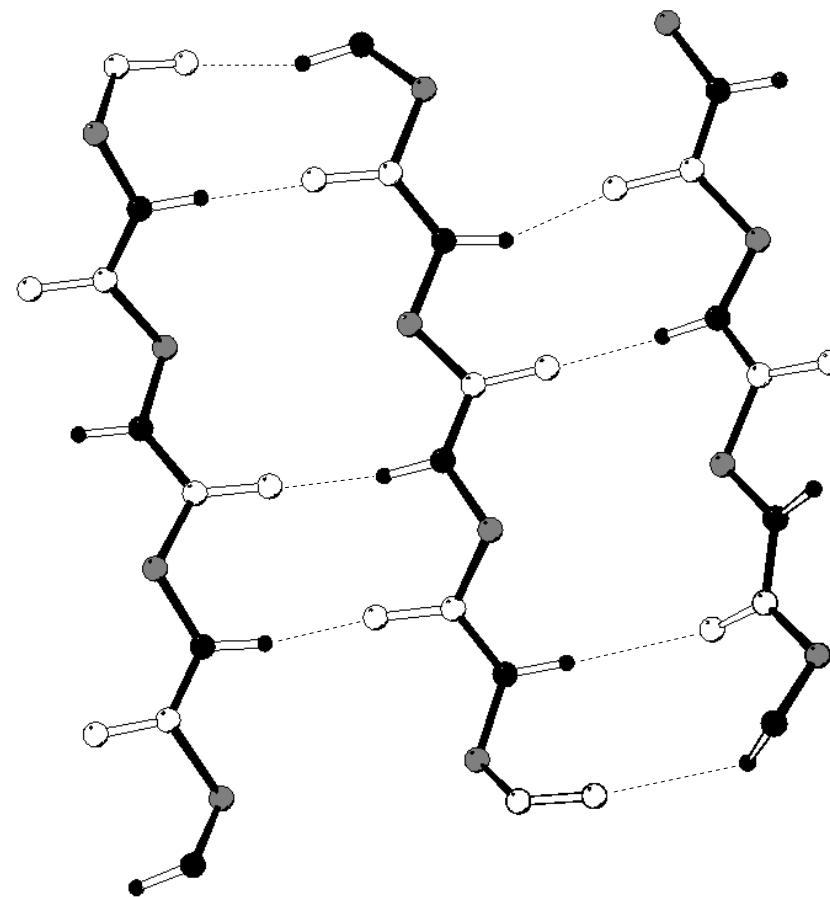
In helices and sheets, polar groups form hydrogen bonds

Hélice α



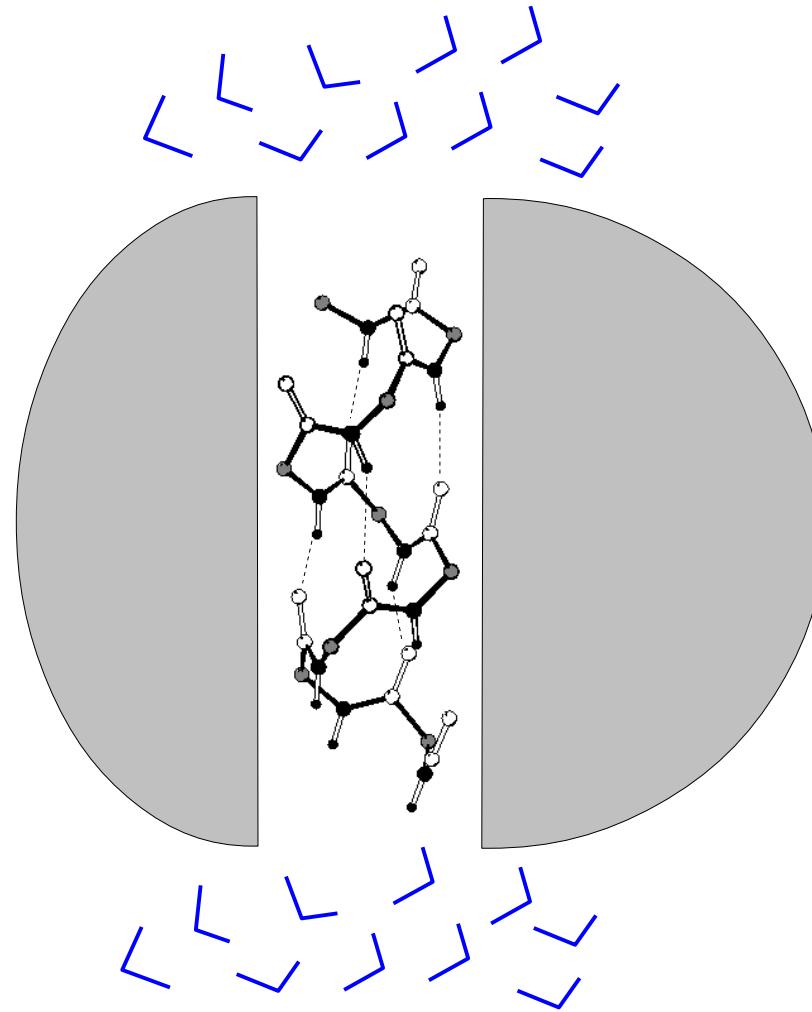
3.6 residues
per turn

Feuillet β (antiparallel)



Pseudo-period of 2

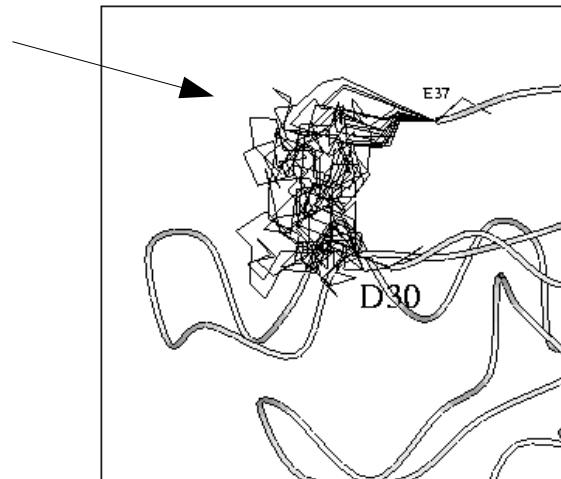
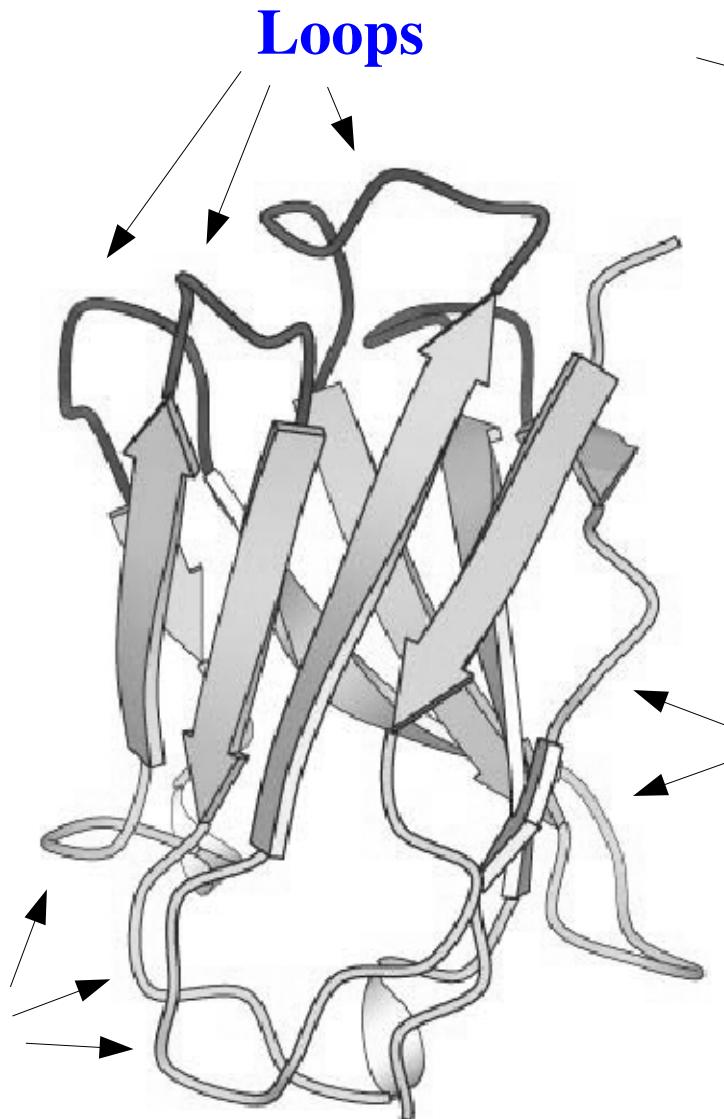
Helices and sheets are extended structures that can *traverse* the hydrophobic core, while forming all possible hydrogen bonds



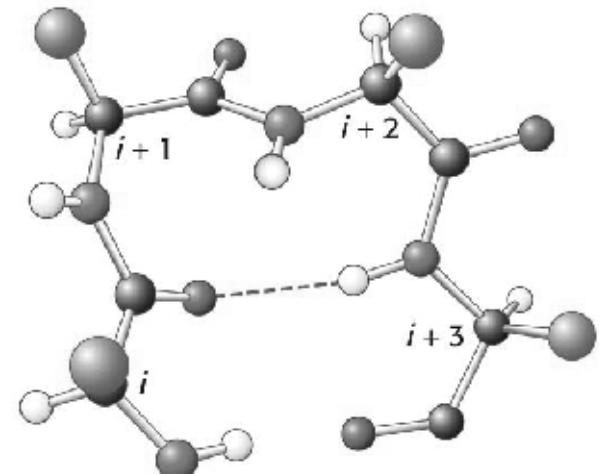
Helices +
sheets =
over 60% of
amino acids

Predicted by Corey & Pauling *before* the experimental structures

Loops are flexible, less conserved, and harder to predict

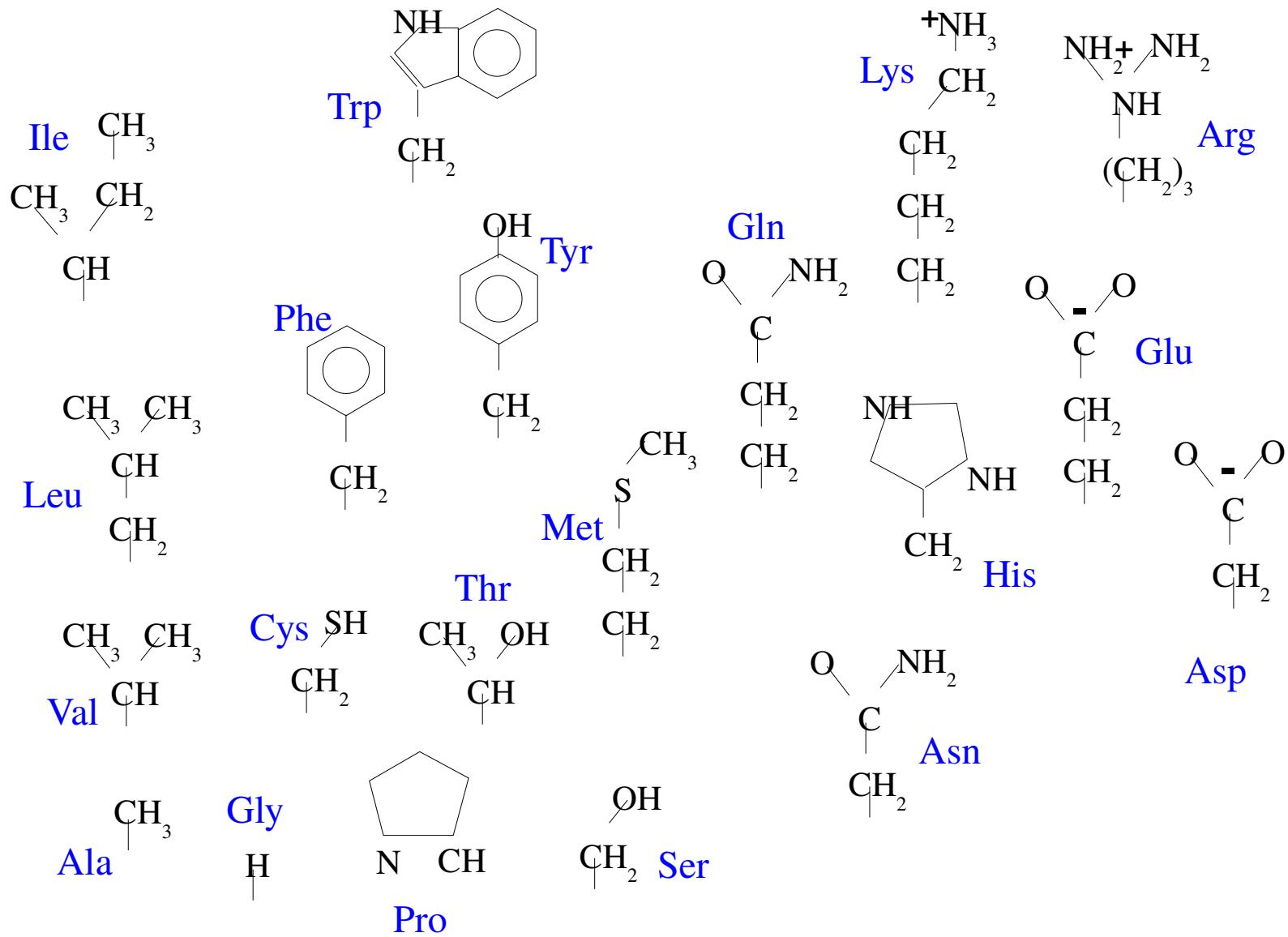


turn

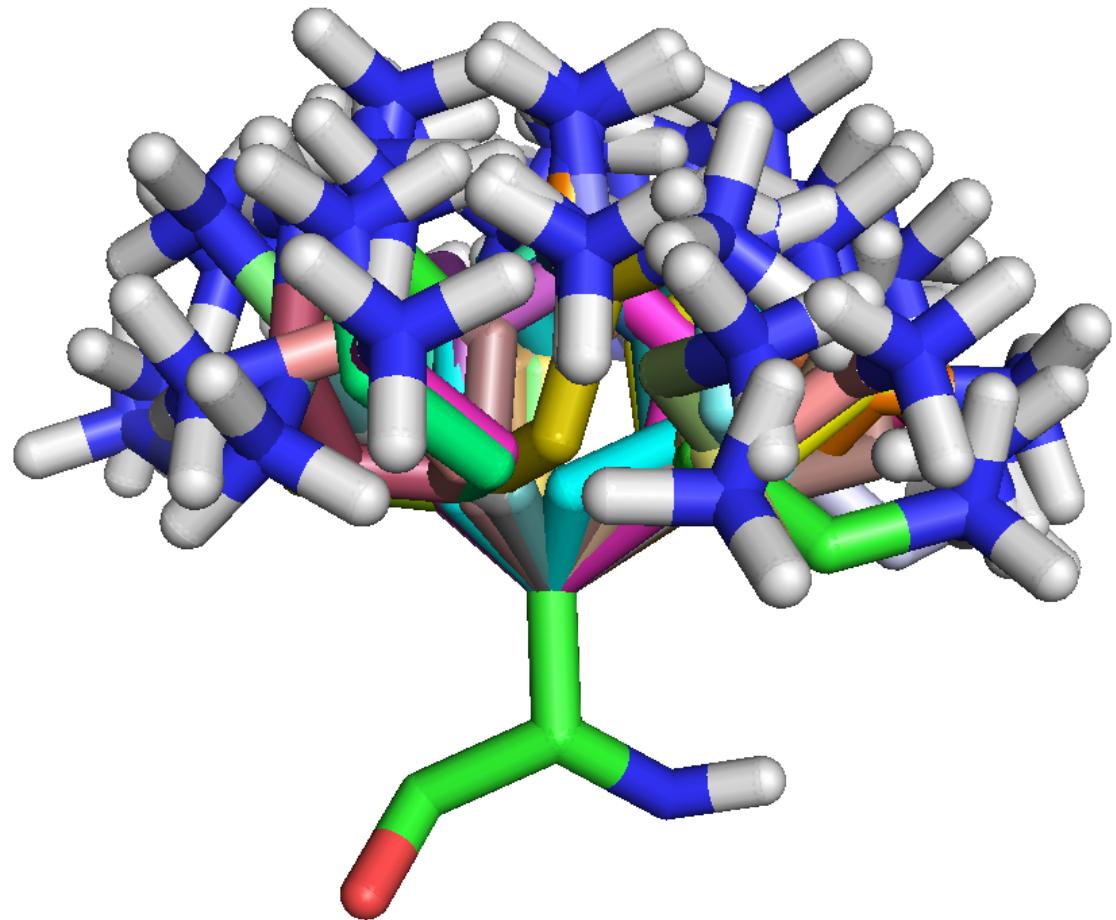
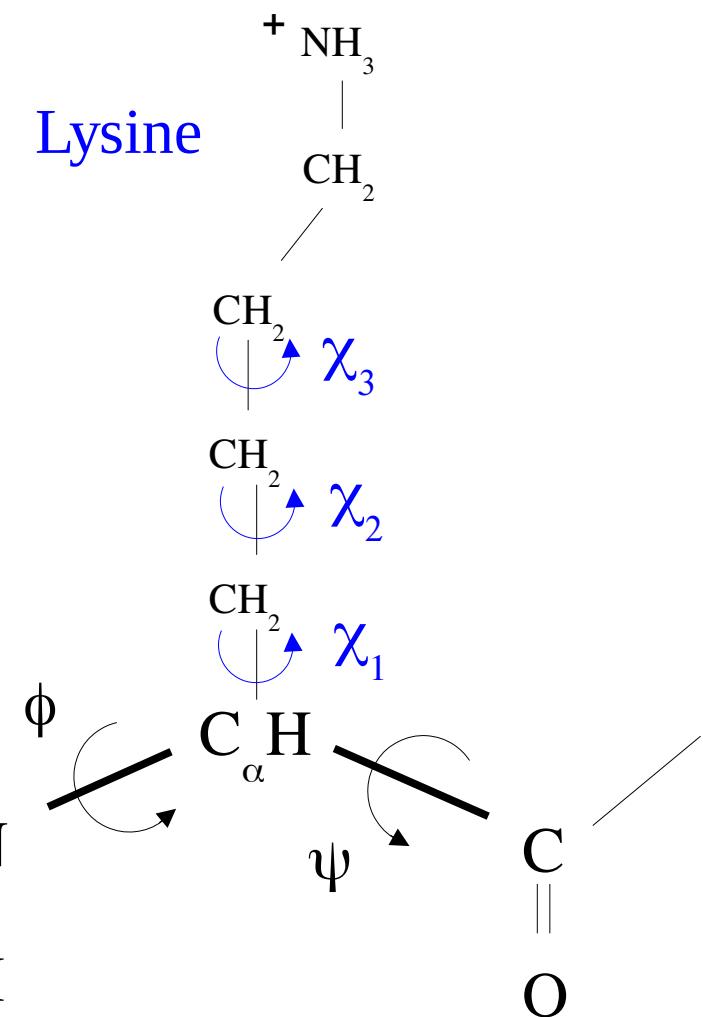


~ 1/3 of amino acids

The amino acids have a range of physico-chimical and structural properties

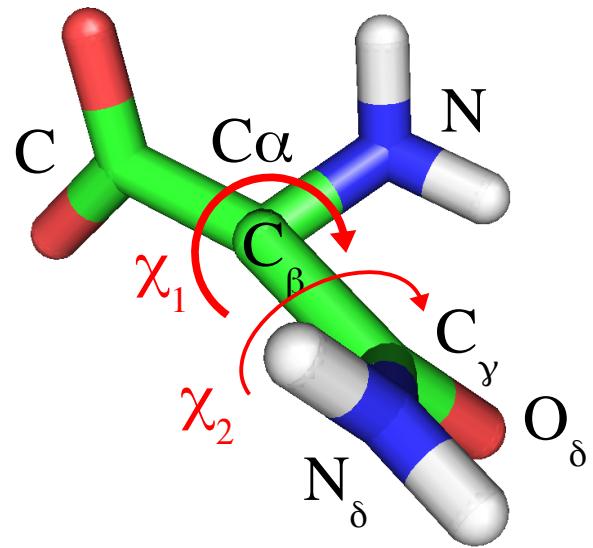


Side chains have flexible torsion angles

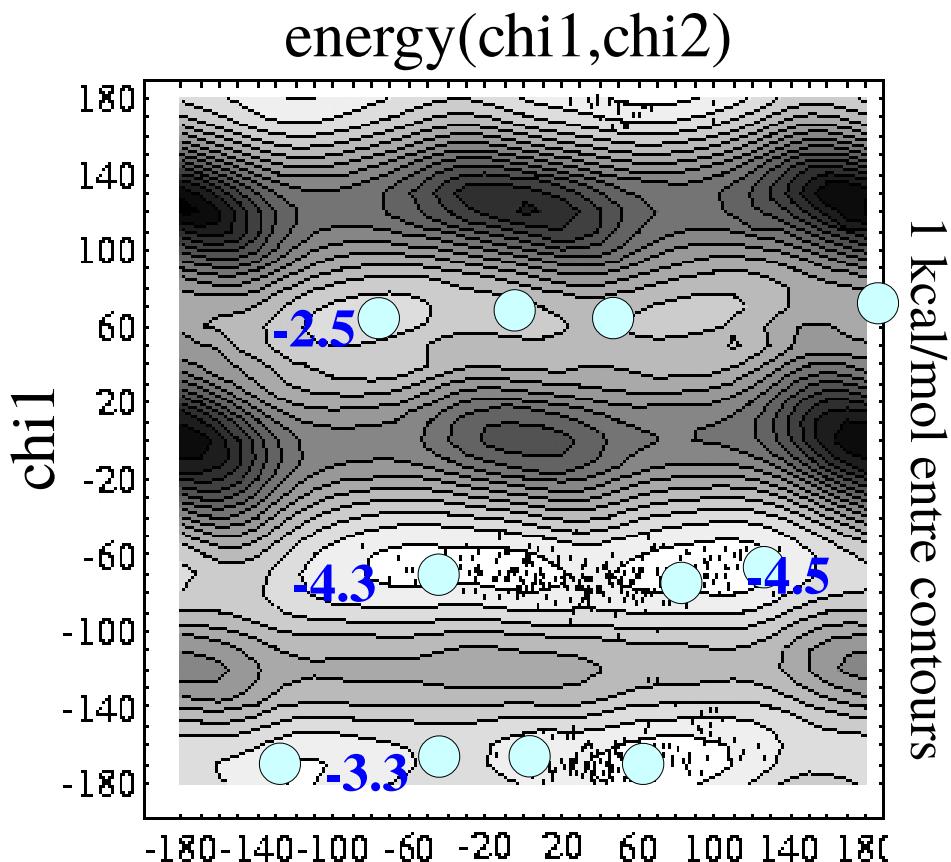


53 conformations of Lys seen in
X-ray structures

Side chains have preferred conformations, or “rotamers”

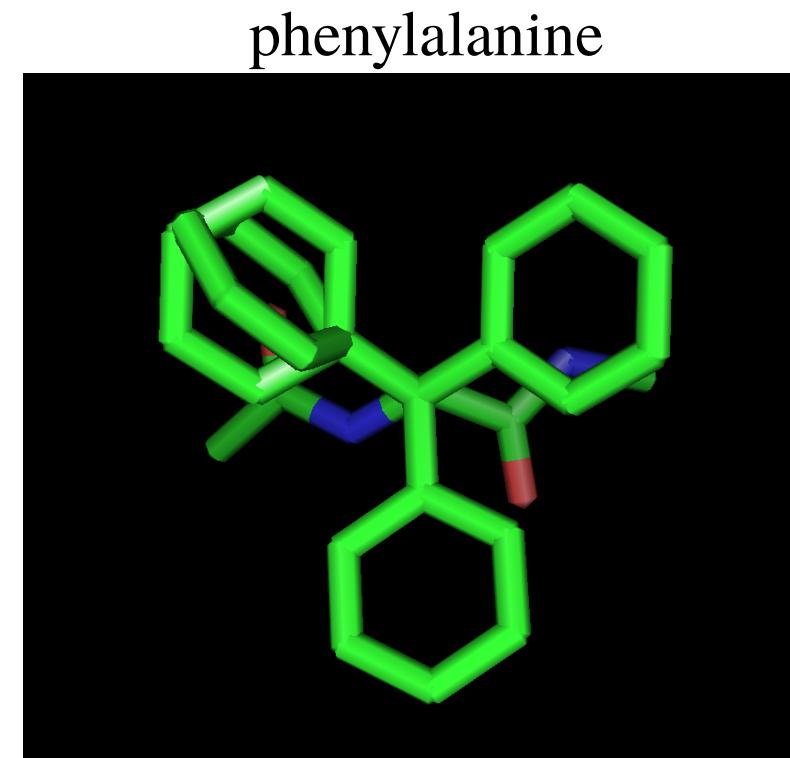
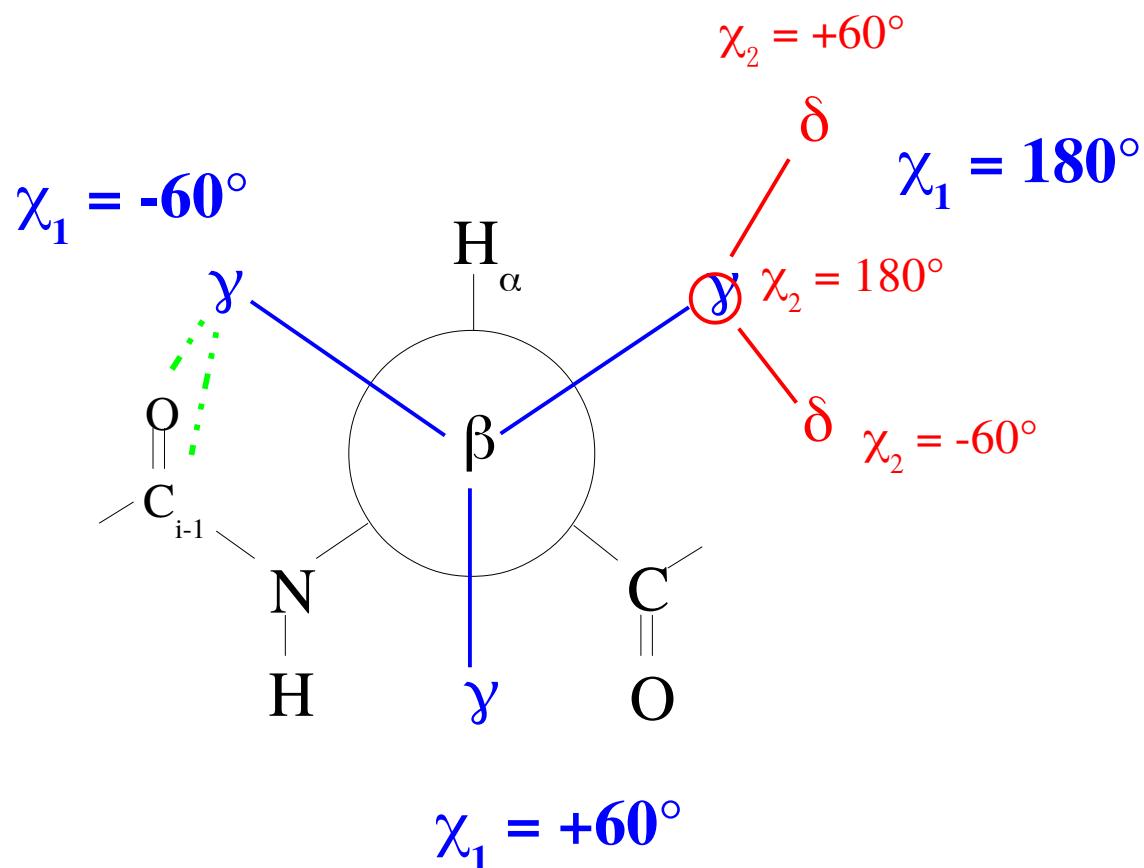


asparagine



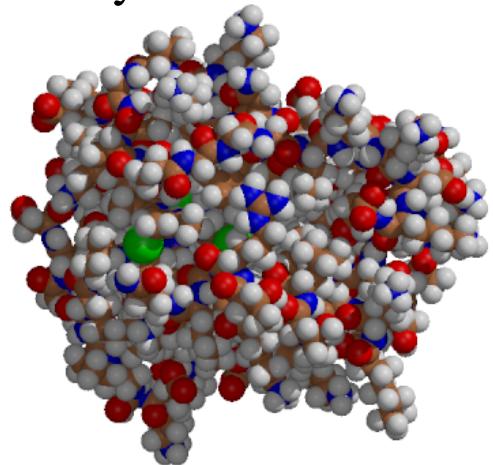
- conformations typical of proteins
- conformations from a simulation

Discrete library of side chain conformations

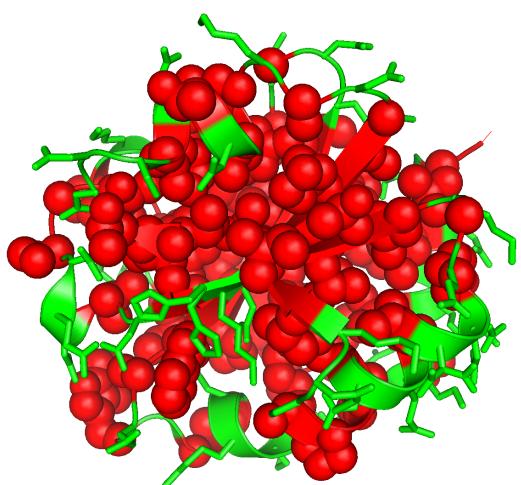
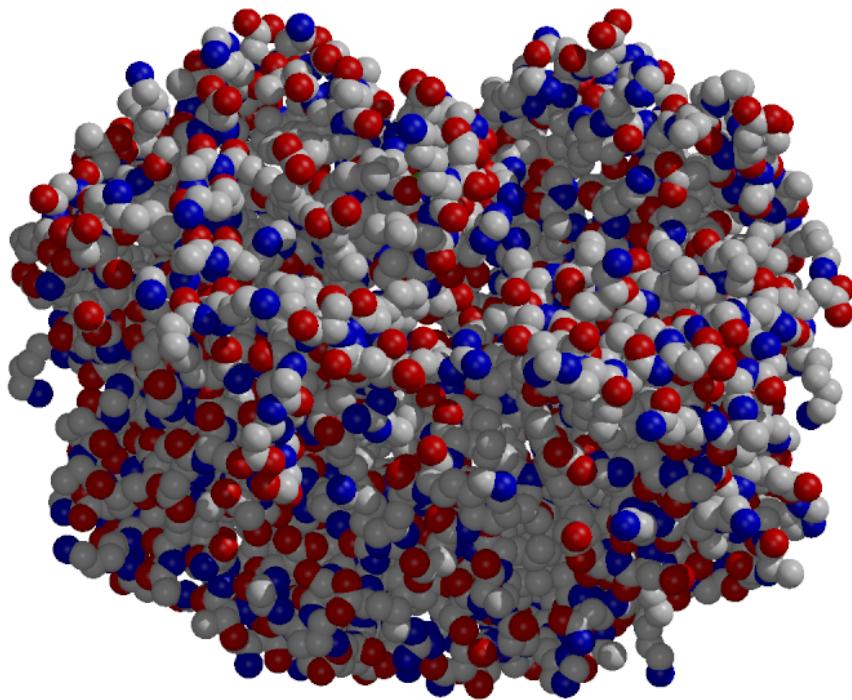


Tertiary structure is globular, with a polar surface and a hydrophobic core

Cytochrome c



Hemoglobin

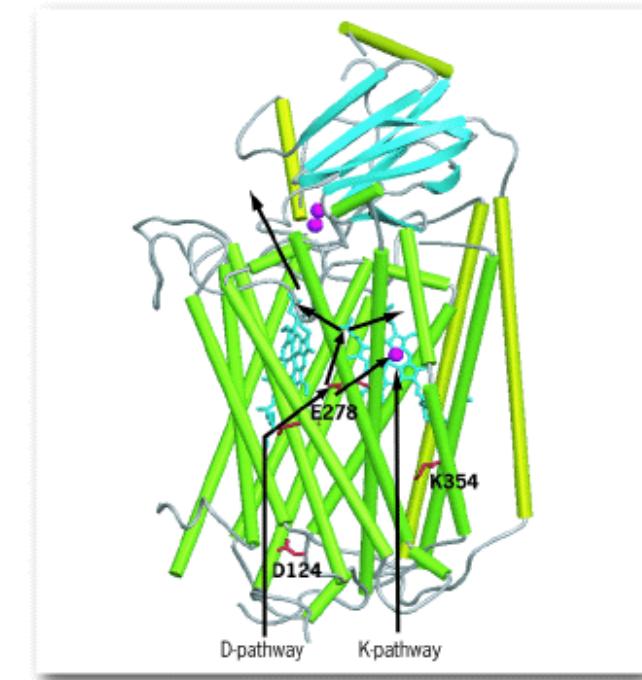
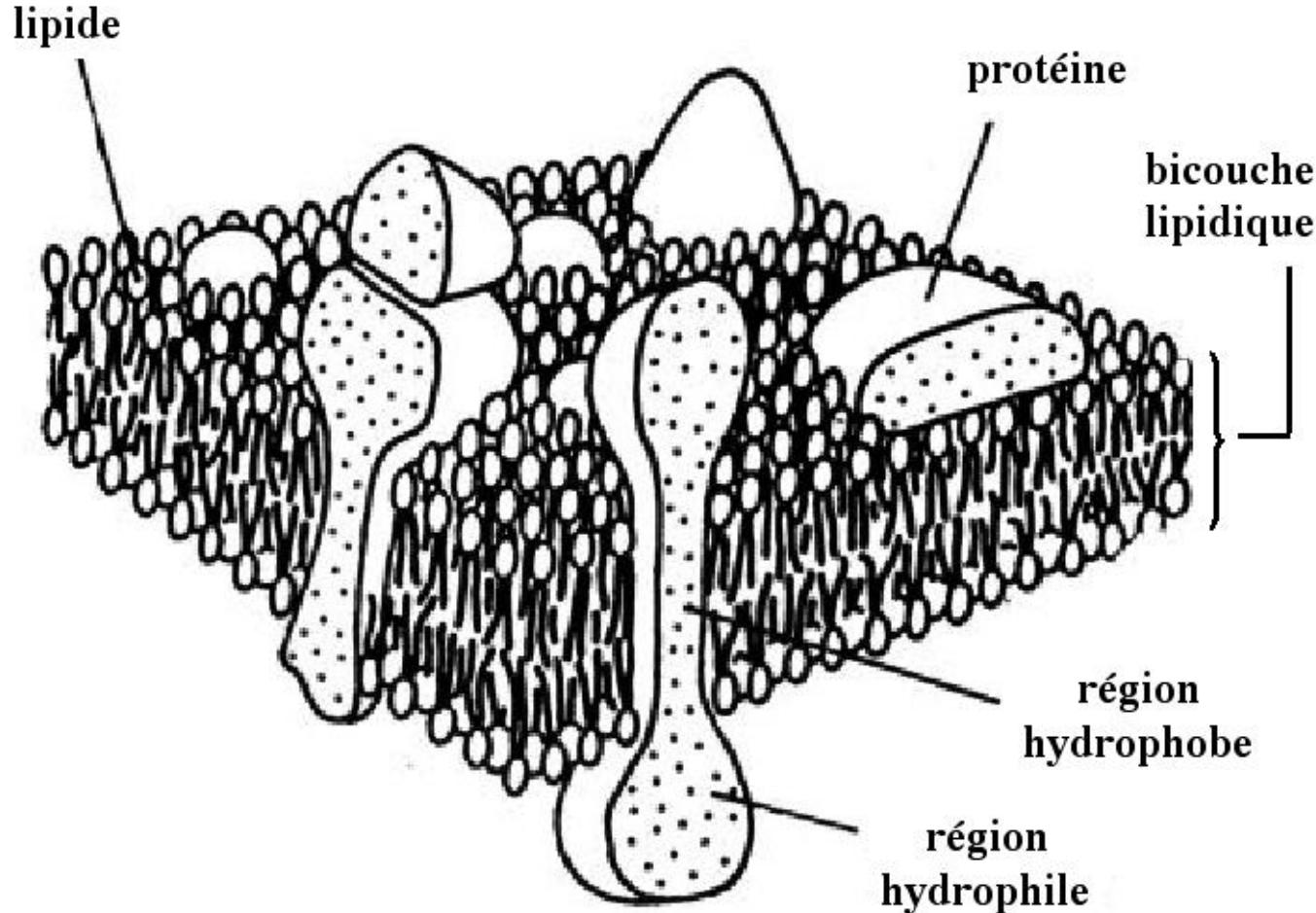


Thioredoxin:
red atoms belong to
hydrophobic side
chains



water

Membrane proteins are a separate case



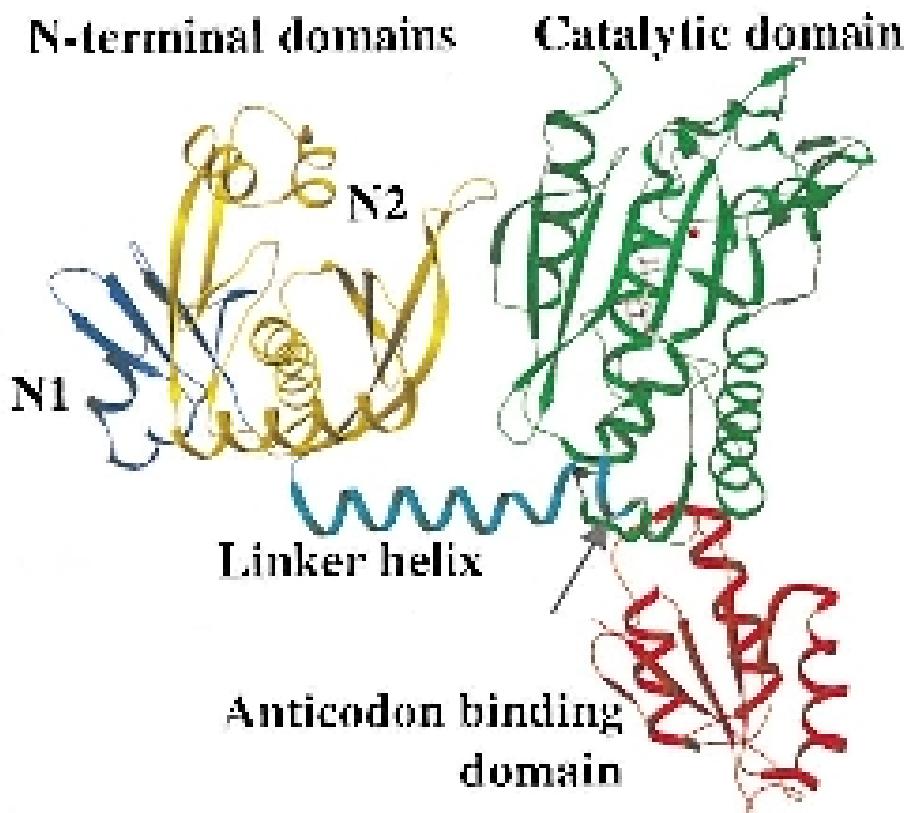
Cytochrome oxidase

~ 30% of our genome humain, ~ 50% of drugs

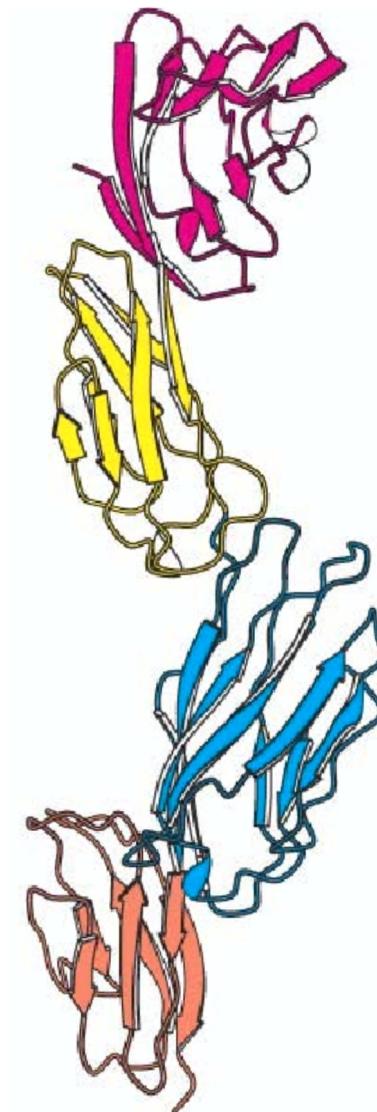
Few known structures (~100).

Few architectures.

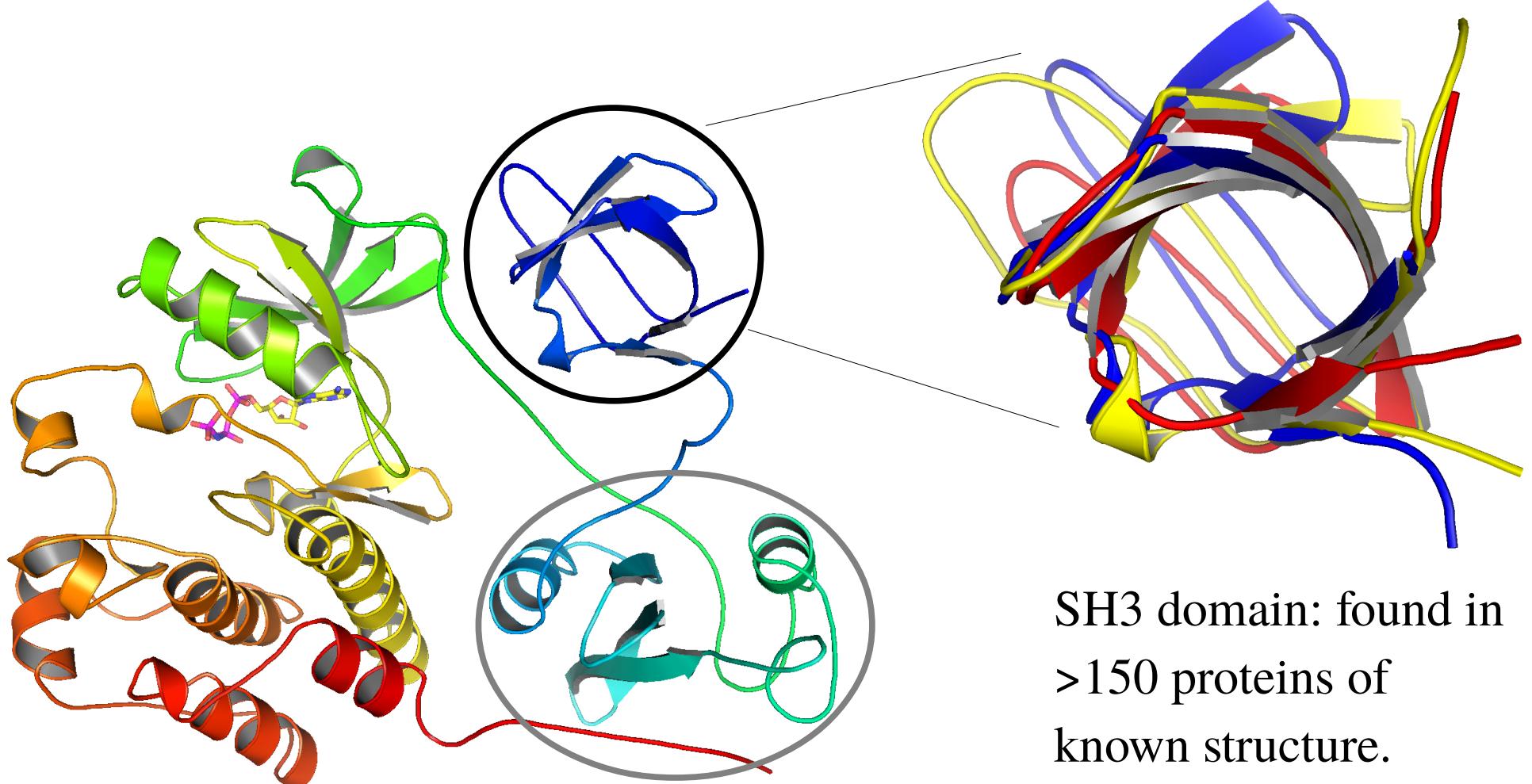
Quaternary structure is a 3D arrangement of structural domains



100-300 amino acids



One domain can appear in many proteins

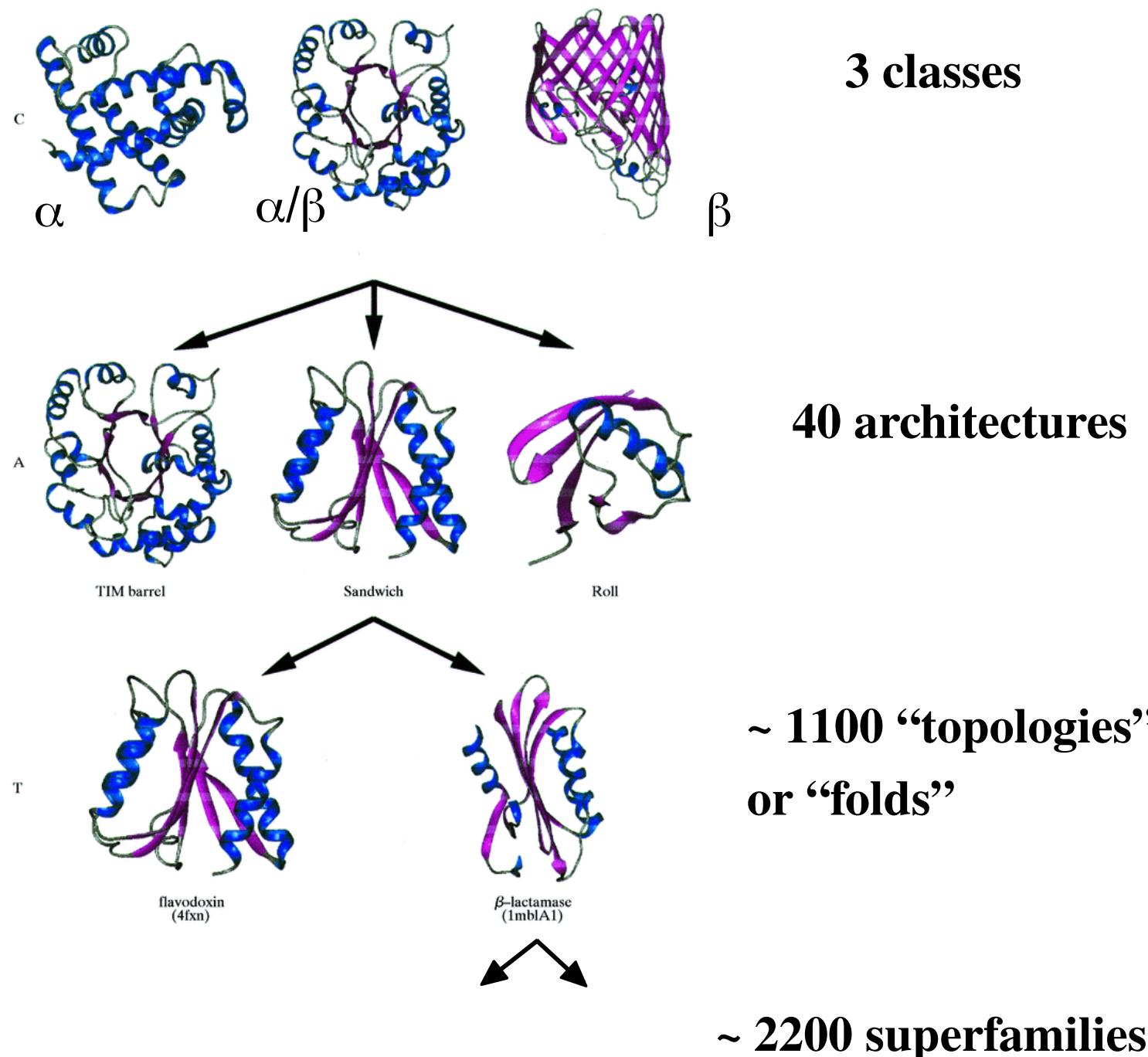


Tyrosine kinase c-Src

SH3 domain: found in
>150 proteins of
known structure.

23 in yeast; do not cross-react

Structural domains: hierarchical classification



Classifying a Protein in
the **CATH Database**
of Domain Structures
Acta Cryst (1998) D54:1155
Orengo, Martin, Hutchinson,
Jones, Jones, Michie,
Swindells, Thornton

Similarities reflect evolution

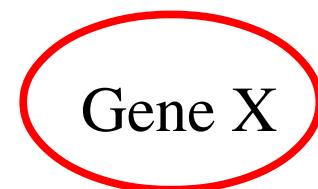
Orthologs



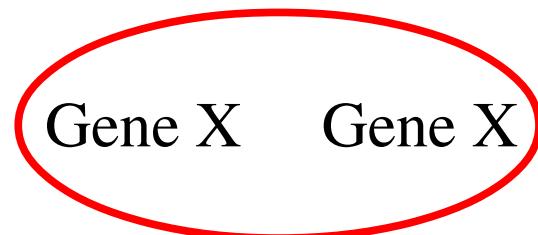
speciation



Paralogs



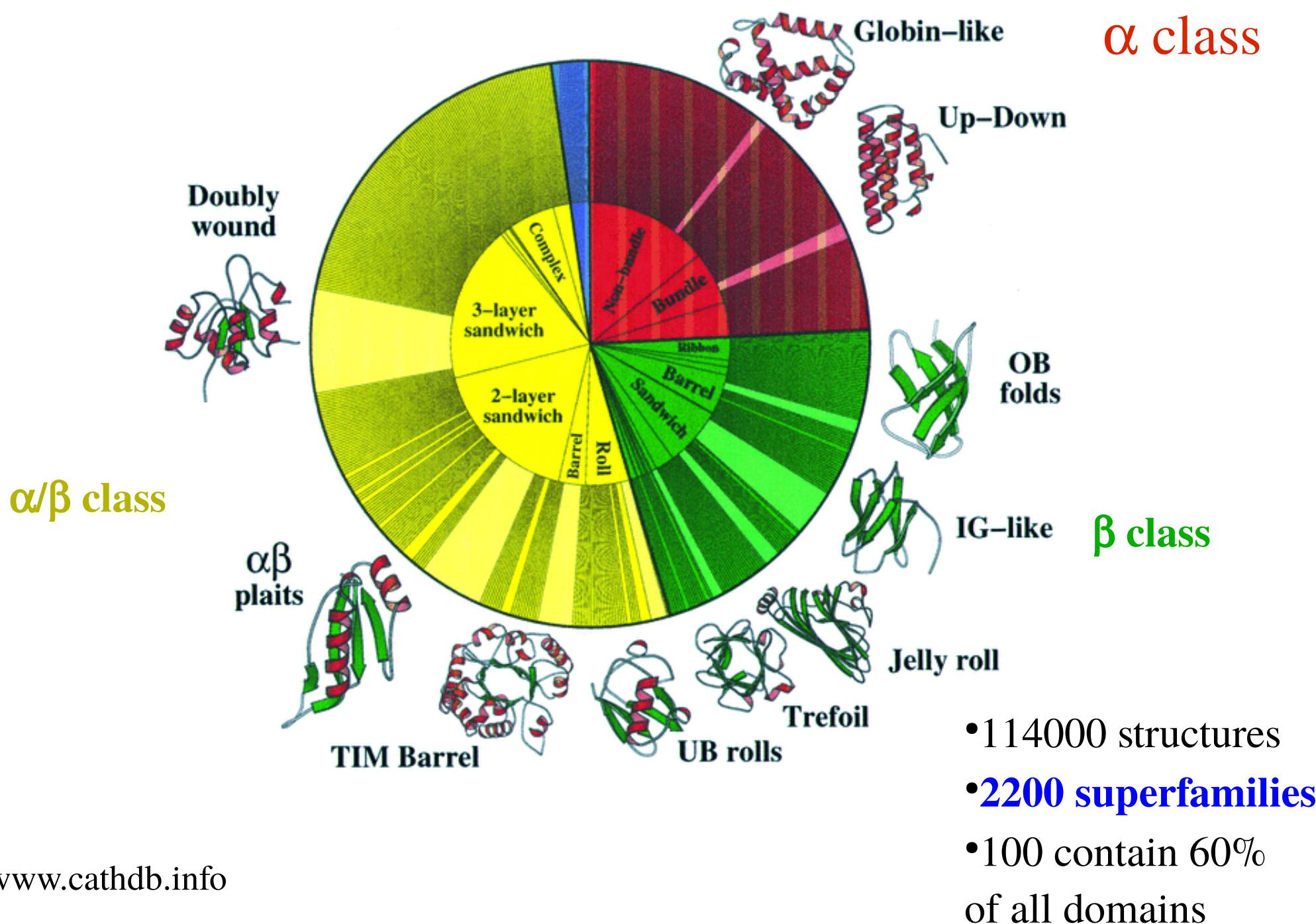
Gene duplication



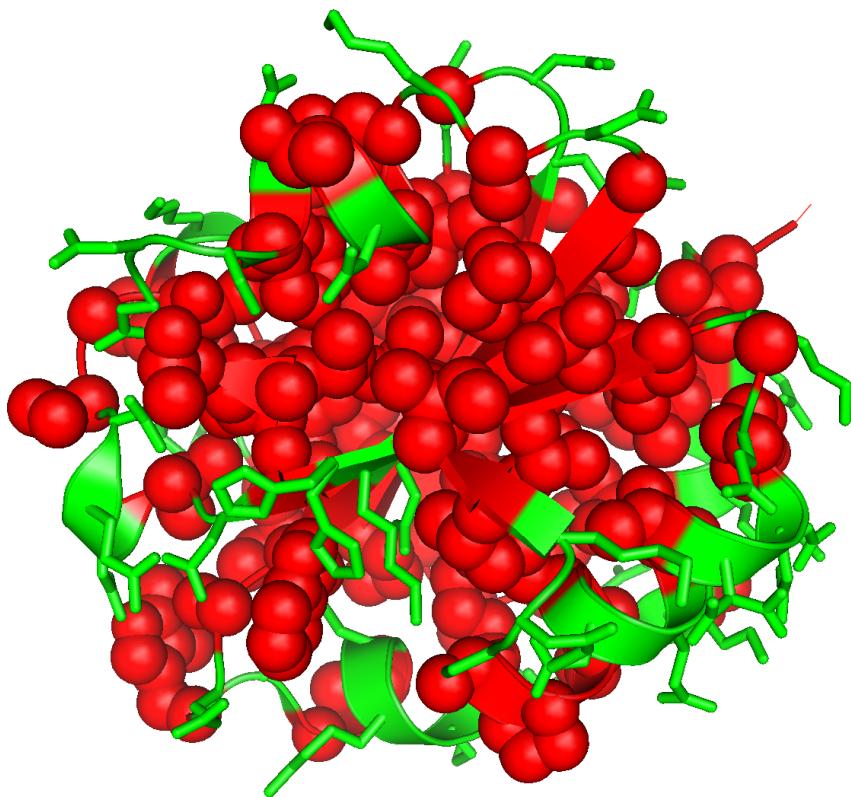
Divergence



Structural classification



Protein stability



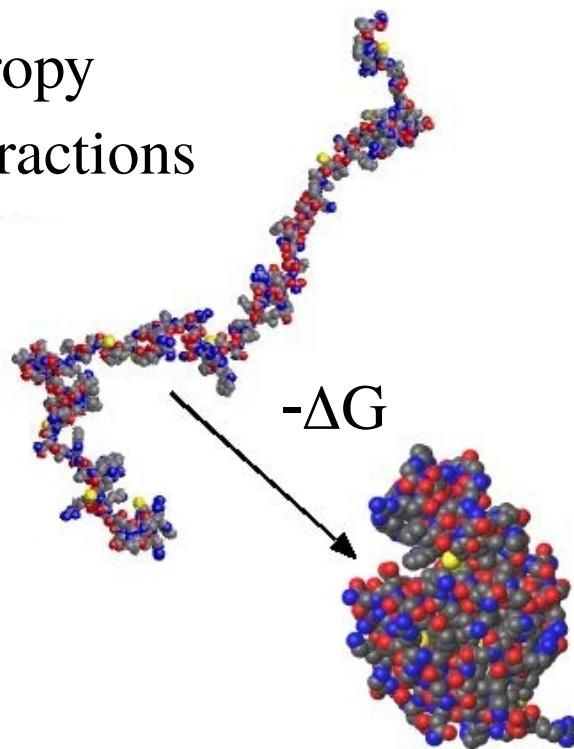
Proteins are marginally stable

$$\Delta G = 5-15 \text{ kcal/mol} = 10-30 \text{ kT}$$

~ 0.1 kT per amino acid

Large conformational entropy

Many protein-solvent interactions



Cancelling effects

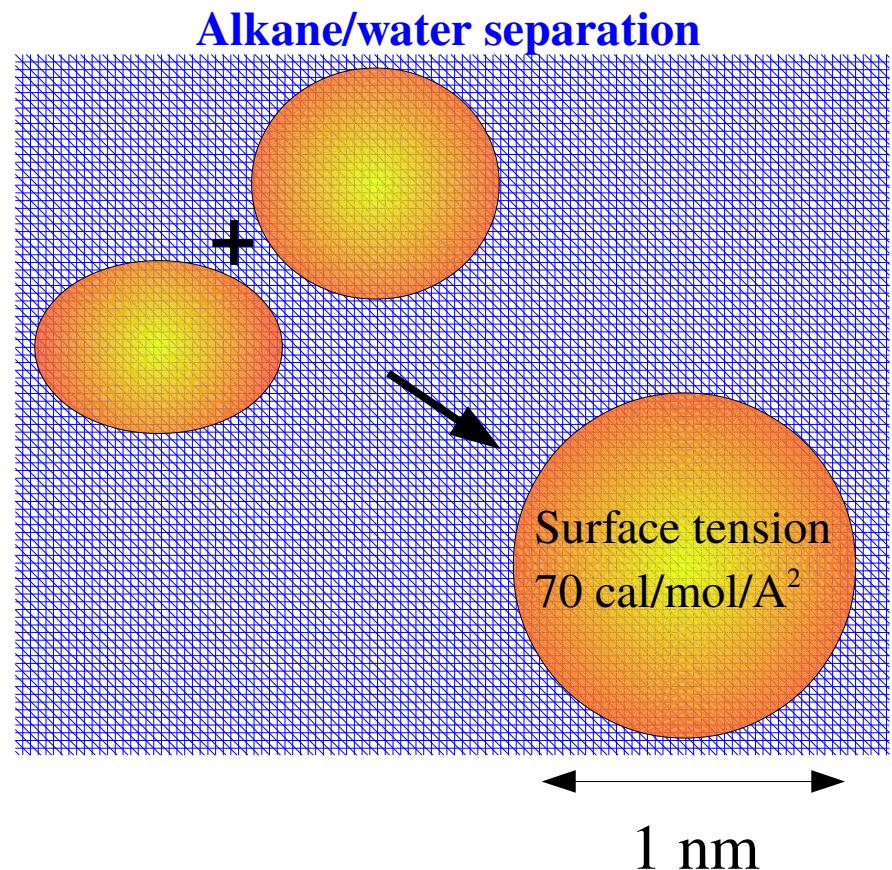
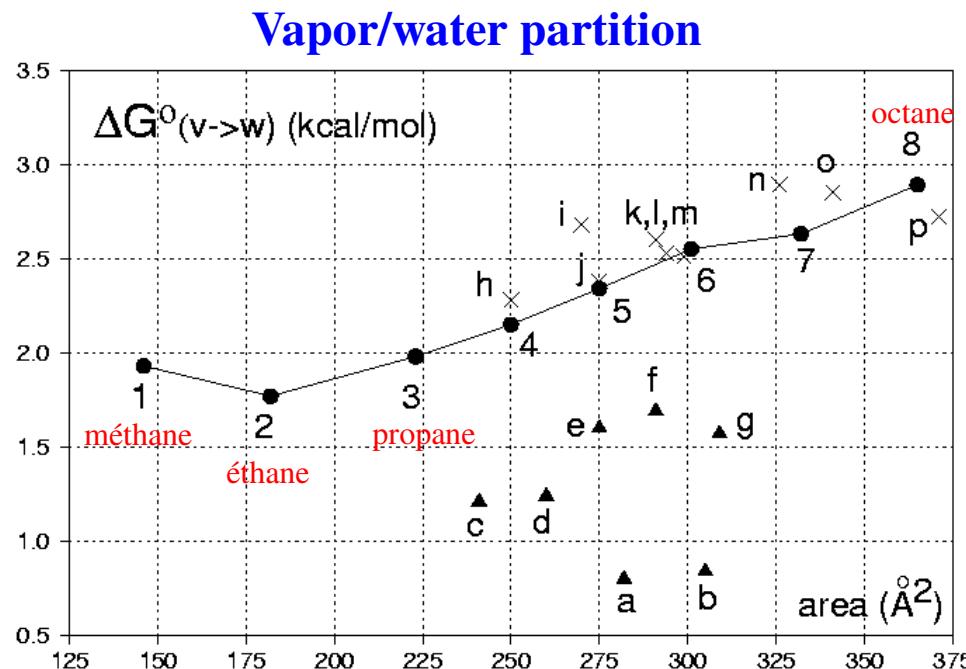
Low conformational entropy

Fewer protein-solvent interactions

Protein-protein interactions

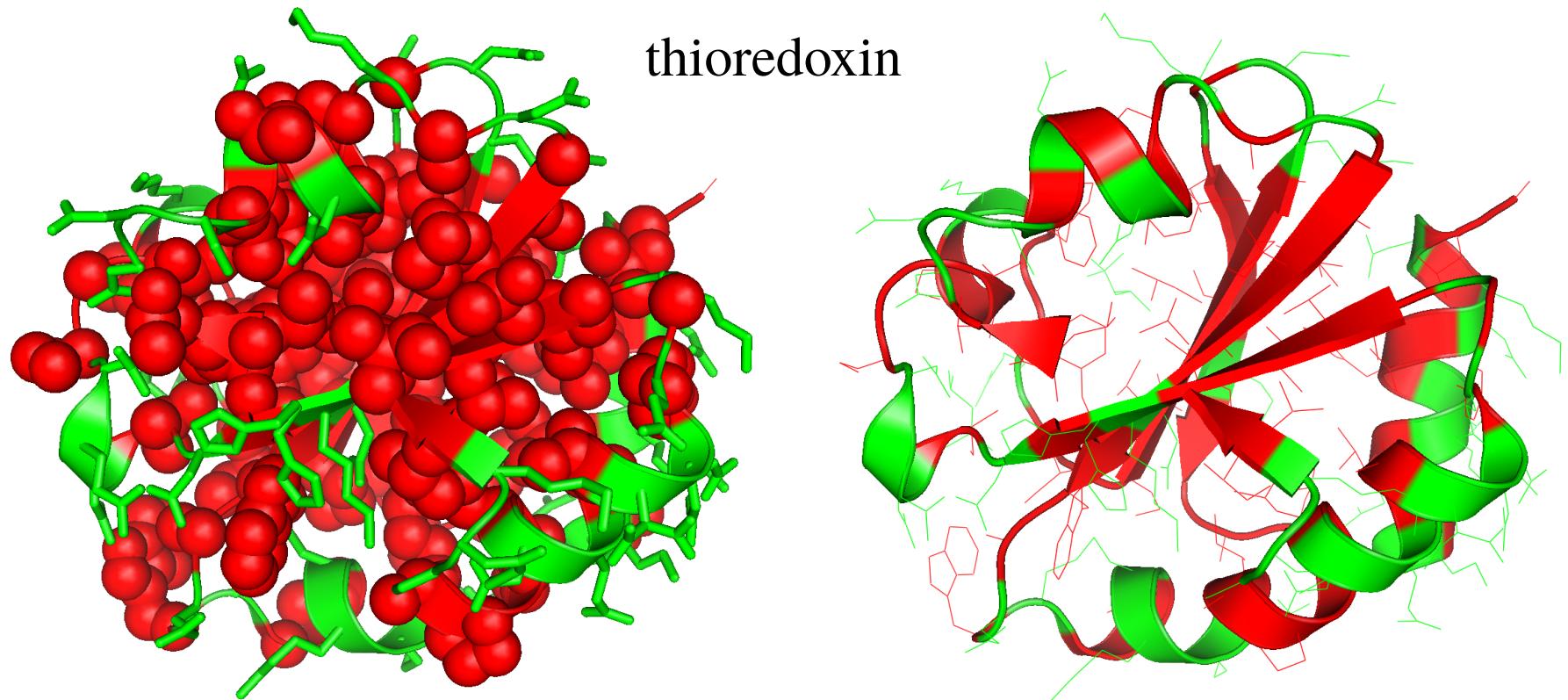
Folding is driven by the hydrophobic effect

Kauzmann, 1959



Saturated alkanes don't like water

Proteins tend to bury apolar groups, with very compact packing

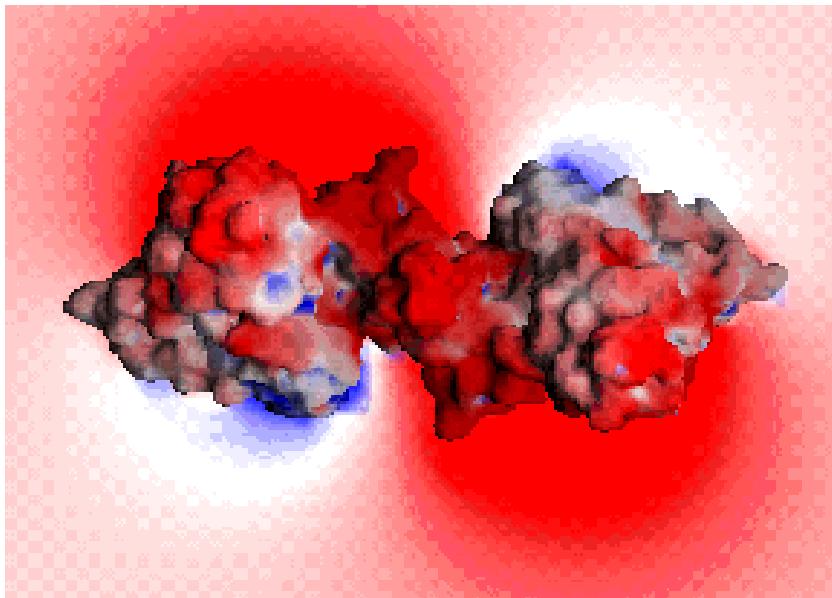


Fraction of occupied volume: ~ 0.74 (average over known structures)

Electrostatic interactions govern molecular recognition

25% of amino acids are ionic

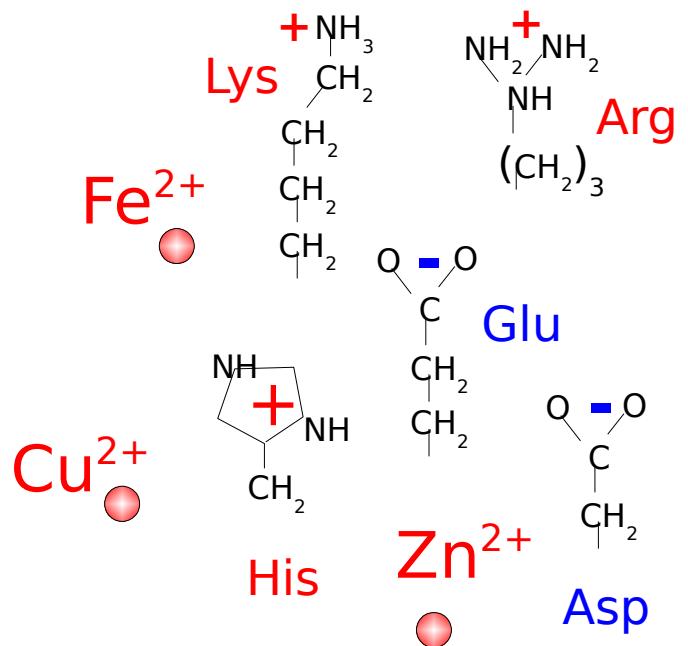
90% of salt bridges are at the surface



Acetylcholine esterase

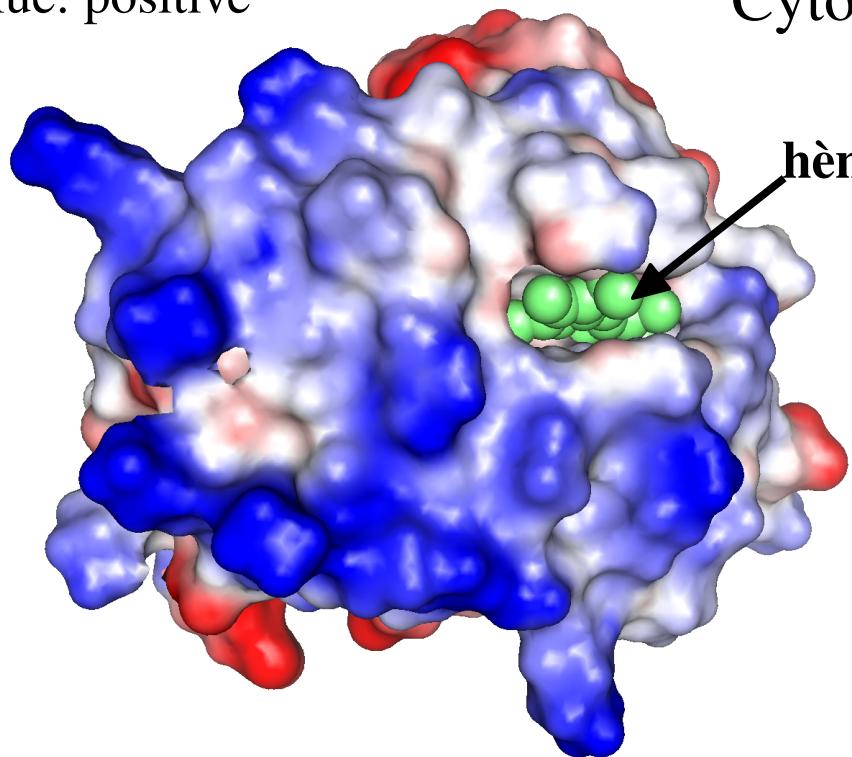
Red = negative electrostatic potential

Substrate acetylcholine is positive



Electrostatic interactions govern molecular recognition

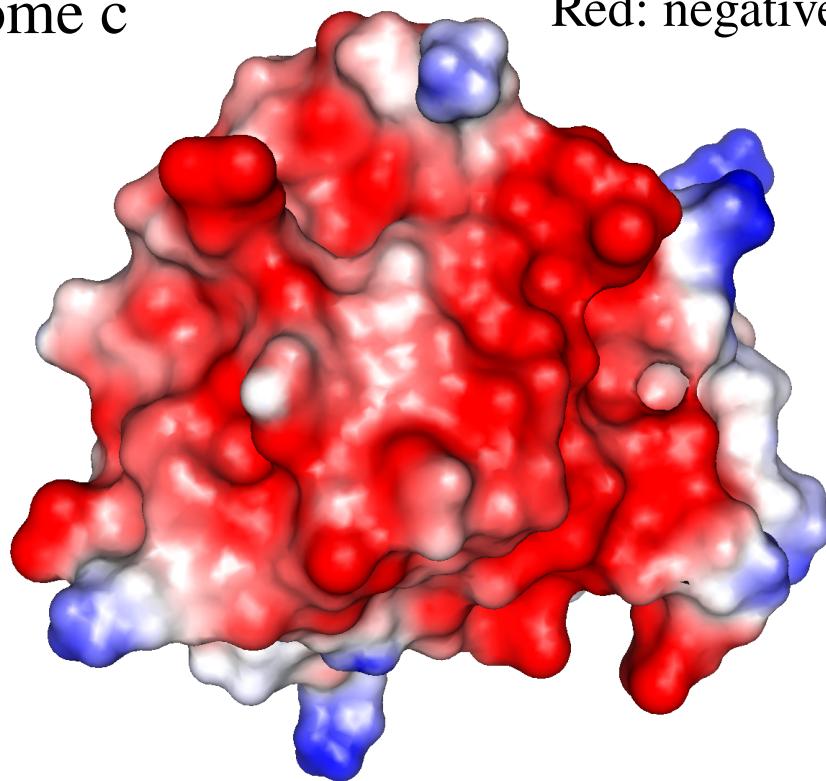
Blue: positive



Front

Cytochrome c

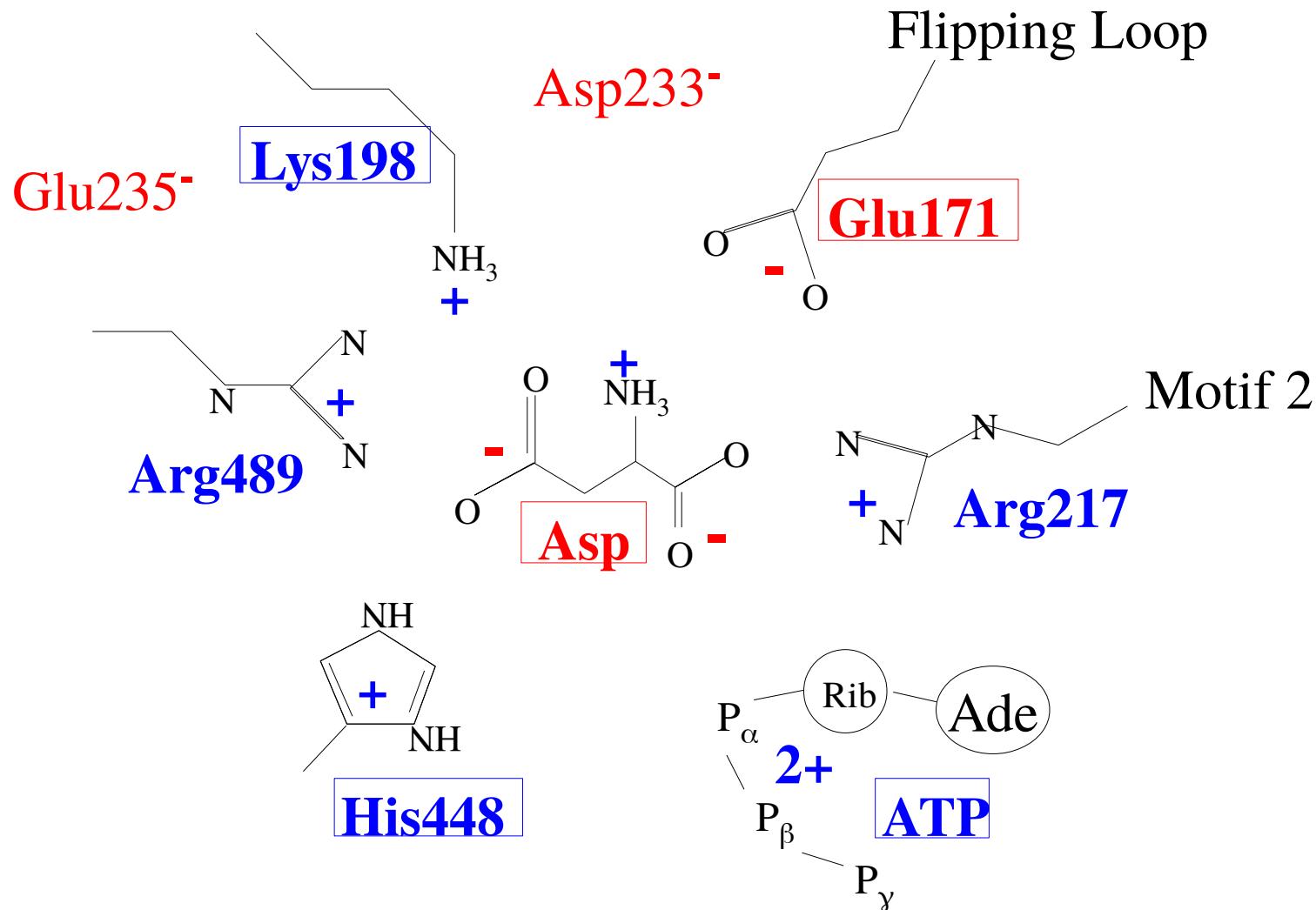
Red: negative



Back

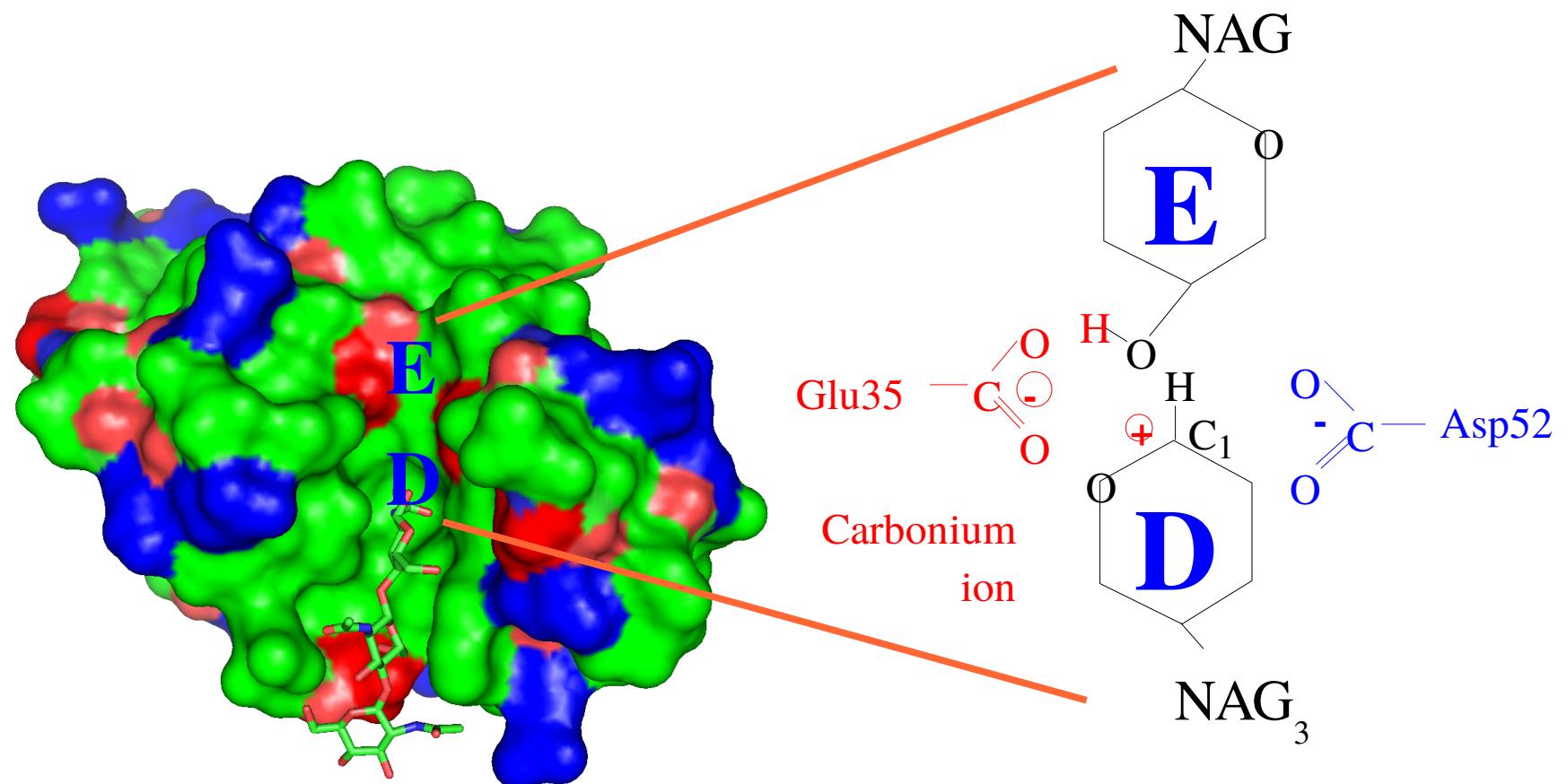
Cytochrome c interacts with the mitochondrial membrane, which is negative, and with negative regions on cytochrome bc1 and cytochrome oxidase.

Electrostatic interactions govern molecular recognition



Aspartyl-tRNA synthetase active site

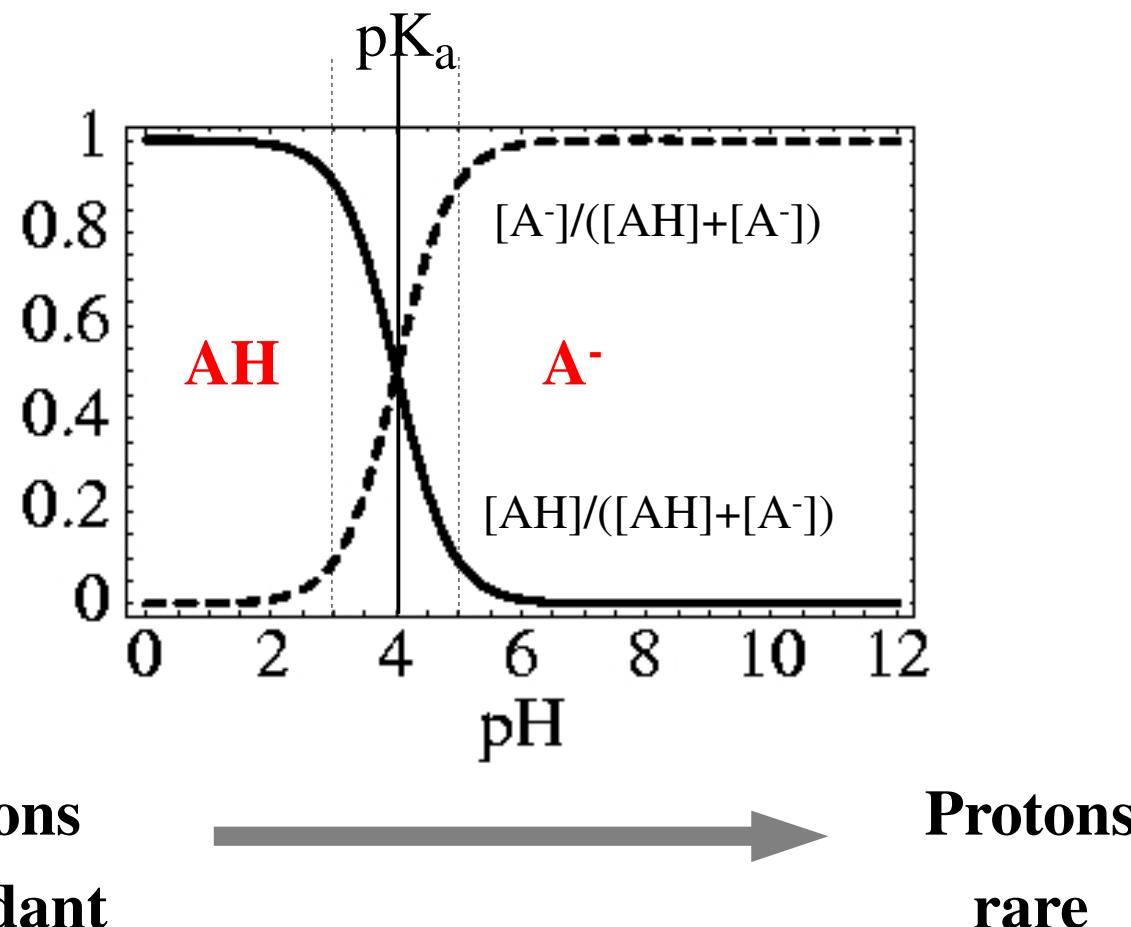
Electrostatic interactions govern molecular recognition



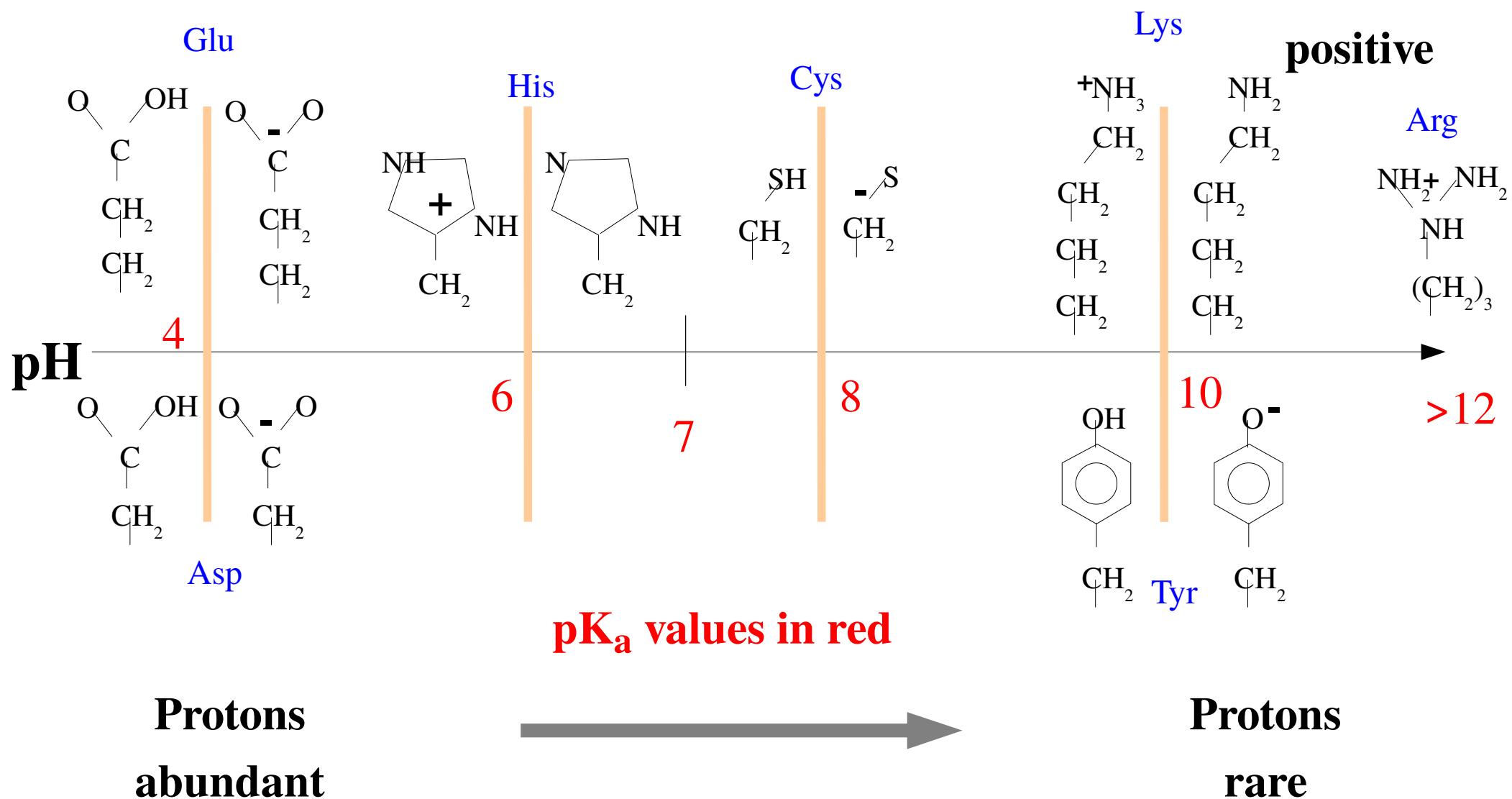
Lysozyme active site

Certain amino acids are ionized at neutral pH

Acid/base equilibrium:



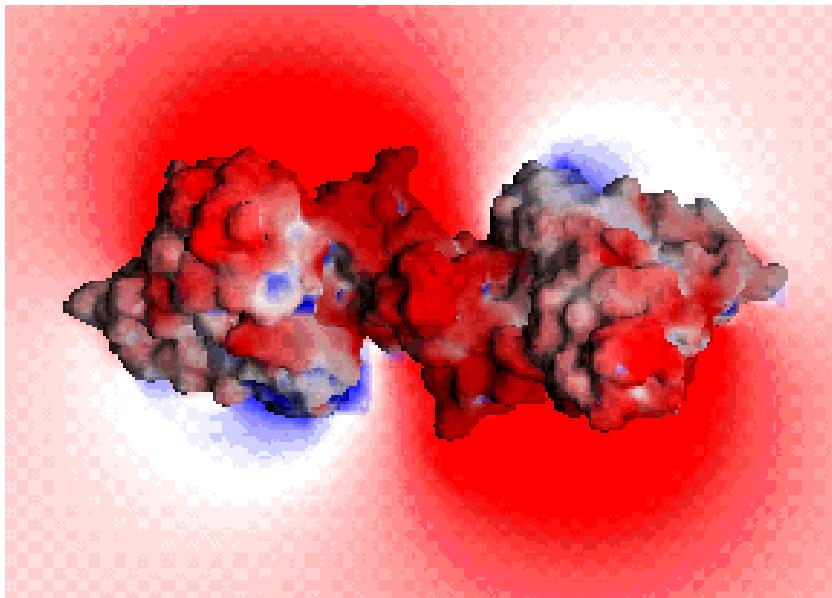
Certain amino acids are ionized at neutral pH



Electrostatic interactions govern molecular recognition

25% of amino acids are ionic

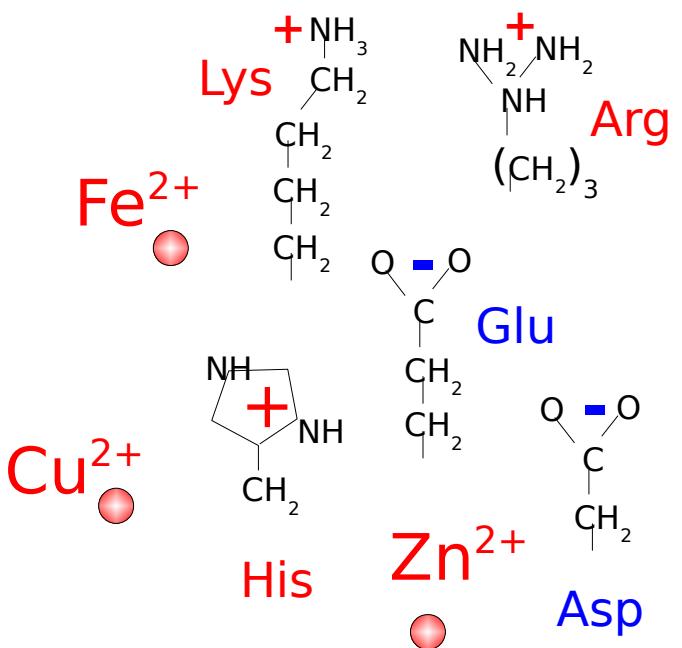
90% of salt bridges are at the surface



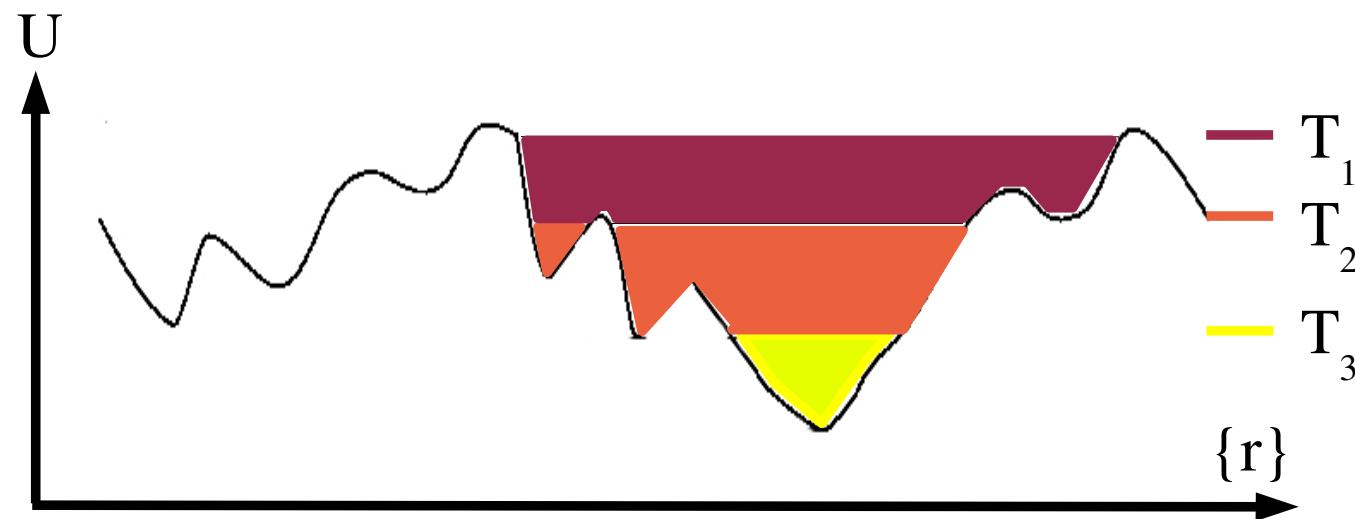
Acetylcholine esterase

Red = negative electrostatic potential

Substrate acetylcholine is positive

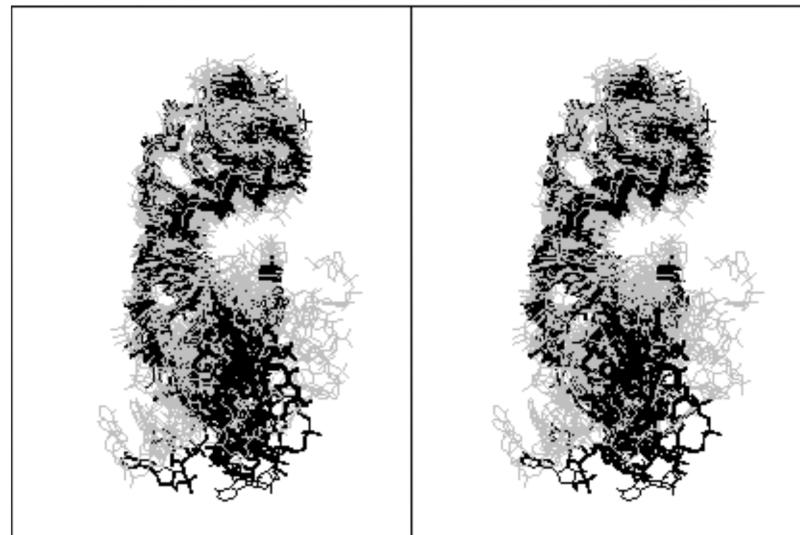


3D structure: molecular mechanics and molecular dynamics



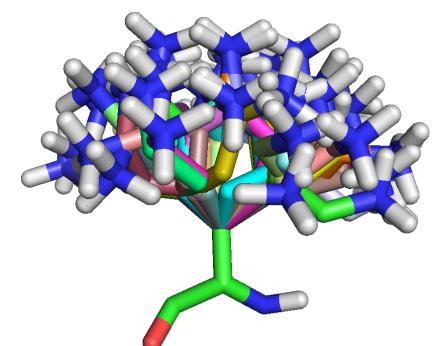
Biomolecular structures are flexible and mobile

U C
U G
C....G
U....A
C....G
G....C
 $^3\text{G}....\text{U}_{16}$
G....C
 $^1\text{G}....\text{C}_{18}$
A
C
C
 $\text{A}_{3\text{ter}}$

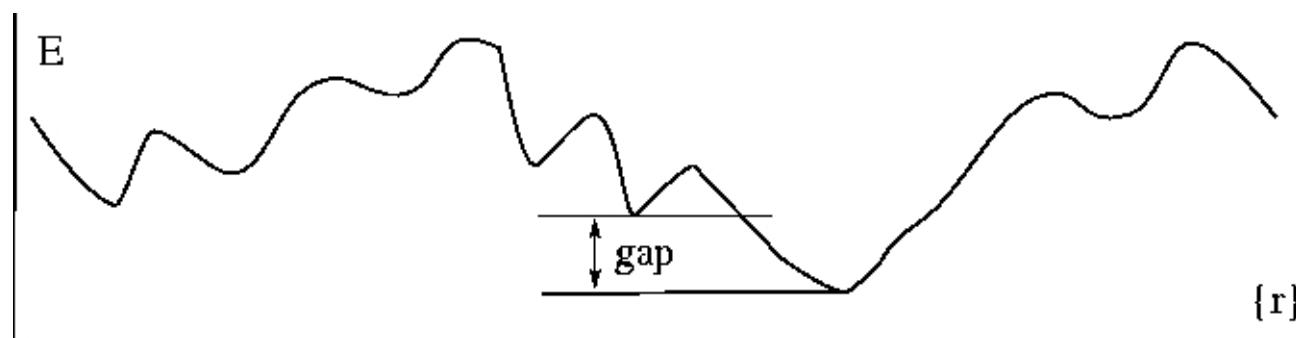
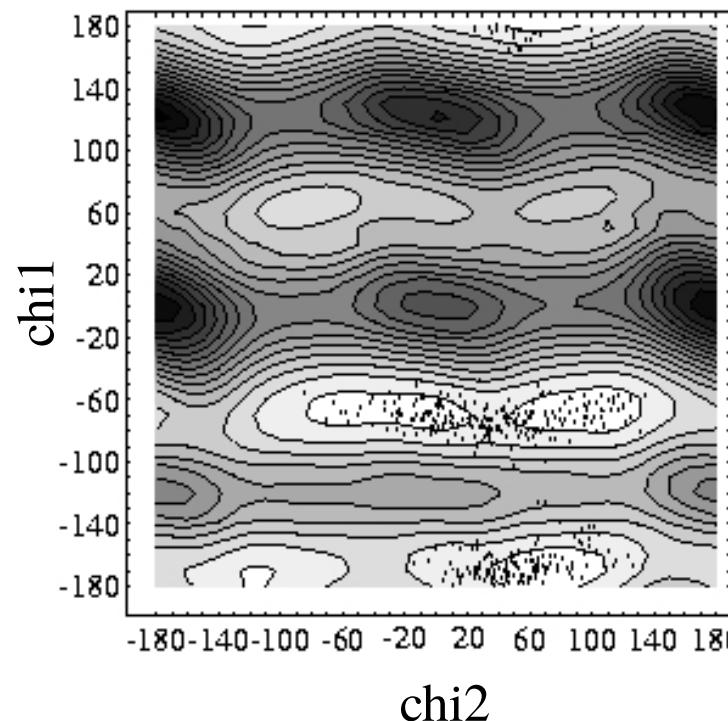
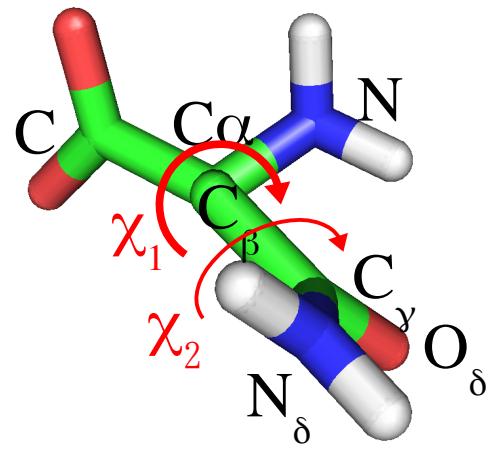


RNA stem-loop
(stereo)

Lysine

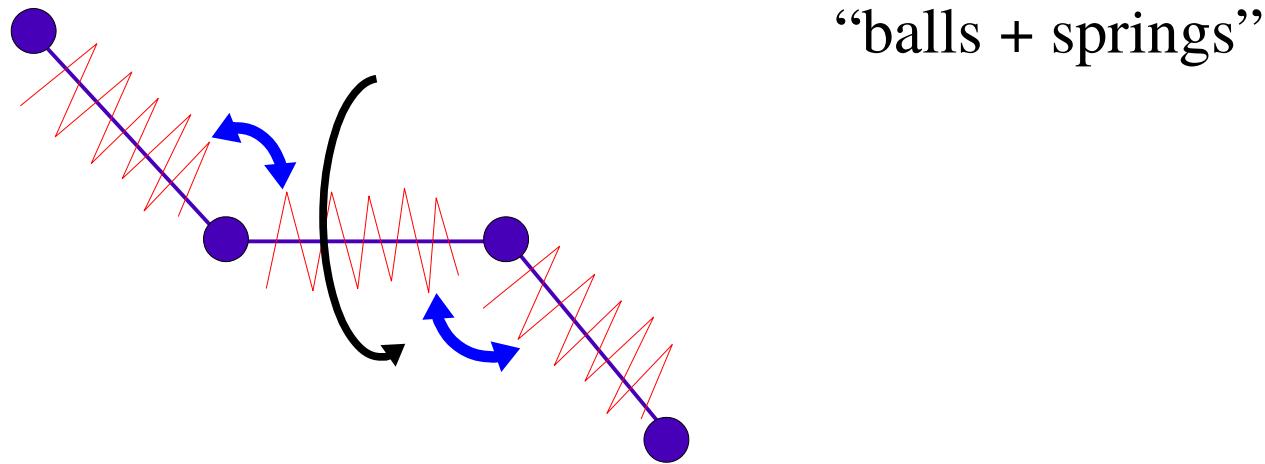


We characterize them by an energy surface



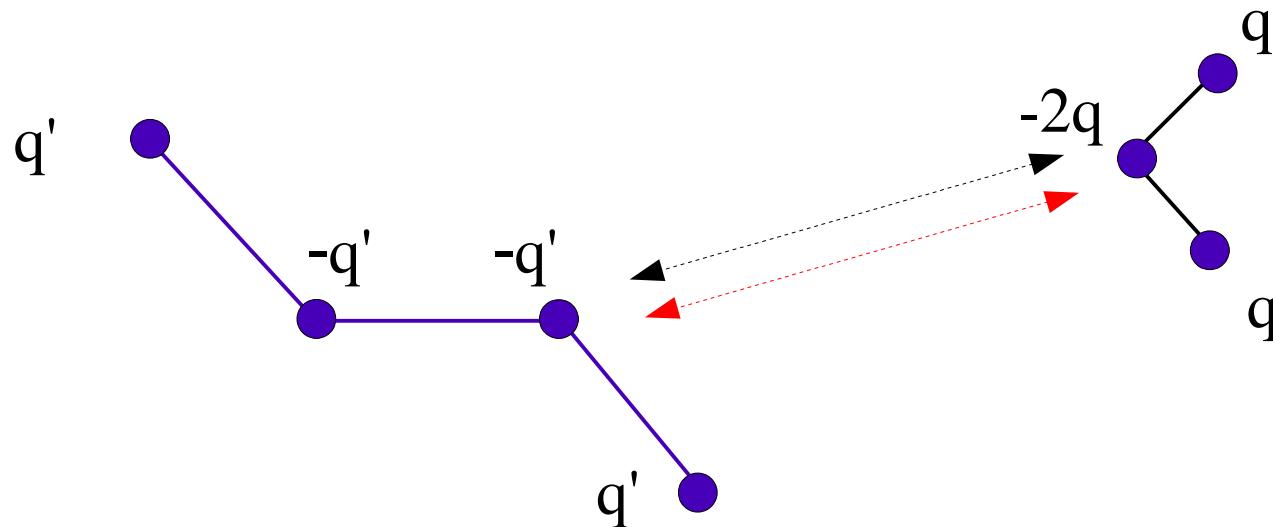
Structure prediction = find low energy structures

We introduce a semi-empirical energy function

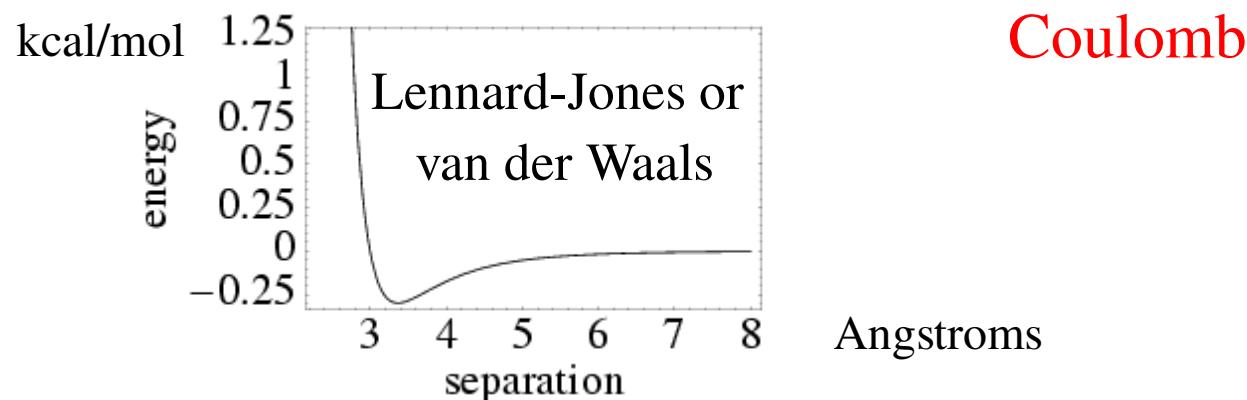


$$U = \sum_{\text{bonds}} k_b (b - b_0)^2 + \sum_{\text{angles}} k_a (a - a_0)^2 + \sum_{\text{torsions}} k_t [1 + \cos(nt - \tau)]$$

We introduce a semi-empirical energy function



$$U = \sum_{ij} [A_{ij} / r_{ij}^{12} - B_{ij} / r_{ij}^6] + \sum_{ij} q_i q_j / r_{ij}$$



Rôle of van der Waals interactions: Leu → Val mutations

11268 *Biochemistry, Vol. 32, No. 42, 1993*

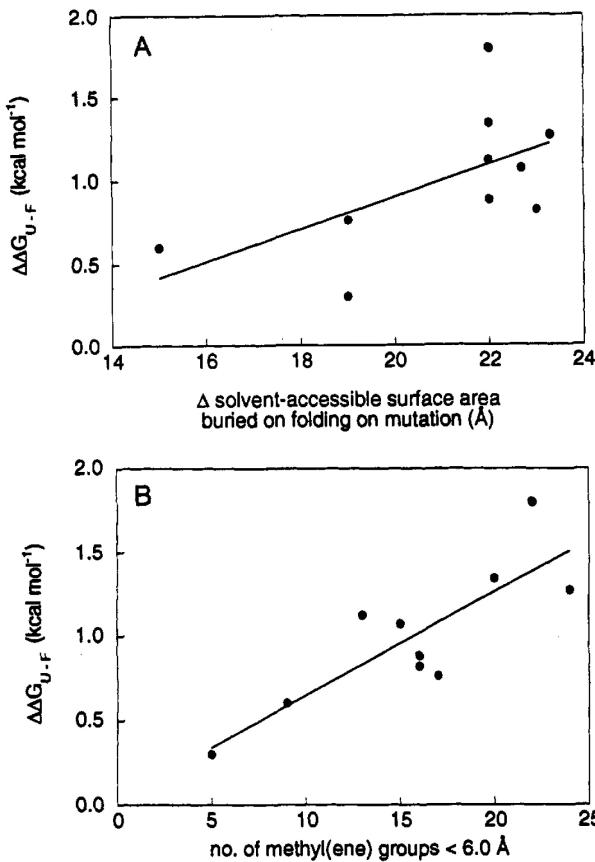
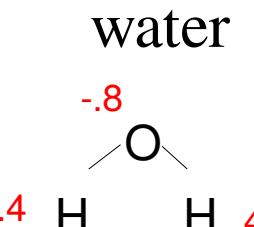
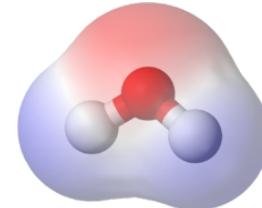
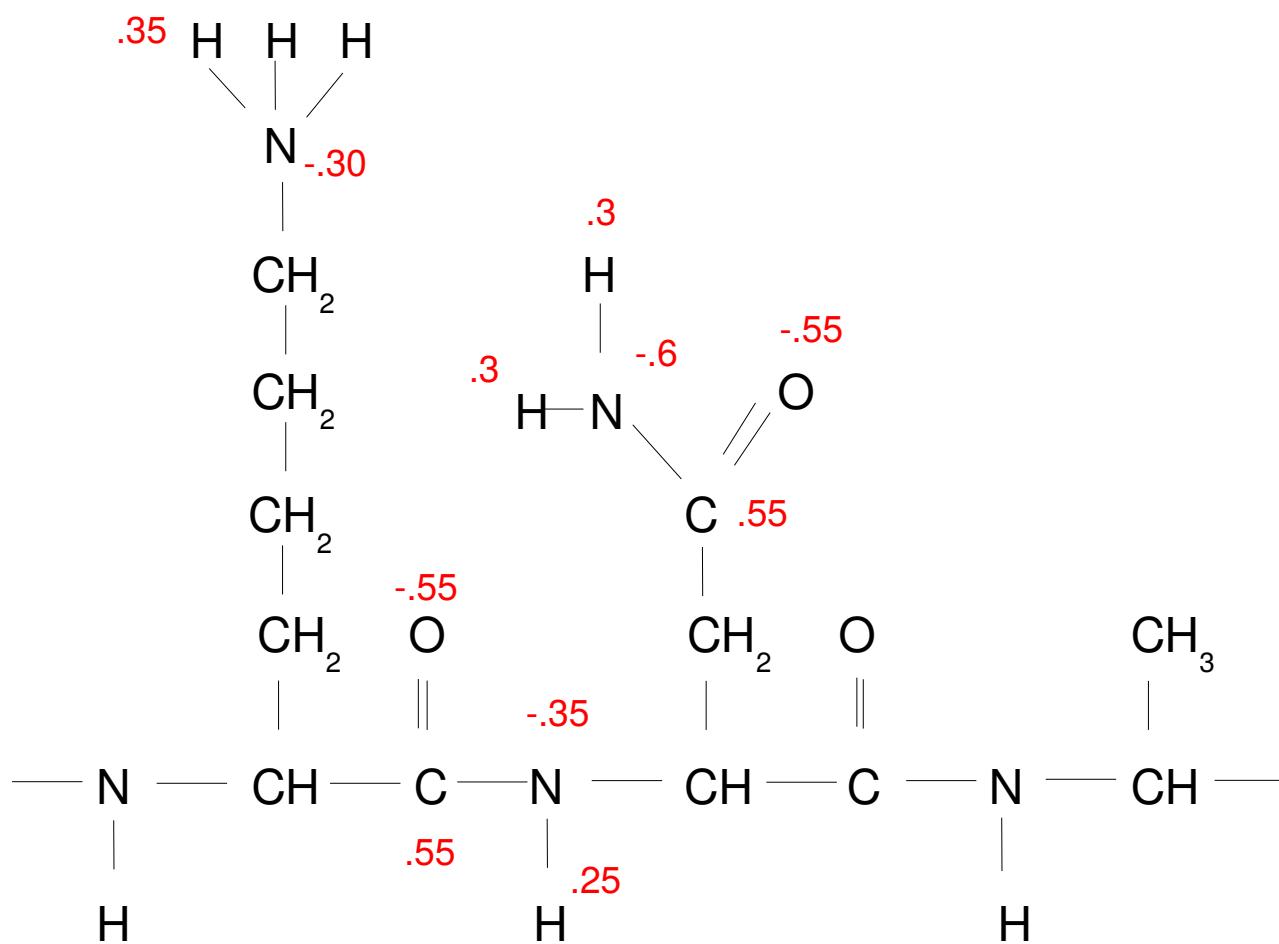


FIGURE 6: (A) Correlation between the difference in the solvent-accessible surface area that is buried on folding between wild-type and mutant side chains and the changes in the free energy of unfolding for Ile → Val mutations only for CI2 and barnase. The solid line shows the best fit of the data to a linear equation (correlation coefficient = 0.59). (B) Correlation between the number of side-chain methylene groups, in a radius of 6 Å of the group deleted from wild-type, and the changes in the free energy of unfolding for Ile → Val mutations only for CI2 and barnase. The solid line shows the best fit of combined CI2 and barnase data to a linear equation (correlation coefficient = 0.83).

Molecular mechanics energy

$$U = \sum_{\text{bonds}} k_b (b - b_0)^2 + \sum_{\text{angles}} k_a (a - a_0)^2 + \sum_{\text{torsions}} k_t [1 + \cos(nt - \tau)]$$

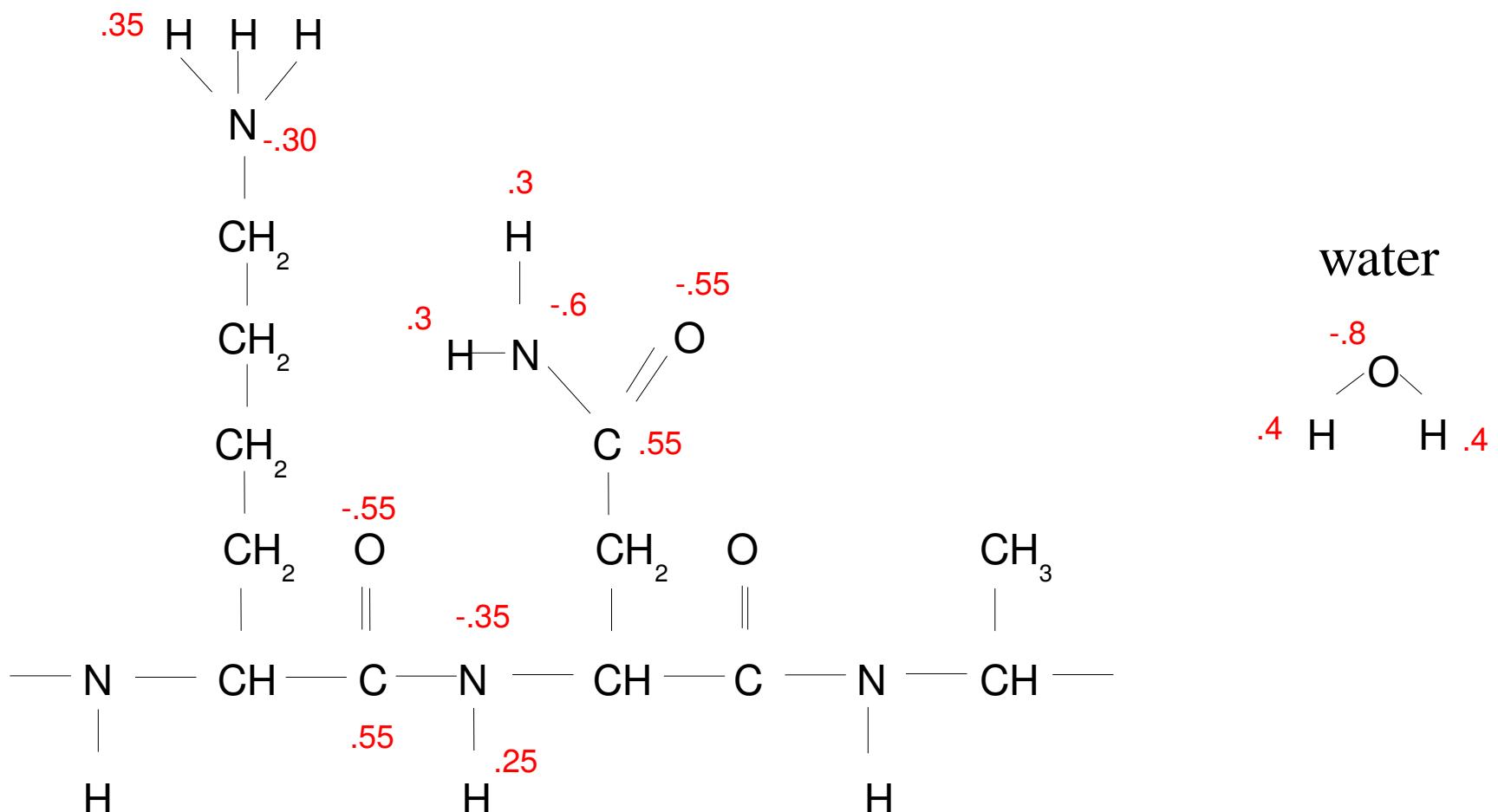
$$+ \sum_{ij} [A_{ij} / r_{ij}^{12} - B_{ij} / r_{ij}^6 + q_i q_j / r_{ij}]$$



Molecular mechanics energy

$$U = \sum_{\text{bonds}} k_b (b - b_0)^2 + \sum_{\text{angles}} k_a (a - a_0)^2 + \sum_{\text{torsions}} k_t [1 + \cos(nt - \tau)]$$

$$+ \sum_{ij} [A_{ij}/r_{ij}^{12} - B_{ij}/r_{ij}^6 + q_i q_j / r_{ij}]$$



Notion of parameter *transferability*

Fichier de “topologie”, extrait

```
RESidue ALA
GROUP
ATOM N      TYPE=NH1      CHARge=-0.35      END
ATOM H      TYPE=H        CHARge= 0.25      END
ATOM CA     TYPE=CH1E     CHARge= 0.10      END
GROUP
ATOM CB     TYPE=CH3E     CHARge= 0.00      END
GROUP
ATOM C       TYPE=C        CHARge= 0.55      END
ATOM O       TYPE=O        CHARge=-0.55      END

BOND N      CA
BOND CA     C
BOND C      O
BOND N      H
BOND CA     CB

IMPRoper   CA      N      C      CB      !tetrahedral CA

END {ALA}
```

{Fichier toph19.pro}

Fichier des paramètres d'énergie, extrait

bond C C 450.0 1.38! B. R. GELIN THESIS AMIDE AND DIPEPTIDES
bond C CH1E 405.0 1.52! EXCEPT WHERE NOTED. CH1E,CH2E,CH3E, AND CT
bond C CH2E 405.0 1.52! ALL TREATED THE SAME. UREY BRADLEY TERMS ADDED
bond C CH3E 405.0 1.52
bond C CR1E 450.0 1.38

:

angle C C C 70.0 106.5! FROM B. R. GELIN THESIS WITH HARMONIC
angle C C CH2E 65.0 126.5! PART OF F TERMS INCORPORATED. ATOMS
angle C C CH3E 65.0 126.5! WITH EXTENDED H COMPENSATED FOR LACK
angle C C CR1E 70.0 122.5! OF H ANGLES.

:

NONBonded H 0.0498 1.4254 0.0498 1.4254 0. 1.
NONBonded HA 0.0450 2.6157 0.0450 2.6157 0. 1 ! charged group.

NONBonded C 0.1200 3.7418 0.1000 3.3854 7.116 1. ! carbonyl carbon
NONBonded CH1E 0.0486 4.2140 0.1000 3.3854 9.944 1. !
NONBonded CH2E 0.1142 3.9823 0.1000 3.3854 17.626 1. ! extended carbons

Molecular mechanics energy

- No explicit electrons
- Partial charges on atoms
- No chemical reactions
- big systems: >100.000 atoms
- Conformational energies
- Explicit or implicit solvent representation
- explicit force calculation
- Molecular dynamics

Parametrization: 1970-present

CHARMM M. Karplus, A. Mackerell and coll. (Harvard)

AMBER P. Kollman, D. Case and coll. (UCSF, Scripps)

OPLS W. Jorgensen and coll. (Yale)

GROMOS H. Berendsen, W. van Gunsteren and coll. (Groningen, Zürich)

Plus récent: AMOEBA Jay Ponder (St Louis)

Atomic charges : crystals of small molecules
simple liquids
quantum chemistry calculations

Lennard-Jones: crystals of small molecules
simple liquids

bonds, angles, torsions: quantum chemistry calculations
small molecule spectroscopy

Software

CHARMM19 (M Karplus et al, Harvard)

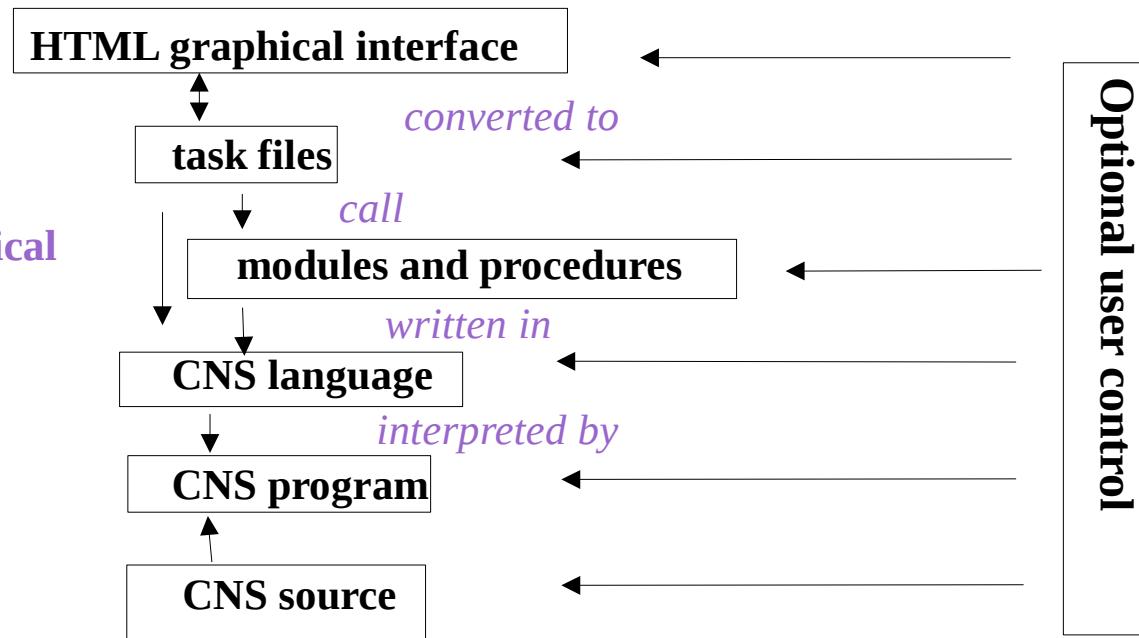
XPLOR (A Brunger, Yale)

CNS ('Crystallography and NMR System') ~220,000 lines

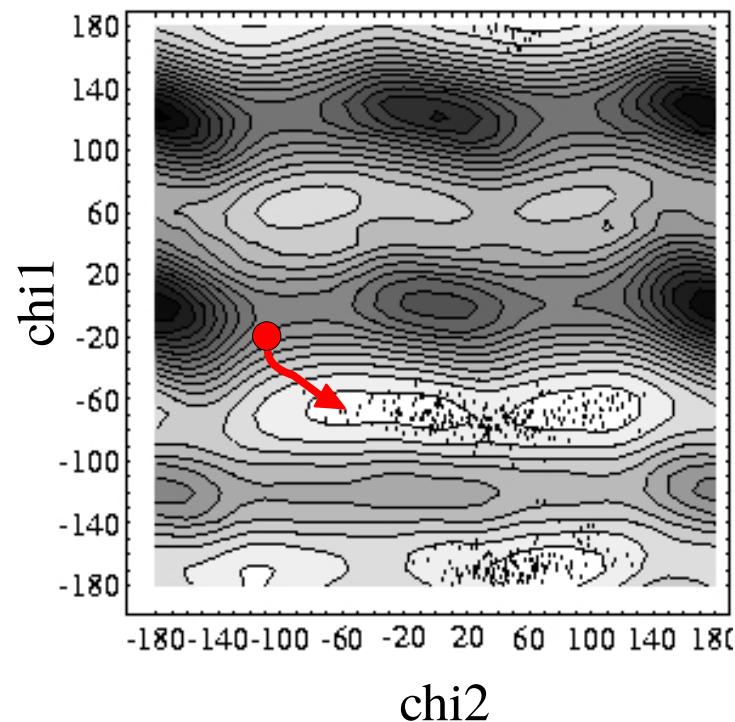
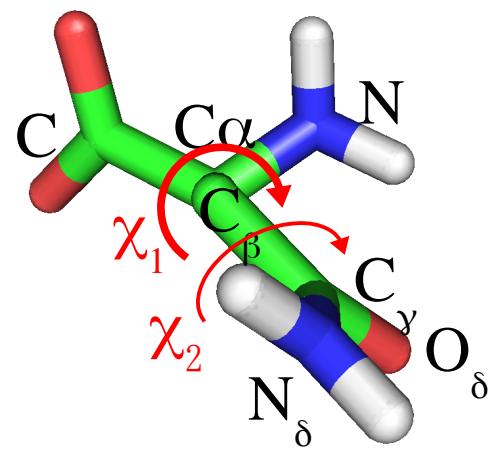
NIH-XPLOR
(M Clore,
C Schweiters,
J Kuszewski)

A Brunger, P Adams, G Clore, W Delano, P Gros, R Grosse-Kunstleve, J Jiang, J Kuszewski, M Nilges, N Pannu, R Read, L Rice, T Simonson, G Warren (1998) Acta Cryst D54, 905.

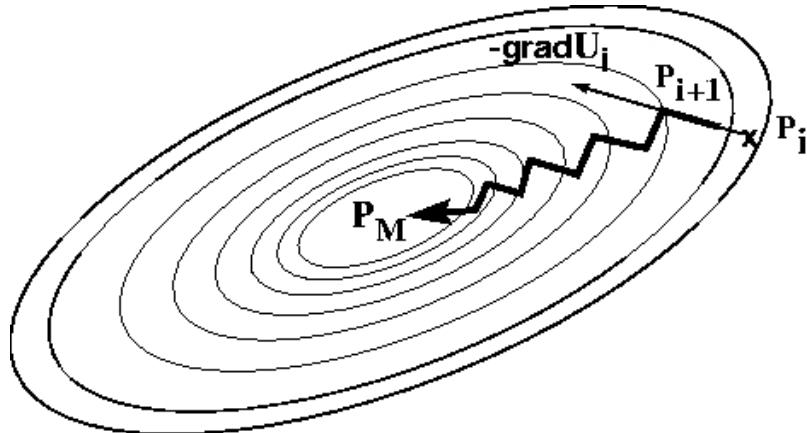
Hierarchical
structure:



Finding low energy structures: energy minimization

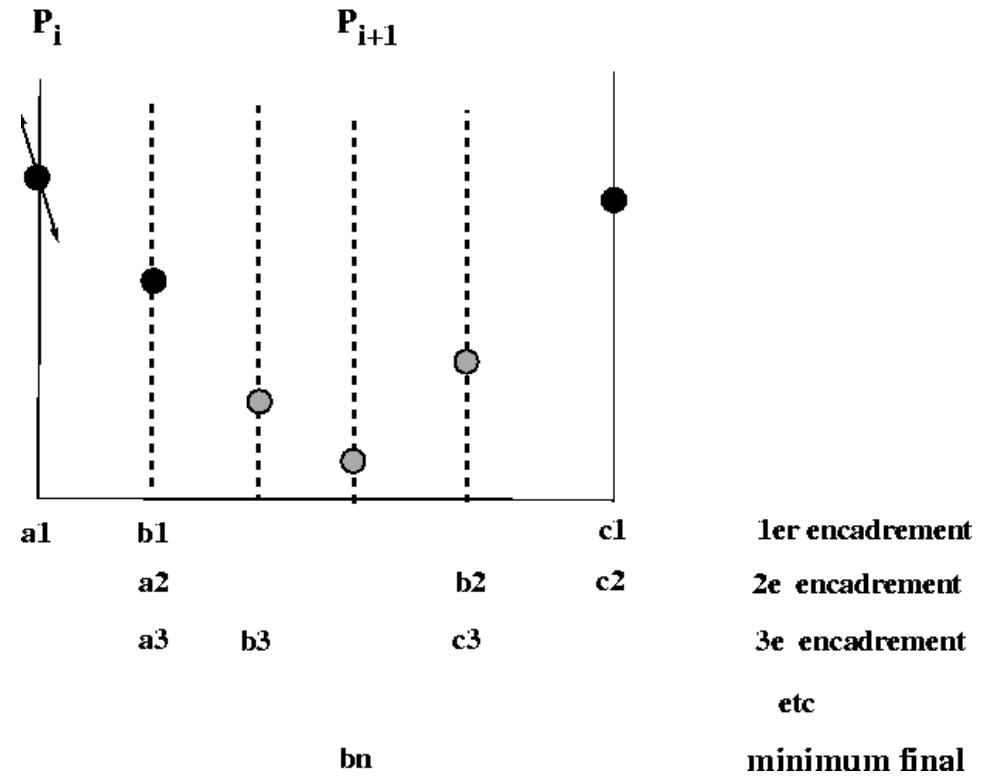


Steepest descent minimization



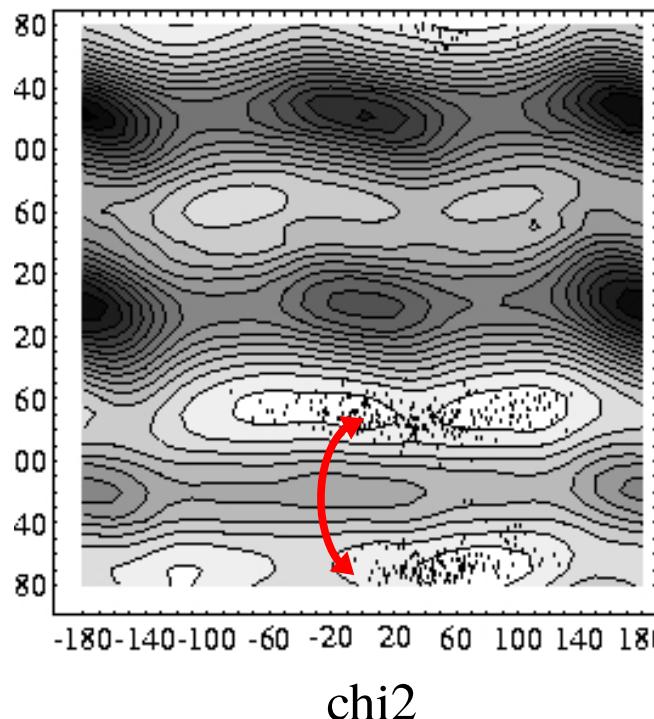
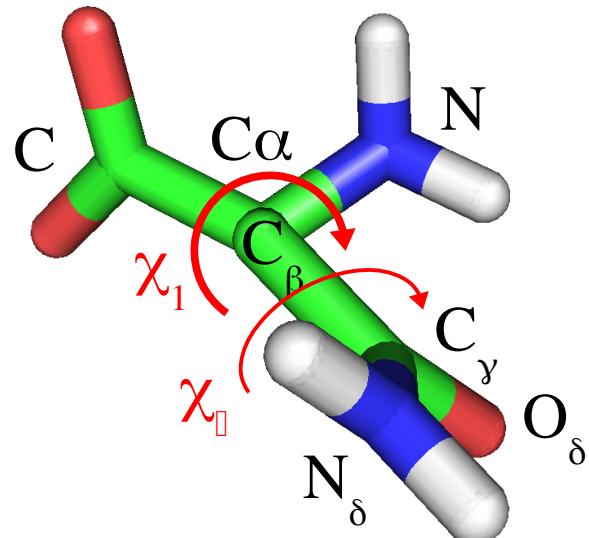
(1) **Method:** starting from P_i , move along the energy gradient: $-\text{grad } U(P_i)$ until a minimum, P_{i+1} . Repeat.

(2) **Finding P_{i+1} :**
successive interpolations



Cf Numerical Recipes; Press et al.

Exploring conformations with molecular dynamics



Can cross
energy barriers
(not too high)

Solve equations of motion numerically:

$$m_i \ddot{\gamma}_i = F_i = -\nabla U_i$$

3 equations
per atom

Verlet algorithm

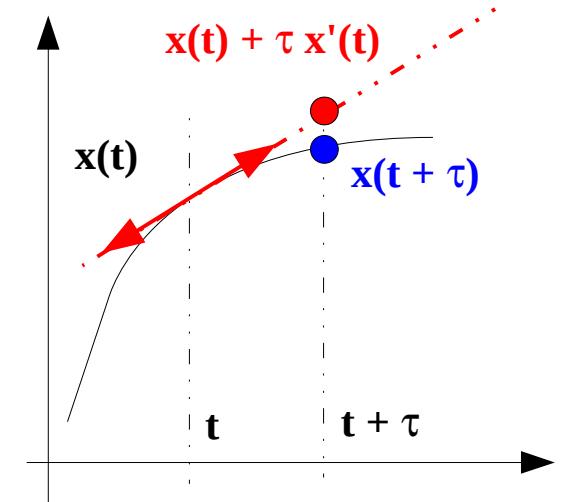
Consider each coordinate of each atom i as a function of time t

The change between two instants t and $t+\tau$ can be approximated:

$$x_i(t+\tau) \approx x_i(t) + \tau x'_i(t)$$

$$y_i(t+\tau) \approx \dots$$

$$z_i(t+\tau) \approx \dots$$



Verlet algorithm

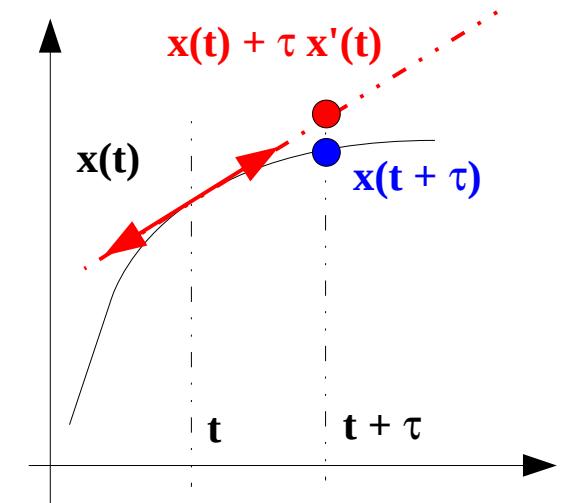
Consider each coordinate of each atom i as a function of time t

The change between two instants t and $t+\tau$ can be approximated:

$$x_i(t+\tau) \approx x_i(t) + \tau x'_i(t) + (\tau^2/2) x''_i(t)$$

$$y_i(t+\tau) \approx \dots$$

$$z_i(t+\tau) \approx \dots$$



Verlet algorithm

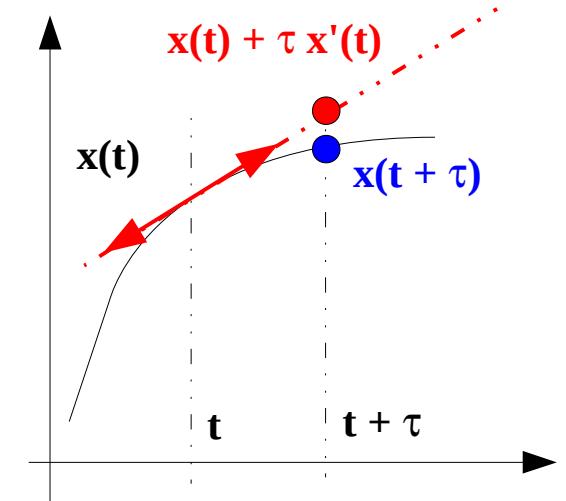
On soustrait les valeurs pour $t+\tau$ et $t-\tau$:

$$x_i(t+\tau) \cong x_i(t) + \tau x'_i(t) + (\tau^2/2) x''_i(t)$$

$$x_i(t-\tau) \cong x_i(t) - \tau x'_i(t) + (\tau^2/2) x''_i(t)$$

$$x_i(t+\tau) + x_i(t-\tau) \cong 2 x_i(t) + \tau^2 x''_i(t)$$

$$x_i(t+\tau) \cong 2 x_i(t) - x_i(t-\tau) + (\tau^2 / m_i) F_{xi}(t)$$



$$\leftarrow x''_i(t) = F_{xi}(t) / m$$

- τ is the **integration step**. It should be smaller than the system's "fastest natural vibrations: $\sim 10^{-15}$ seconds
- At each step we perform one force calculation

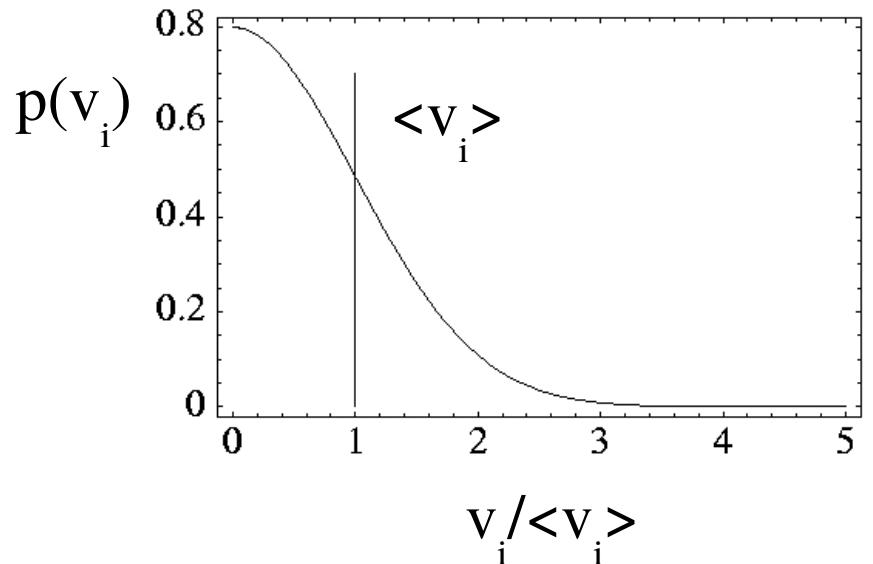
Link between temperature and atomic velocities

$$p(v_i) = A \exp(-m_i v_i^2 / 2kT)$$

T temperature

k Boltzmann constant

Maxwell-Boltzmann distribution



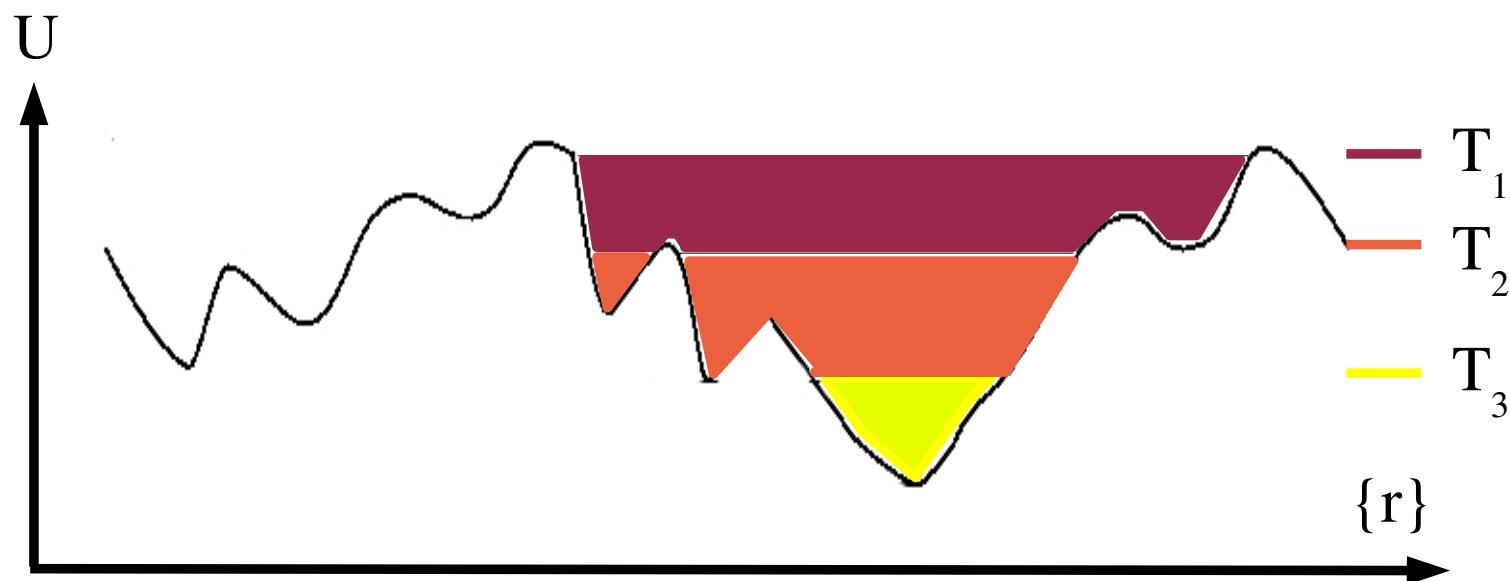
- Assign random initial velocities
- Rescale them from time to time

Recap: to simulate the dynamics of a biomolecule

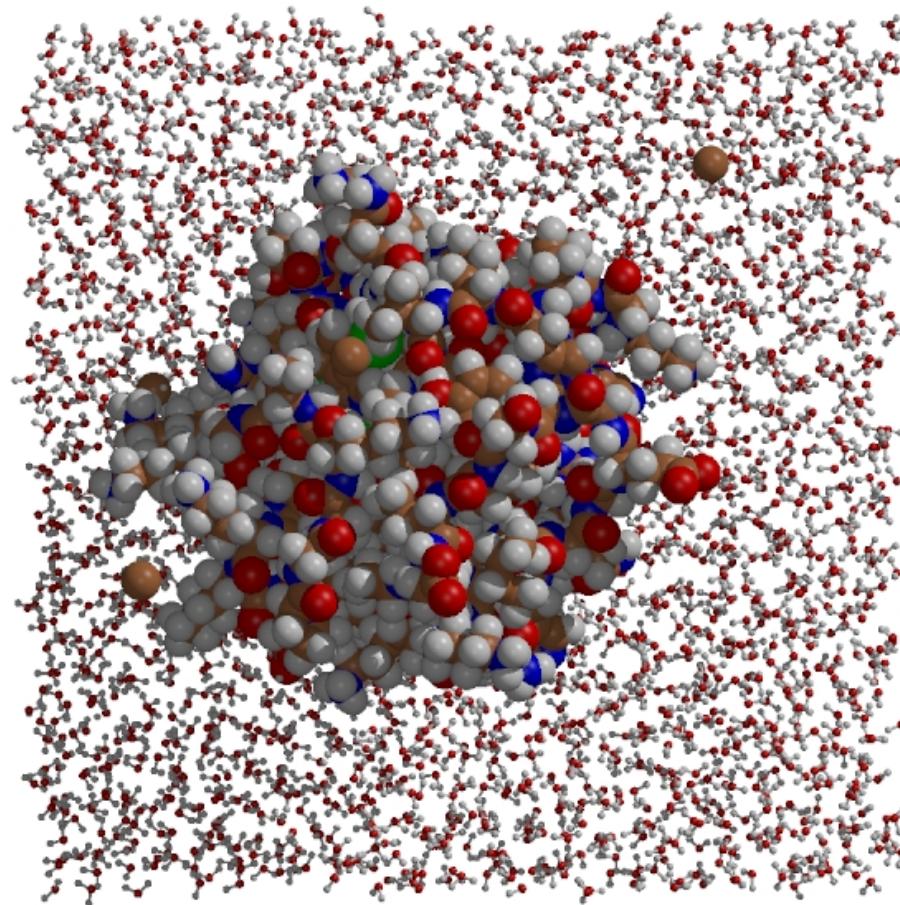
- “Molecular mechanics” “ball and stick” energy function
- Adjustable parameters (force constants, charges, etc) from small molecules, considered “transferable”
- “Parabolic” approximation of $x_i(t)$
- Recursive Verlet algorithm to solve equations of motion
- (Simple) relation between atomic velocities and temperature
- Strategy to model solvent....

Simulated annealing: Optimization on a rough energy surface

- Explore with molecular dynamics
- High temperature; lower gradually
- Refine structures from NMR or crystallography

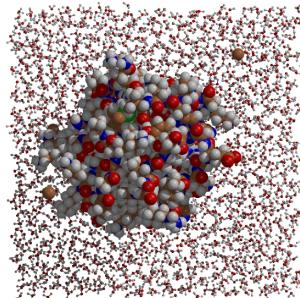


Explicit solvent

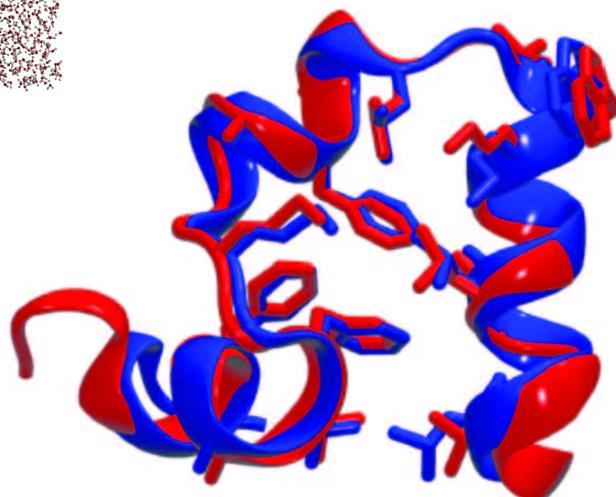


- Water box
- Periodic boundary conditions

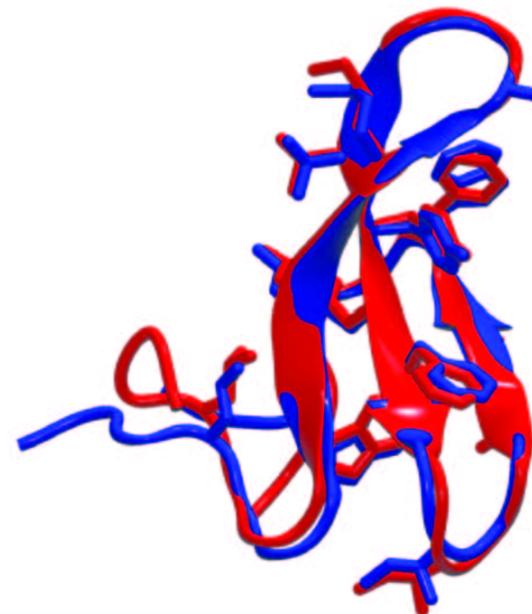
In silico protein folding



vilin



FiP35



Bleu: crystallography
Red: simulations

Domaine WW
35 acides aminés
Se replie
expérimentalement en
14 µs
Par simulation en 10 µs
4000 molécules d'eau

Shaw et al (2010)
Science, 330:341

Nobel chimie 2013
Martin Karplus



Structural Bioinformatics; editors PE Bourne, H Weissig; Wiley, 2003

Computational Biochemistry & Biophysics, editors Becker, Mackerell, Roux, Watanabe; Marcel Dekker, 2001

Bioinformatics: from genomes to drugs; editor T Lengauer; Wiley, 2002

Molecular modelling: principles and applications; A Leach; Prentice Hall, 2001

Introduction to Genomics; A Lesk, Oxford University Press, 2007

Bioinformatique: génomique et post-génomique; Dardel & Képès, Ecole Polytechnique, 2002

X-PLOR version 3.8, A System for X-ray crystallography and NMR; Brunger, Yale University Press, 1992

Simonson (2005) Médecine Sciences 21:609-612; Peut-on prédire la structure des protéines?

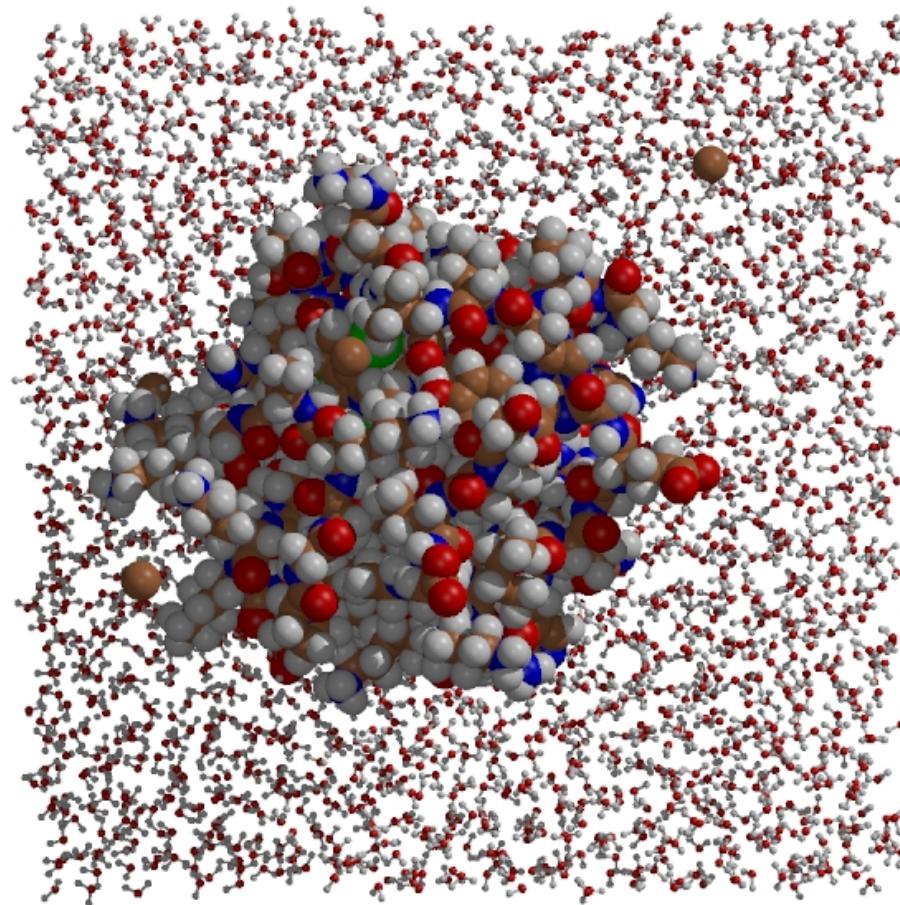
CNS; Brunger, Adams, Clore, Delano, Gros, Grosse-Kunstleve, Jiang, Kuszewski, Nilges, Pannu, Read, Rice, Simonson, Warren (1998) Acta Cryst D54, 905.

Brunger, Adams, DeLano, Gros, Grosse-Kunstleve, Jiang, Pannu, Read, Rice, Simonson (2001)

The structure determination language of the Crystallography and NMR System.

International Tables of Crystallography, Volume F. Editors: M.G. Rossmann and E. Arnold; Dordrecht: Kluwer Academic Publishers, the Netherlands.

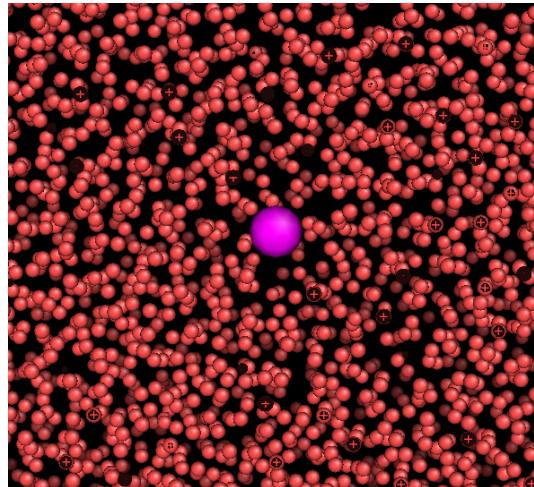
Explicit solvent



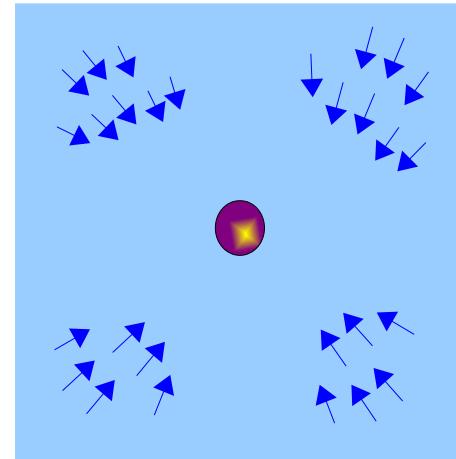
- Water box
- Periodic boundary conditions

Implicit solvent

“explicit”



“implicit”

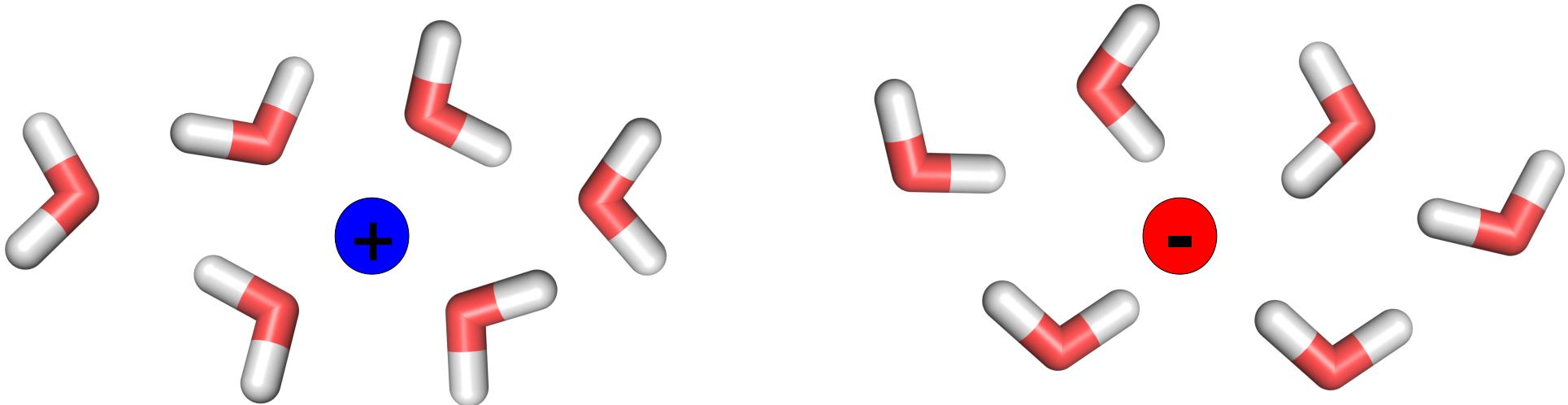


Dielectric continuum

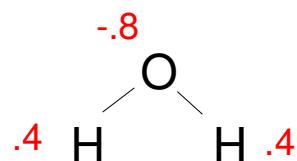
Electrostatic energy of a pair of ions: $U = q q' / \epsilon r$

$\epsilon = 80$ = dielectric constant of water

Ion-ion interactions are greatly reduced in water



From a distance (a nanometer), each ion's charge is partly compensated by nearby water oxygens and hydrogens.



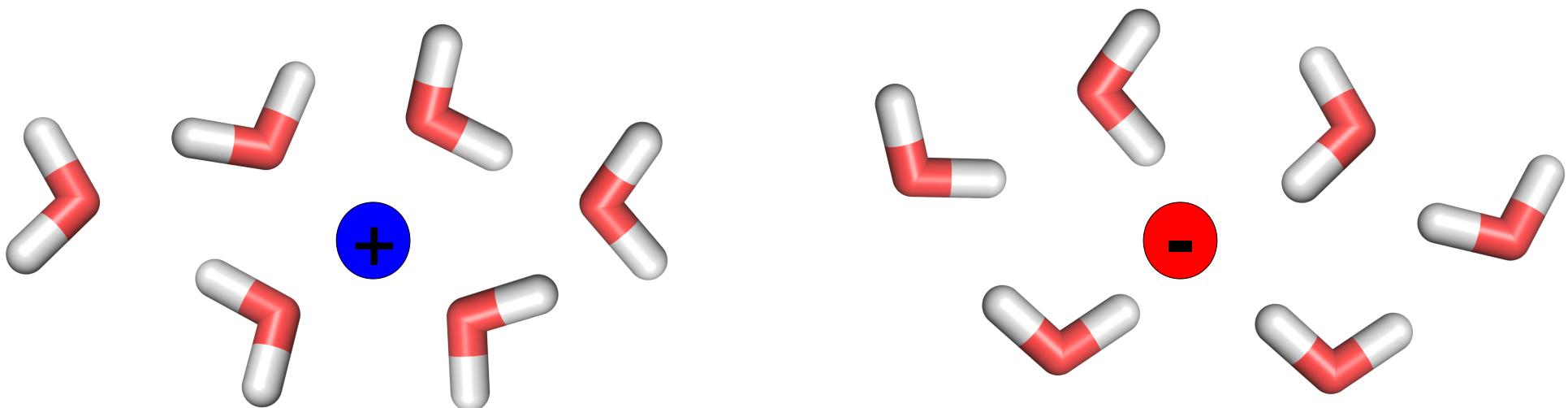
$$U = q q' / \epsilon r$$

Importance of solvent: Mg^{2+} – HPO_4^{2-} association

Association energy in vacuum: **-515 kcal/mol**

In water: **-3.7 kcal/mol**

[standard state conditions]

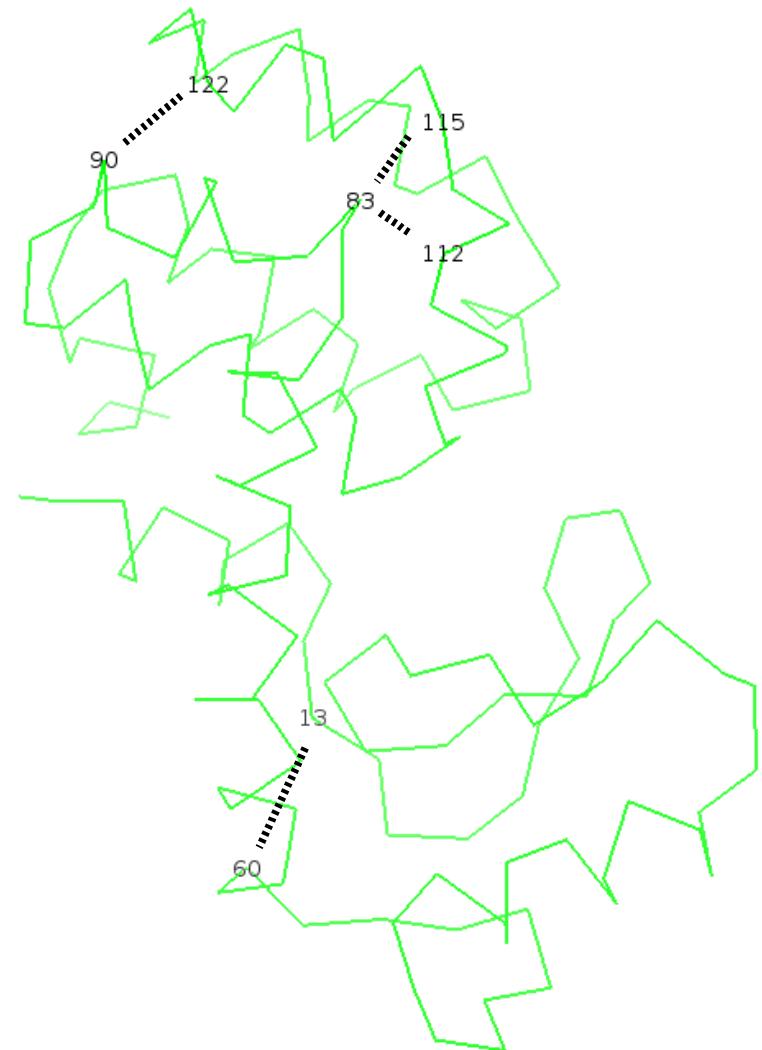


Water has a “screening” effect, reduction by almost 100.

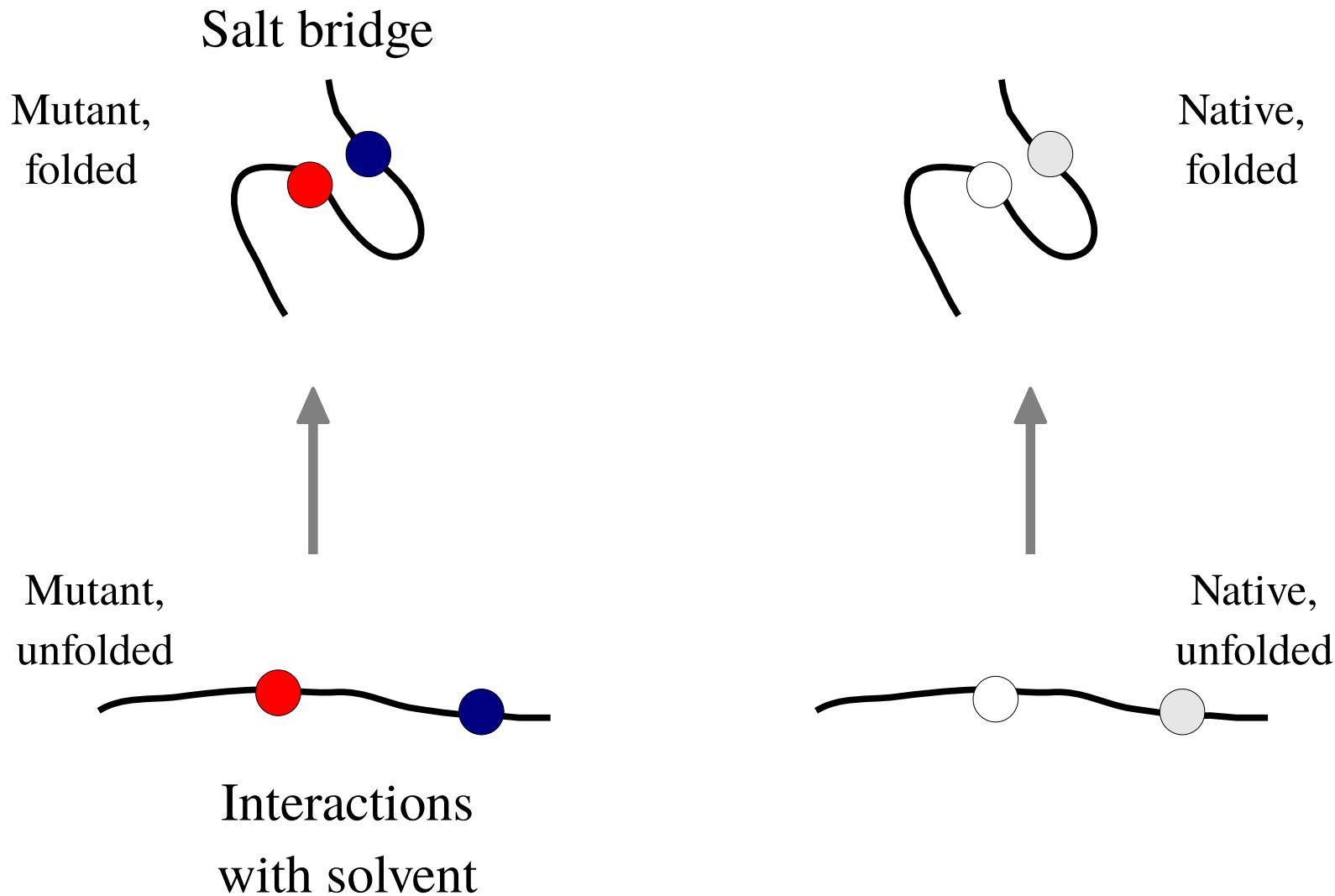
Effect of salt bridges on stability of T4 lysozyme

Dao-Pin, Nicholson, Matthews (1991) Biochemistry, 30:7142

Mutation	$\Delta\Delta G$ (kcal/mol)
K60H/L13D	-2.8
K83H/A112D	-1.5
S90H/Q122D	-2.2
T115E (...K83)	+0.3
S90H	-1.1
K60H	-0.1
K83H	-0.5



Small contribution, compared to enormous vacuum interaction energies

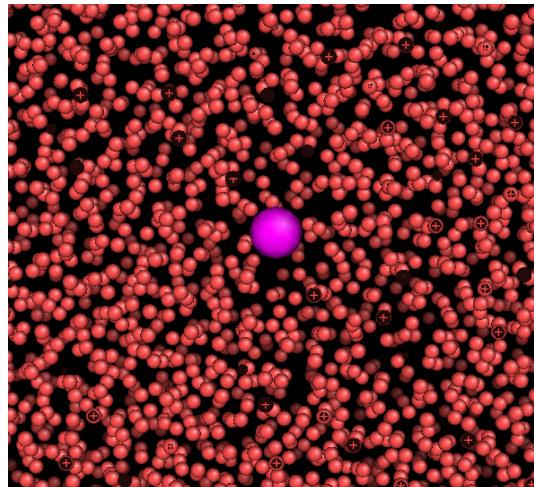


Compensating effects

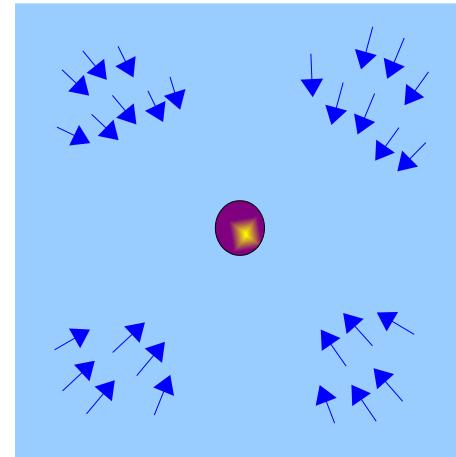
**90% of salt bridges are
at the surface**

Implicit solvent

“explicit”



“implicit”



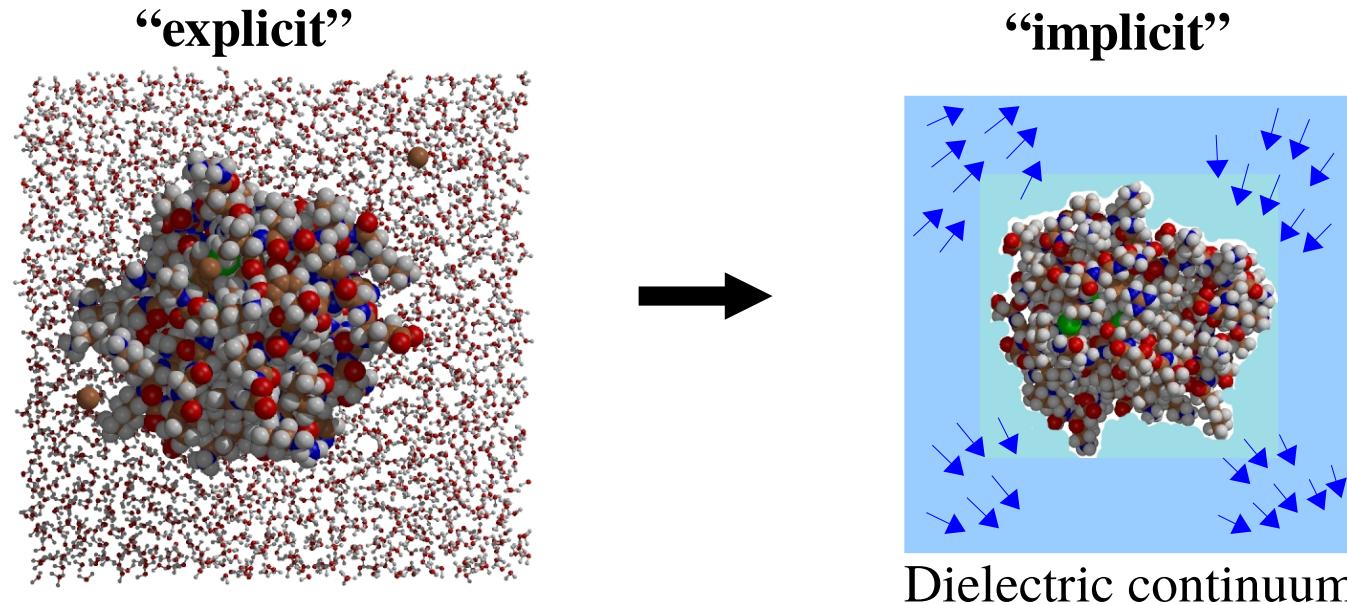
Dielectric continuum

Electrostatic energy of a pair of ions: $U = q q' / \epsilon r$

$\epsilon = 80$ = dielectric constant of water

Fine for a few ions in water.....

Implicit solvent

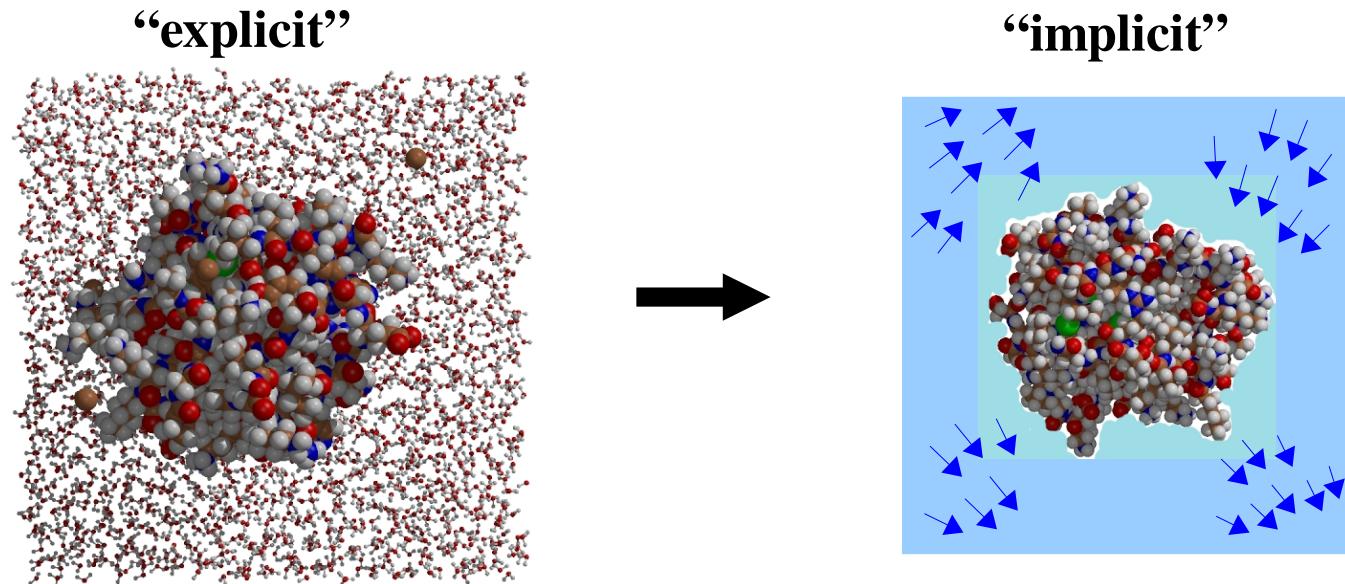


A biomolecule is more complex: heterogeneous system

Initially, we can ignore this complexity: $U = q q' / \epsilon r$

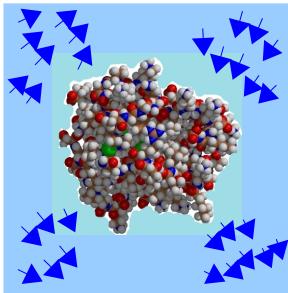
Going further: more sophisticated model

Implicit solvent: “generalized Born” model (GB)



- Buried charges: $E = q q' / r$
- Exposed charges: $E = q q' / \epsilon r$
 $\epsilon = 80$
- Intermediate cases: **interpolation**

“Generalized Born” solvent model (GB)



$$U_{\text{elec}} = U_{\text{Coul}} + U_{\text{GB}}$$

$$U_{\text{Coul}} = \sum_{i,j} q_i q_j / r_{ij}$$
 interaction with no solvent

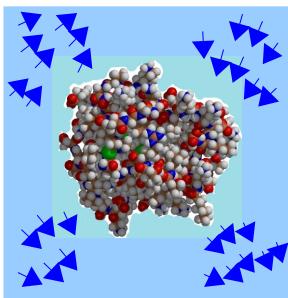
$$U_{\text{GB}} = (-1 + 1/\epsilon) \sum_{i,j} q_i q_j F_{\text{GB}}(i,j)$$
 term due to solvent

$$F_{\text{GB}}(i,j) = 1/(r_{ij}^2 + b_{ij}^2) \exp[-r_{ij}^2/4b_{ij}^2]^{1/2}$$

$$b_{ij} = (b_i b_j)^{1/2}$$

b_i, b_j = distance au solvant des atomes i, j

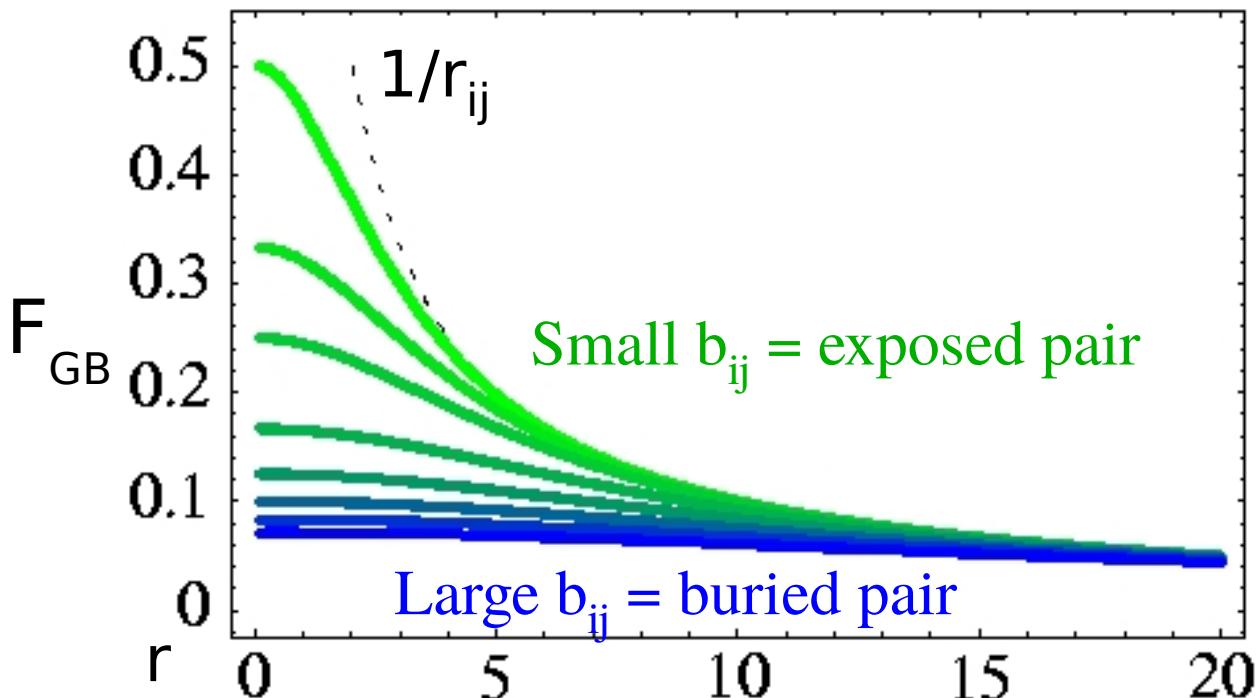
“Generalized Born” solvent model (GB)



$$U_{\text{elec}} = U_{\text{Coul}} + U_{\text{GB}}$$

$$U_{\text{GB}} = (-1 + 1/\varepsilon) \sum_{i,j} q_i q_j F_{\text{GB}}(i,j)$$

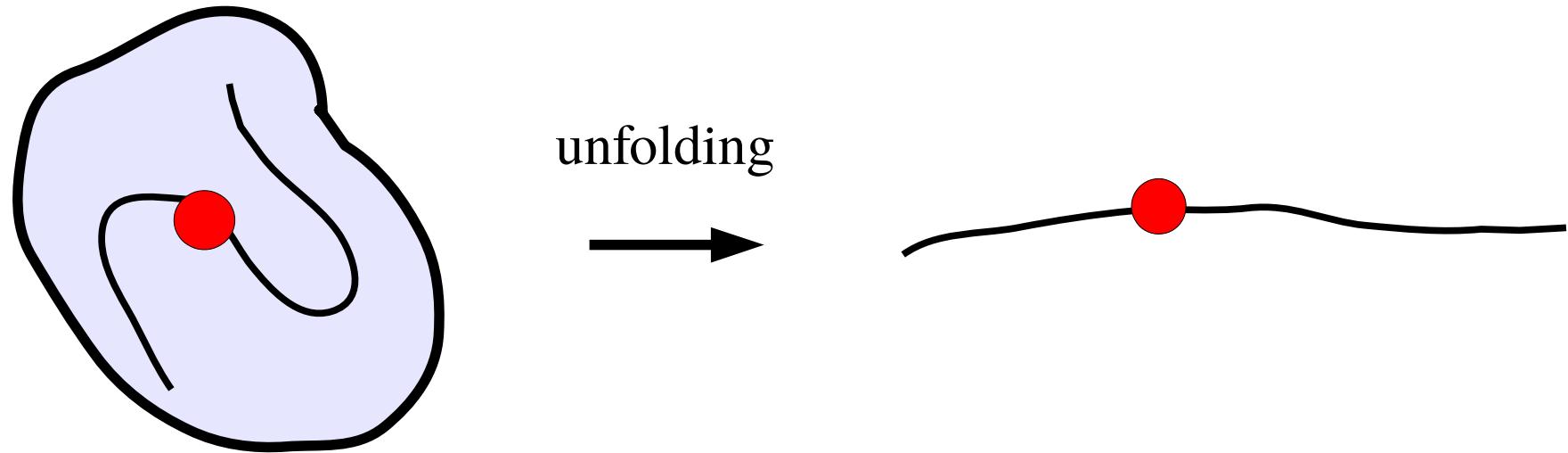
$$F_{\text{GB}}(i,j) = 1/(r_{ij}^2 + b_{ij}^2) \exp[-r_{ij}^2/4b_{ij}^2]^{1/2}$$



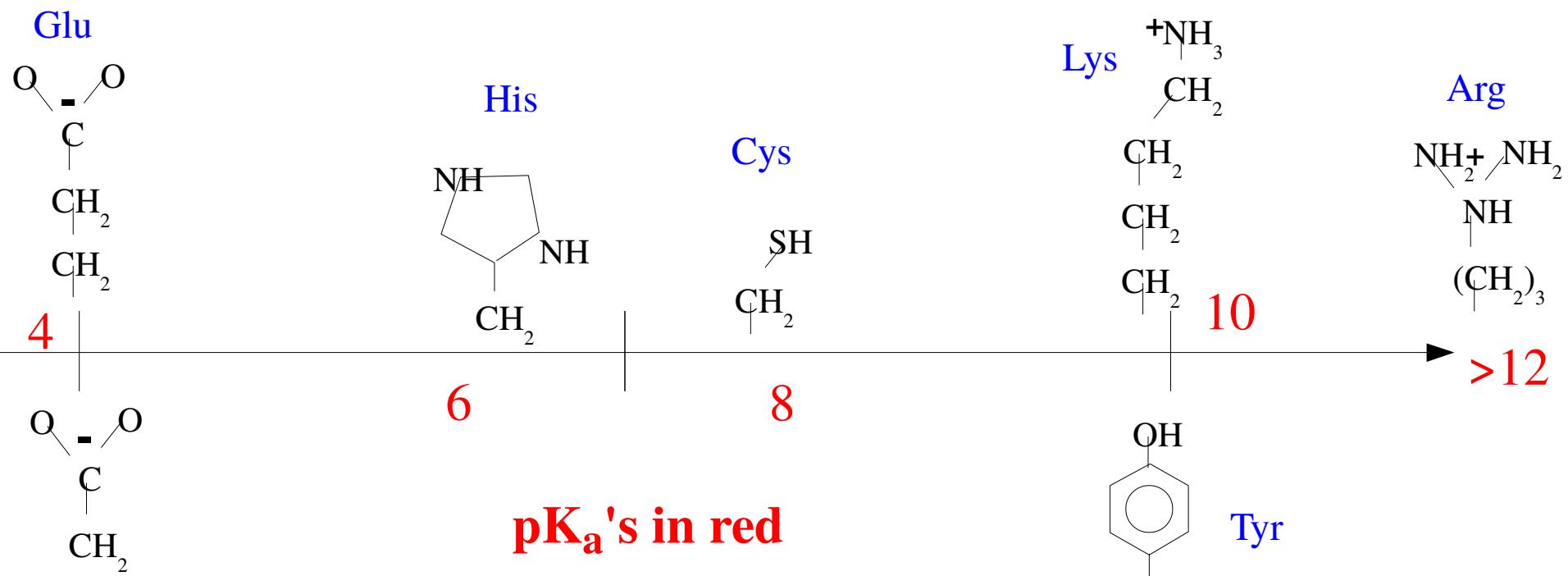
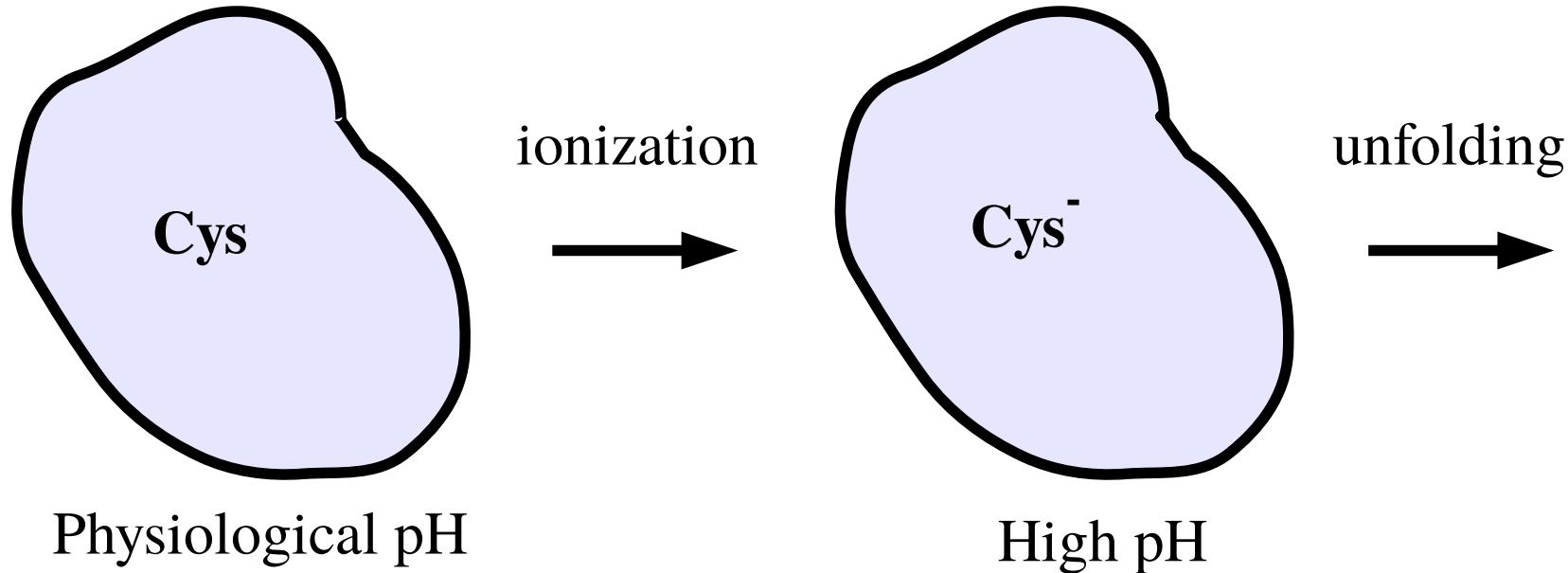
Small b_{ij} : compensation between U_{Coul} et U_{GB} , U_{elec} close to $U_{\text{Coul}}/\varepsilon$

Large b_{ij} : little compensation, U_{elec} close to U_{Coul}

A buried charge would destabilize a protein: why?



A buried charge would destabilize a protein: why?

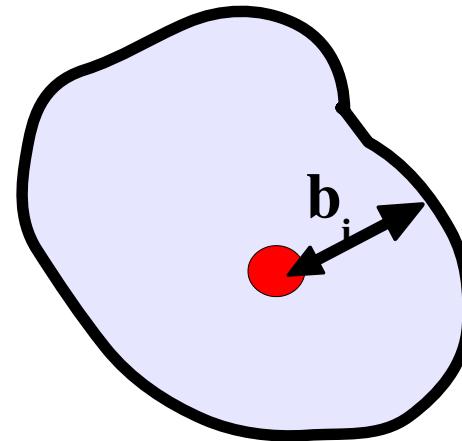


How does GB handle this?

U_{GB} contains “ii” terms

A single charge:

$$U_{GB} = (-1 + 1/\epsilon) q_i^2 / b_i \\ \cong -q_i^2 / b_i$$



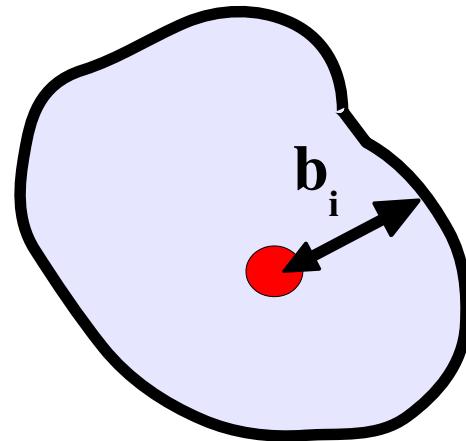
$1/\epsilon$ is small, so $U_{GB} = -q_i q_i / b_i$:

$q_i \leftrightarrow$ solvent interaction = as if q_i sees a charge $-q_i$ at a distance of b_i

“Self” energy term GBSE in XPLOR

U_{GB} contains “ii” terms

To get ii term, let $r_{ij} \rightarrow 0$:



$$U_{GB} = (1/\varepsilon - 1) \sum_{ij} q_i q_j / (r_{ij}^2 + b_i b_j \exp[-r_{ij}^2/4b_i b_j])^{1/2}$$

$$F_{GB}(r_{ij})$$

$i = j ; r \rightarrow 0 ; F_{GB} \rightarrow 1/b_i ; U_{GB} \rightarrow (1/\varepsilon - 1) q_i^2 / b_i$

Polarity of side chain analogs with GB

Some ions

ion	radius ^a (Å)	ΔG_{exp}^b (kcal/mol)
Li ⁺	0.78	-122.1
Na ⁺	0.98	-98.2
K ⁺	1.33	-80.6
Rb ⁺	1.49	-75.5
Cs ⁺	1.65	-67.8
Mg ²⁺	0.78	-455.5
Mn ²⁺	0.91	-437.8
Ca ²⁺	1.06	-380.8
Sr ²⁺	1.27	-345.9
Ba ²⁺	1.43	-315.5

Cf also Fersht;
 Structure & mechanism
 in protein science;
 Freeman.

	solute	area	Exp.	ACE
1	acetamide	216	-9.70	-13.6
2	acetic acid	212	-6.73	-10.2
3	methane thiol	192	-1.24	-6.6
4	propionamide	247	-9.40	-13.3
5	propionic acid	245	-6.47	-9.3
6	4-methyl imidazole	259	-10.25	-9.7
7	3-methyl indole	331	-5.90	-4.9
8	p-cresol	303	-6.13	-5.3
9	ethyl methyl sulfide	246	-1.49	-3.0
10	methanol	174	-5.10	-8.9
11	ethanol	206	-4.95	-7.5
12	toluene	297	-0.76	-0.9
13	butyl amine	272	-4.38	-1.1
14	ethane thiol	220	-1.30	-5.9
15	dimethyl sulfide	222	-1.54	-3.3
16	N-methyl acetamide	254	-10.07	-9.4
17	acetone	233	-3.80	-3.5
18	butanoic acid	274	-6.35	-9.3
19	1-propanol	236	-4.84	-6.7
20	2-propanol	235	-4.75	-6.3
21	2-methyl-2-propanol	260	-4.52	-5.2
22	methyl acetate	246	-3.31	-2.3
23	ethyl acetate	280	-3.09	-2.7
24	methyl amine	178	-4.57	-4.0
25	ethyl amine	209	-4.56	-3.3
26	propyl amine	241	-4.45	-3.0
27	benzene	275	-0.88	-2.4
28	phenol	253	-6.50	-5.4
29	dimethyl acetamide	272	-8.54	-5.1
rms error				1.33

Solvation free
 energies
 (transfer vapor
 →water) in
 kcal/mol.

J Comp Chem,
 1999, 20:322

Exc: expliquer le dépliement des protéines à fort et faible pH

Exercice:

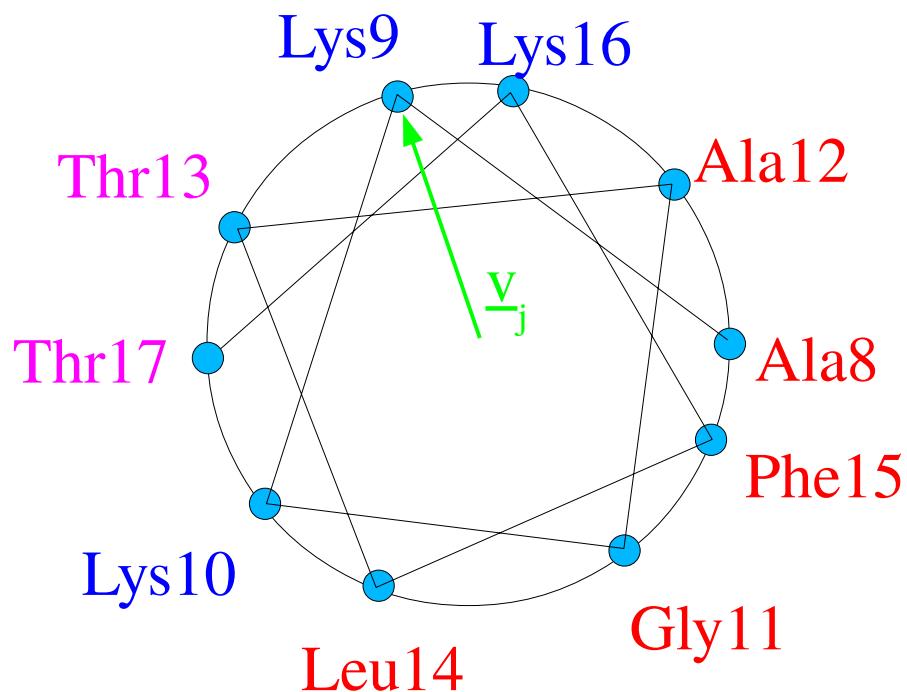
En utilisant XPLOR et le modèle GB, estimer l'énergie (libre) d'une charge ponctuelle placée sur un atome enfoui au centre du Trpcage, toutes les autres charges atomiques étant mises à zéro.

Comparer à ce que vous obtenez avec une charge placée sur un acide aminé seul (le reste de la protéine étant enlevée).

On admet que quand le Trpcage se replie, l'énergie libre varie de 0 à -5 kcal/mol. Quelle est, à l'équilibre, la fraction de protéine repliée?

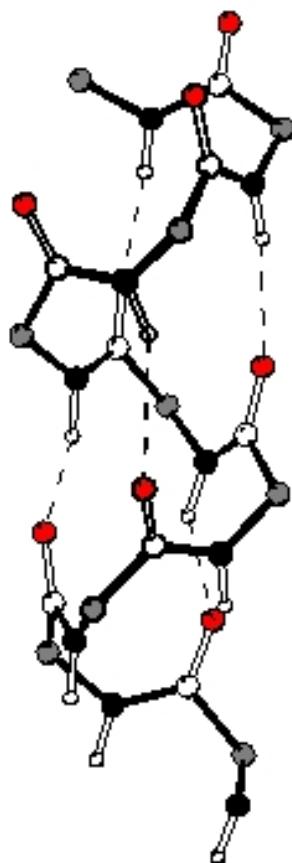
Nous avons vu ci-dessus l'effet d'une “mutation” sur l'énergie de la protéine repliée et celle de la protéine dépliée. (On utilisera le calcul fait pour un acide aminé isolé pour représenter la protéine dépliée; est-ce que c'est raisonnable?) Pour la protéine “mutée”, quelle est la fraction de protéine repliée à l'équilibre?

Secondary structure prediction

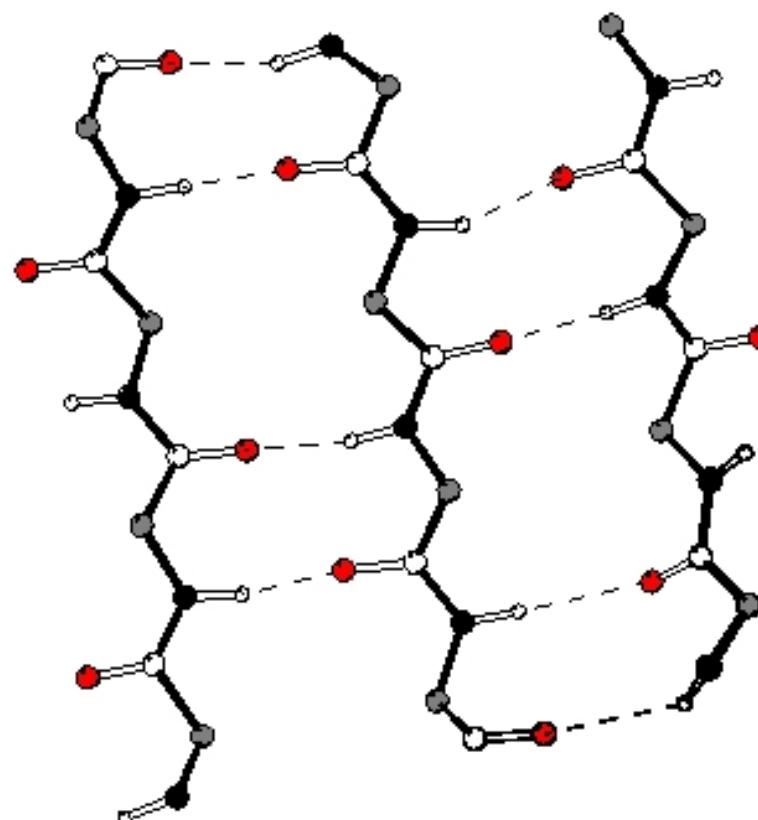


Secondary structure prediction

Hélice α



Feuillet β



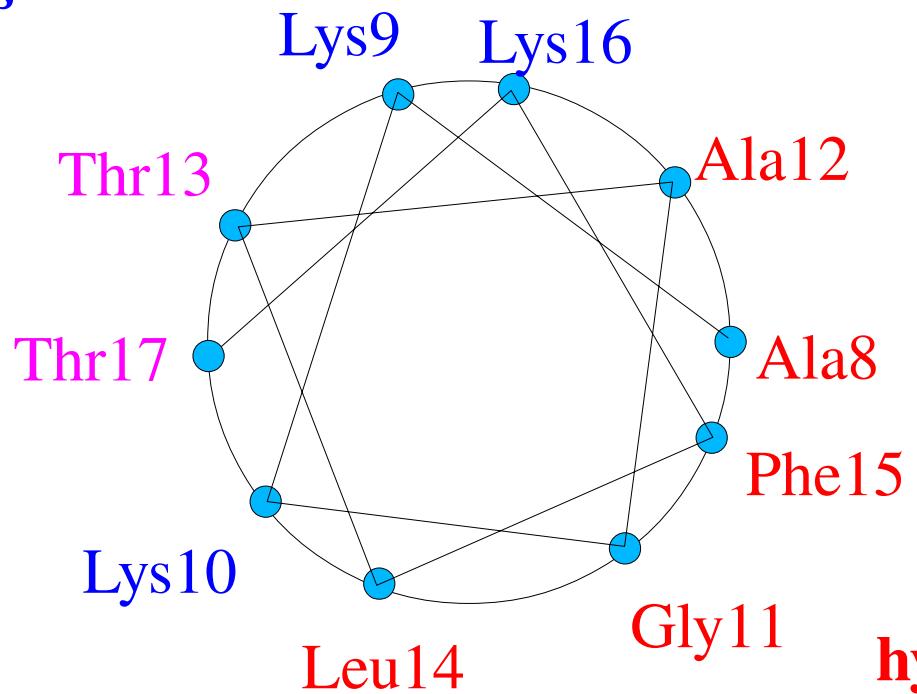
3.6 residues per turn

Pseudo-periodicity 2

Secondary structure prediction

Helices are often amphipathic, with a buried and an exposed side

hydrophiles

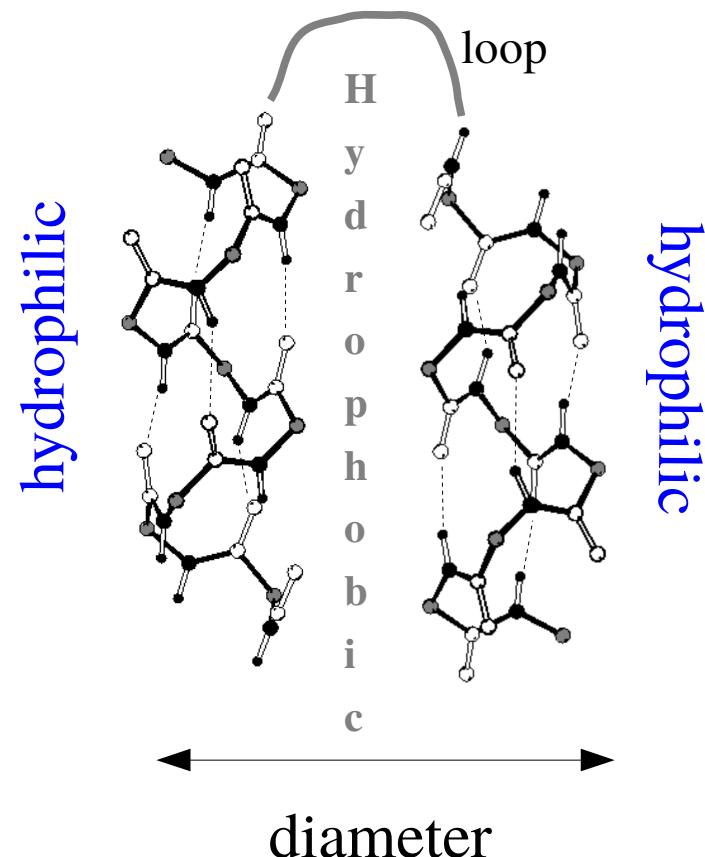


“Helical wheel”

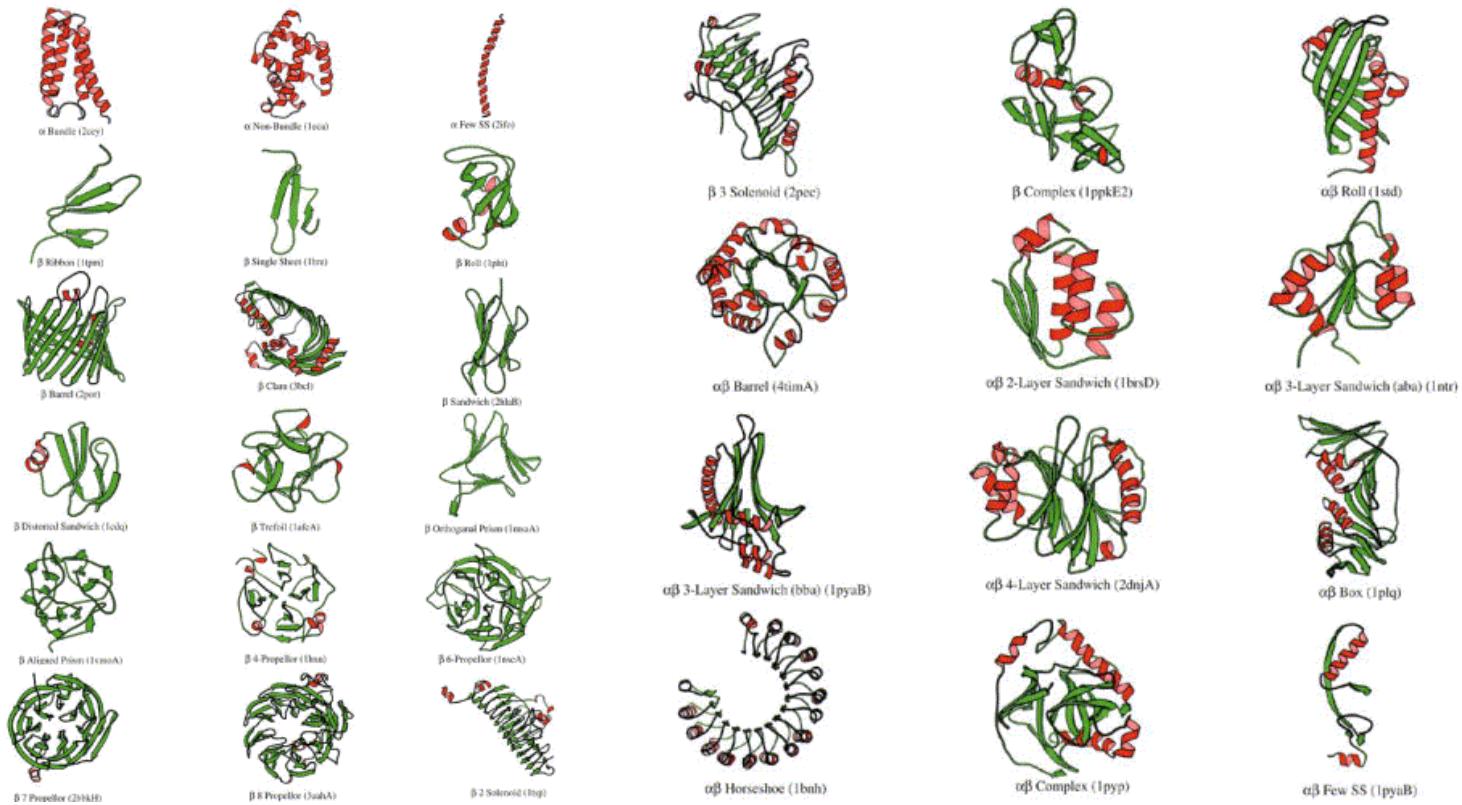
Simplified, amphipathic helix

In proteins, the diameter of a domain often corresponds to a few secondary structure elements

Helices are often amphipathic, with a buried and an exposed side

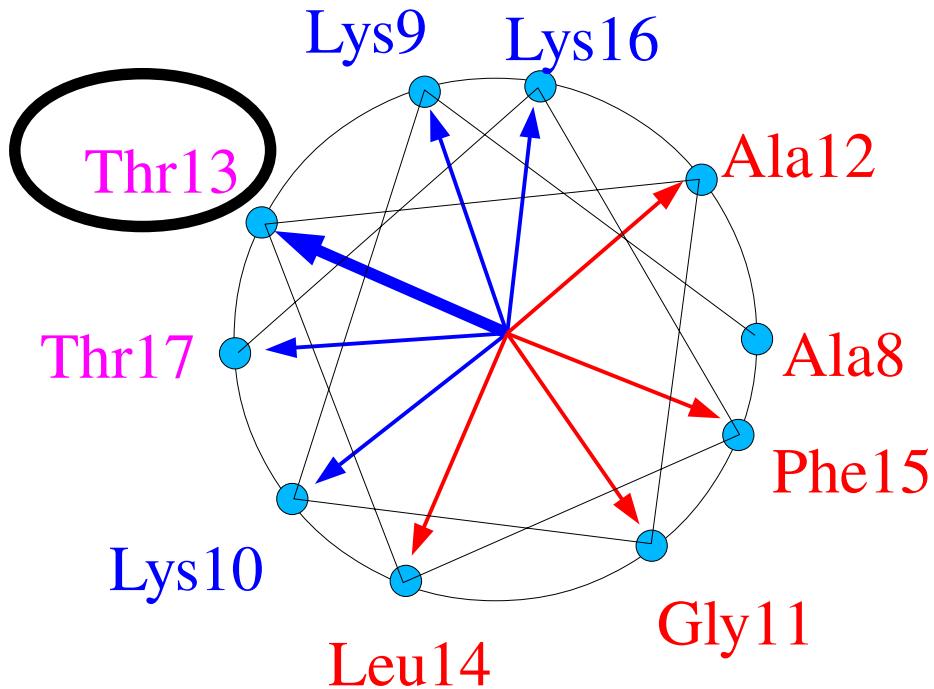


In proteins, the diameter of a domain often corresponds to a few secondary structure elements

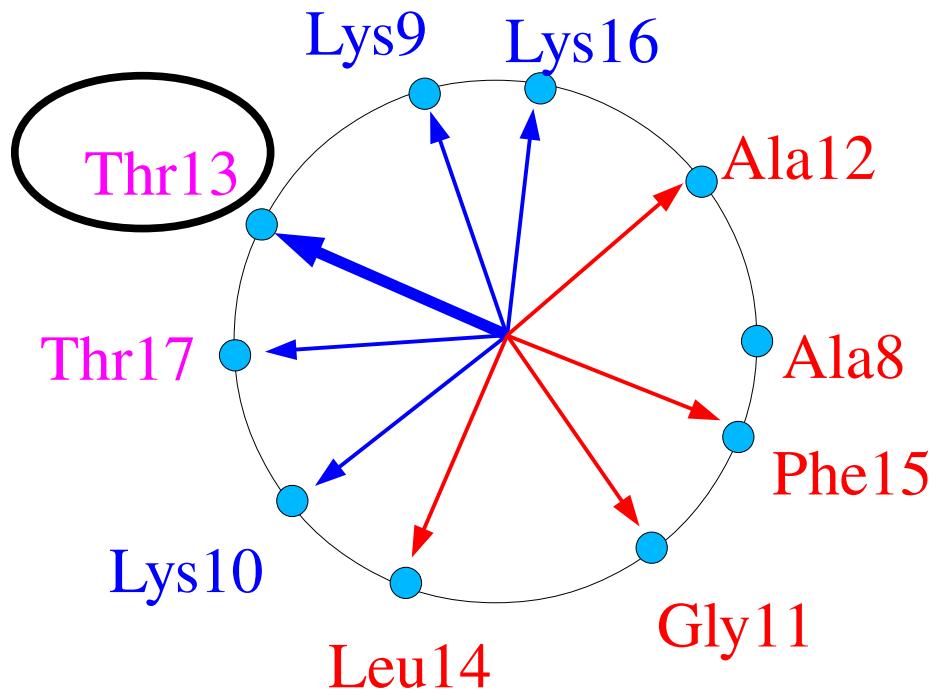


30 “architectures” from CATH

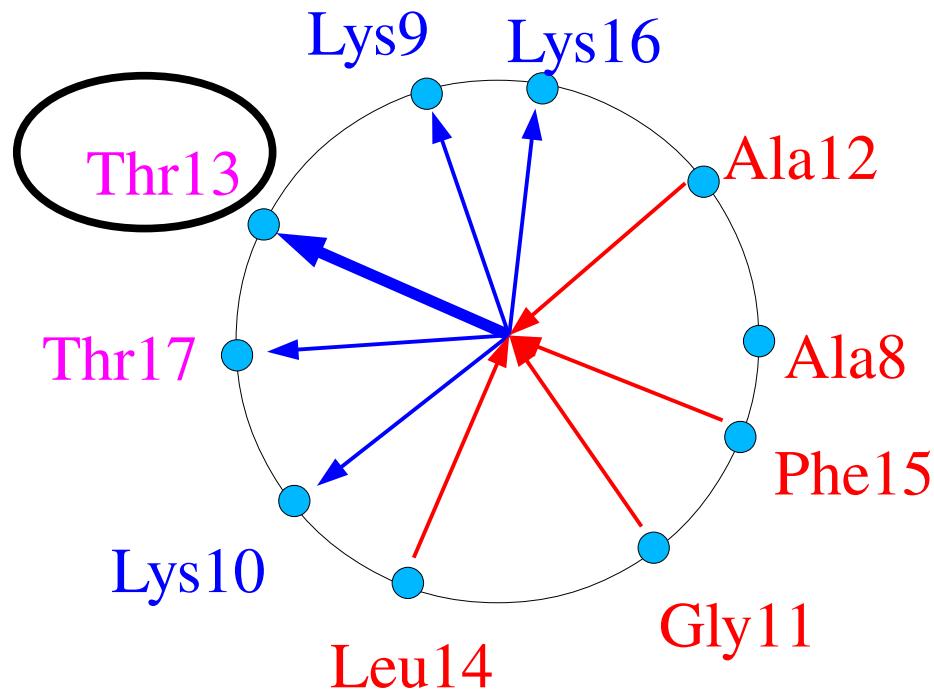
Consider an amino acid, say **Thr13**. We assume it is in a helix, shown below. Its neighbors 9, 10, 11, 12, **13**, 14, 15, 16, 17 define radial vectors:



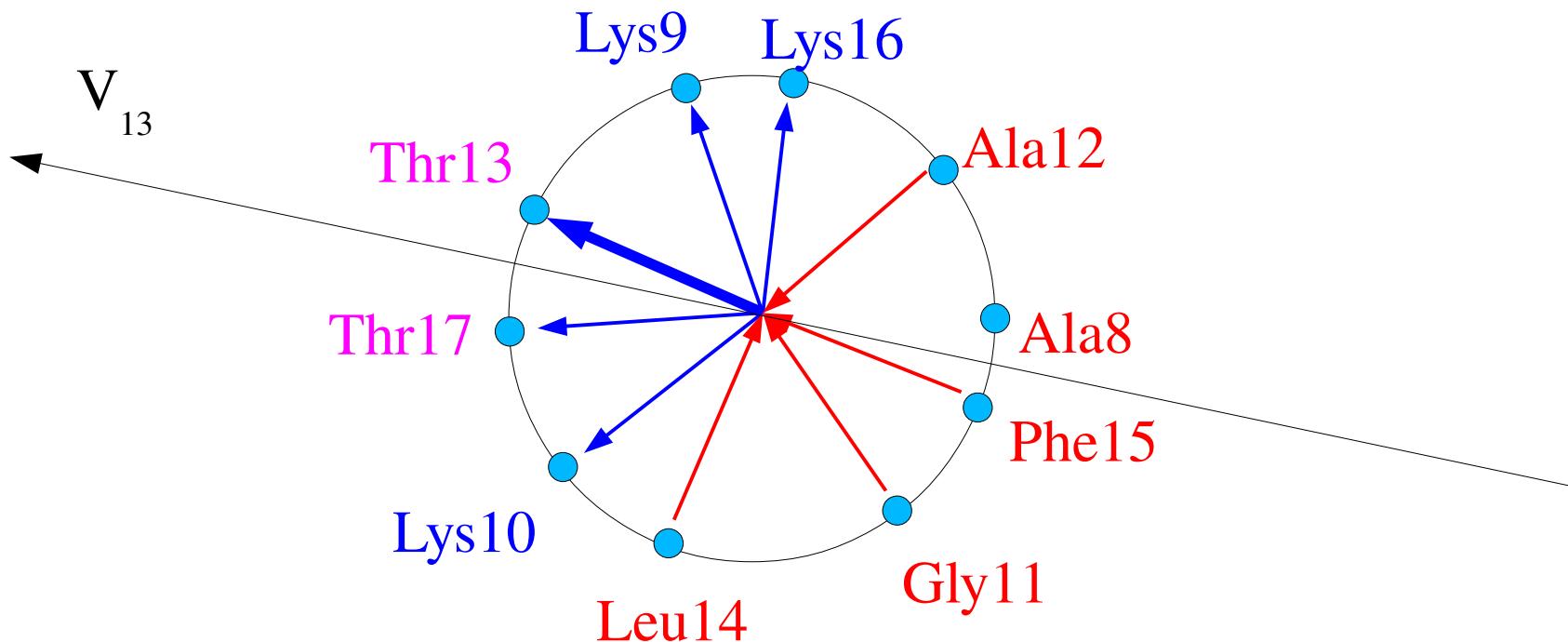
Consider an amino acid, say **Thr13**. We assume it is in a helix, shown below. Its neighbors 9, 10, 11, 12, **13**, 14, 15, 16, 17 define radial vectors:



Consider an amino acid, say **Thr13**. We assume it is in a helix, shown below. Its neighbors 9, 10, 11, 12, **13**, 14, 15, 16, 17 define radial vectors. For hydrophobes, reverse the vectors:



Sum the vectors:

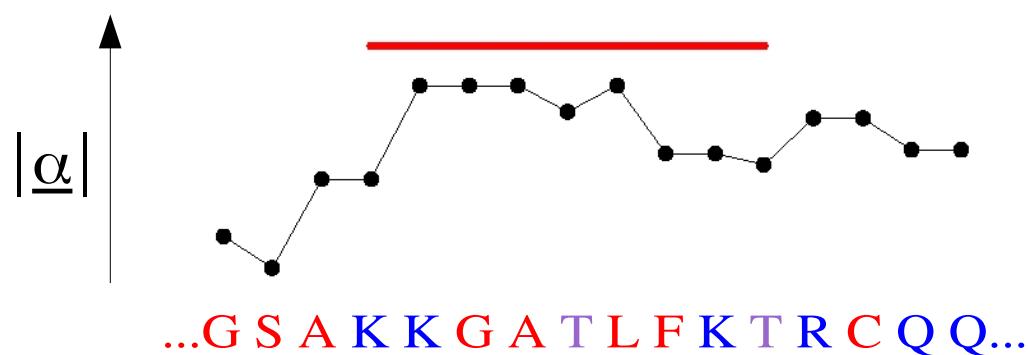
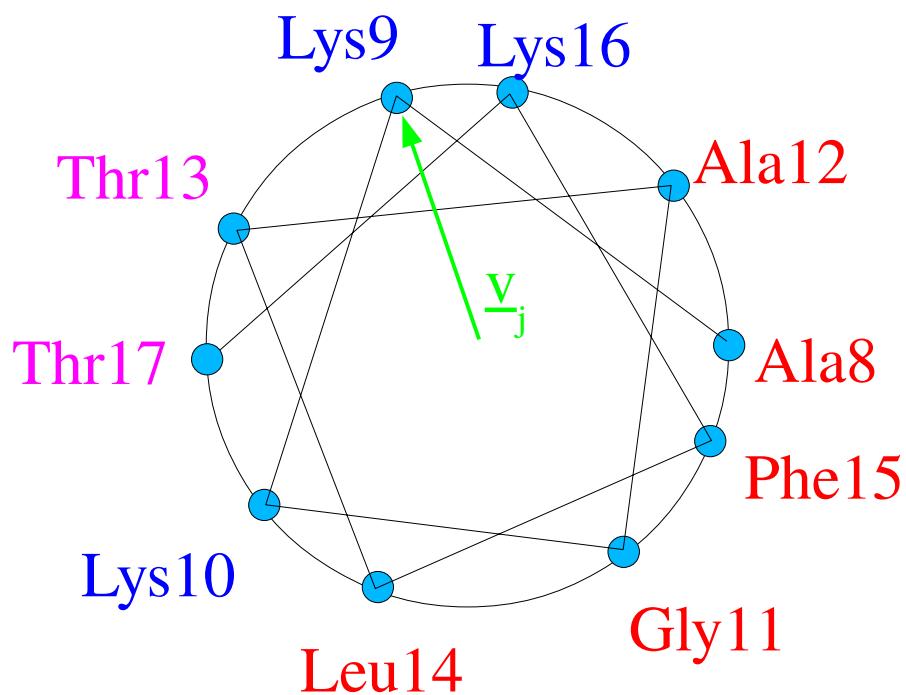


If our hypothesis is correct, we obtain a large vector V_{13} (the v_j add up “constructively”). Else, V_{13} is small.

We have defined a “hydrophobic moment”:

$$\underline{\alpha}_i = \sum_{j=i-m}^{i+m} H_j \underline{v}_j$$

H_j = hydrophobicity,
depends on amino acid type
 $m = 3$ ou 4



The hydrophobic moment reveals amphipathic regions, which have periodicity 3.6.

Beta sheet hypothesis: we define
a different hydrophobic moment

$$\beta_i = \sum_{j=i-m}^{i+m} H_j v_j$$

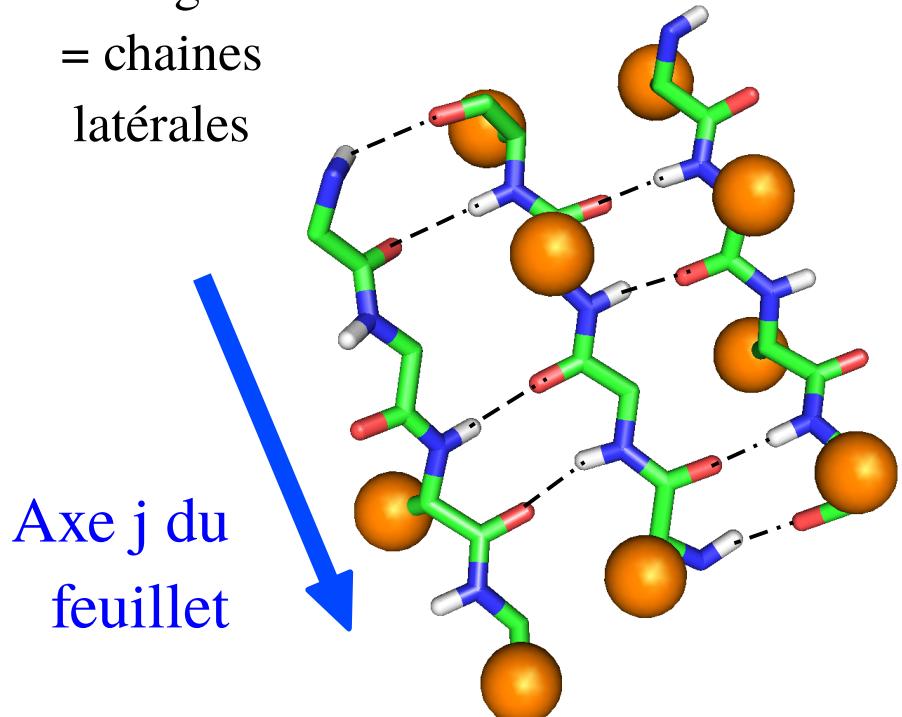
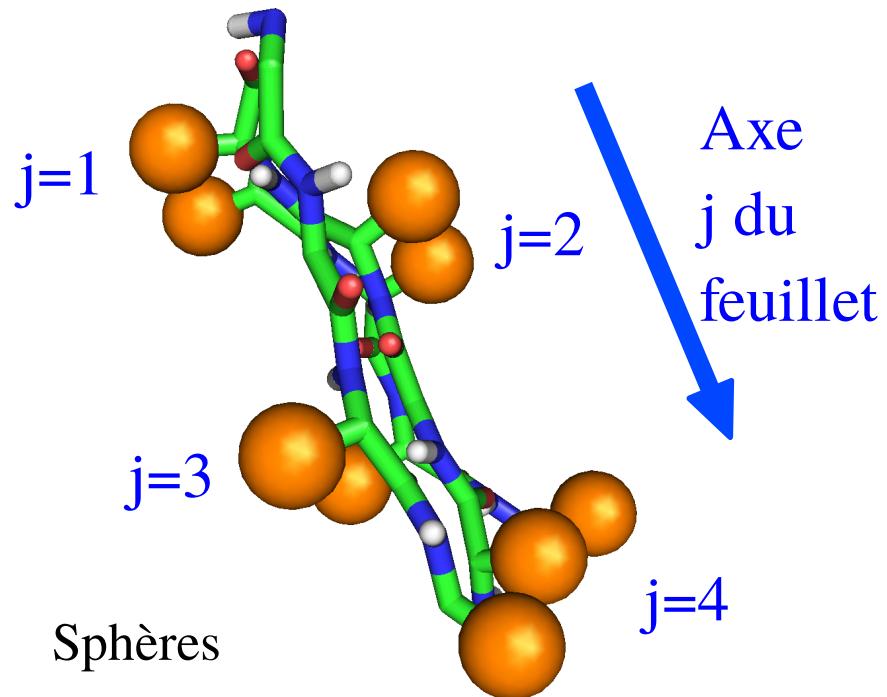
$$v_j = (-1)^j$$

(side chain above or below sheet)

Correct hypothesis:

large $|\beta_i|$

Pseudo-periodicity 2



~ 75% success in general

CYC_MAIZE	1 --ASFSEAPPGNPKAGEKIFKT[KCAQ]C[HTVEKGAGHKQGPNLNGLFGRQSGTTAGYSYSA
CYC_ARUMA	1 --ASF[AEAPPGNPKAGEKIFKT[KCAQ]C[HTVEKGAGHKQGPNLNGLFGRQSGTTAGYSYSA
CYC_ABUTH	1 --ASFQZAPPGBAKAGEKIFKT[KCAQ]C[HTVEKGAGHKQGPNLNGLFGRQSGTTAGYSYSA
CYC_ACENE	1 --ASF[AEAPPGNPAAGEKIFKT[KCAQ]C[HTVDKGAGHKQGPNLNGLFGRQSGTTAGYSYSA
CYC_ALLPO	1 --ATFSZAPPGBZKAGQKIFKL[KCAQ]C[HTVEKGAGHKQGPNLNGLFGRQSGTAAGYSYSA
CYC_WHEAT	1 --ASFSEAPPGNPDAGAKIFKT[KCAQ]C[HTVDAGAGHKQGPNLHGLFGRQSGTTAGYSYSA
CYC_HORSE	1 -----GDVEKGKKIFVQ[KCAQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAPGF[Y]TD
CYC_BOVIN	1 -----GDVEKGKKIFVQ[KCAQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAPGF[SY]TD
CYC_MOUSE	1 -----GDVEKGKKIFVQ[KCAQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAAAGF[SY]TD
CYC_RABBIT	1 -----GDVEKGKKIFVQ[KCAQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAVGF[SY]TD
CYC_ALLMI	1 -----GDVEKGKKIFVQ[KCAQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAPGF[SY]TE
CYC_AP_TPA	1 -----GDIEKGKKIFVQ[KCSQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAEGF[SY]TD
CYC_CHICK	1 -----GDIEKGKKIFVQ[KCSQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAEGF[SY]TD
CYC_ANAPL	1 -----GDVEKGKKIFVQ[KCSQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAEGF[SY]TD
CYC_HUMAN	1 -----GDVEKGKKIFV[KCSQ]C[HTVEKG]GKHTGPNLHGLFGRKTGQAPGYSYTA
CYC_APIME	1 -----GIPAGDPEKGKKIFVQ[KCAQ]C[TIEGG]GK[VGPNL]YGVYGRKTGQAPGYSYTD
CYC7_YEAST	1 MAKESTGFKPGSAKKGATLFKTR[CQ]C[TIEEGG]GPNKVGPNLHGI[FGR]ESGQVKGYSYTD

HHHHHHHHHHHH

EEE

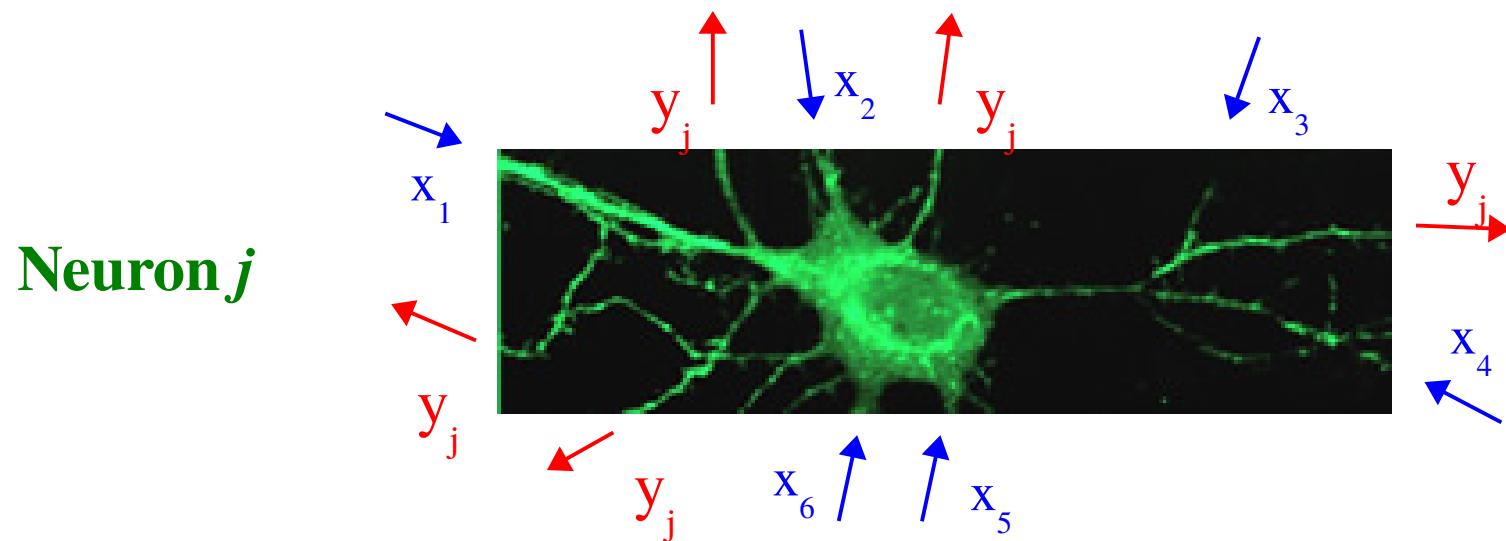
CYC_MAIZE	59 ANKNKAVV[EEN]TLYDYLLNP[KY]IPGTMVFPGLXKPQE[RADLIAYLKEATA-
CYC_ARUMA	59 ANKNMAVIWEE[STLYDYL]LNPK[KY]IPGTMVFPGLXKPQE[RADLIAYLKESTA-
CYC_ABUTH	59 ANKNMAVN[WGEN]TLYDYLLNP[KY]IPGTMVFPGLKKPQDRADLIAYLKZSTA-
CYC_ACENE	59 ANKNMAVN[WGYNT]LYDYLLNP[KY]IPGTMVFPGLKKPQDRADLIAYLKQSTA-
CYC_ALLPO	59 ANKNMAVV[WZZBT]LYDYLLNP[KY]IPGTMVFPGLKKPQDRADLIAYLKESTA-
CYC_WHEAT	59 ANKNKAVV[EEN]TLYDYLLNP[KY]IPGTMVFPGLKKPQDRADLIAYLKKATSS
CYC_HORSE	51 ANKNKGITWKEETLMEYLENP[KY]IPGTMIFAGIKKKTEREDLIAYLKKATNE
CYC_BOVIN	51 ANKNKGITWGEE[TLMEYLENP[KY]IPGTMIFAGIKKKGEREDLIAYLKKATNE]
CYC_MOUSE	51 ANKNKGITWGEE[TLMEYLENP[KY]IPGTMIFAGIKKKGERADLIAYLKKATNE]
CYC_RABBIT	51 ANKNKGITWGEE[TLMEYLENP[KY]IPGTMIFAGIKKKDERADLIAYLKKATNE]
CYC_ALLMI	51 ANKNKGITWGEE[TLMEYLENP[KY]IPGTMIFAGIKKKPERADLIAYLKEATSN]
CYC_AP_TPA	51 ANKNKGITWGEE[TLMEYLENP[KY]IPGTMIFAGIKKKSERADLIAYLKDATSK]
CYC_CHICK	51 ANKNKGITWGEE[TLMEYLENP[KY]IPGTMIFAGIKKKSERV[DLIAYLKD]ATSK]
CYC_ANAPL	51 ANKNKGITWGEE[TLMEYLENP[KY]IPGTMIFAGIKKKSERADLIAYLDATAK]
CYC_HUMAN	51 ANKNKGITWGEE[TLMEYLENP[KY]IPGTMIFAGIKKKSERADLIAYLKKATNE]
CYC_APIME	55 ANKGKGITWNE[TLFEYLENP[KY]IPGTMIFAGLKKPQE[RADLIAYLIEQASK-]
CYC7_YEAST	61 ANINKNVKWDDEDSMS[EYLTPNPK[KY]IPGTMIFAGLKKPQE[RADLIAYLIEQASK-]

HHHHHH EEE HHHHHHHHHHHH HHHHH

HHHHHHHHHHHHHHHHHH

A method inspired by neural networks

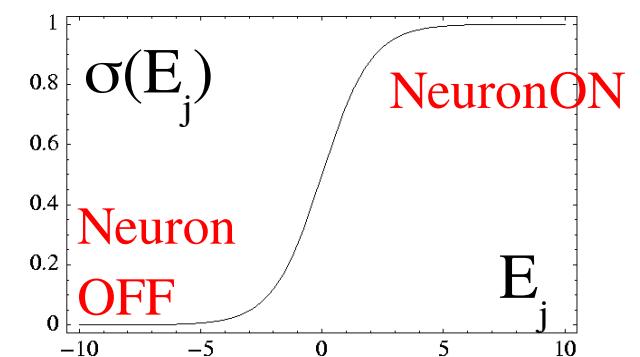
A neuron j receives input signals x_1, x_2, \dots, x_n from nearby cells $1, 2, \dots, n$; sums them, and converts them into an output signal y_j :



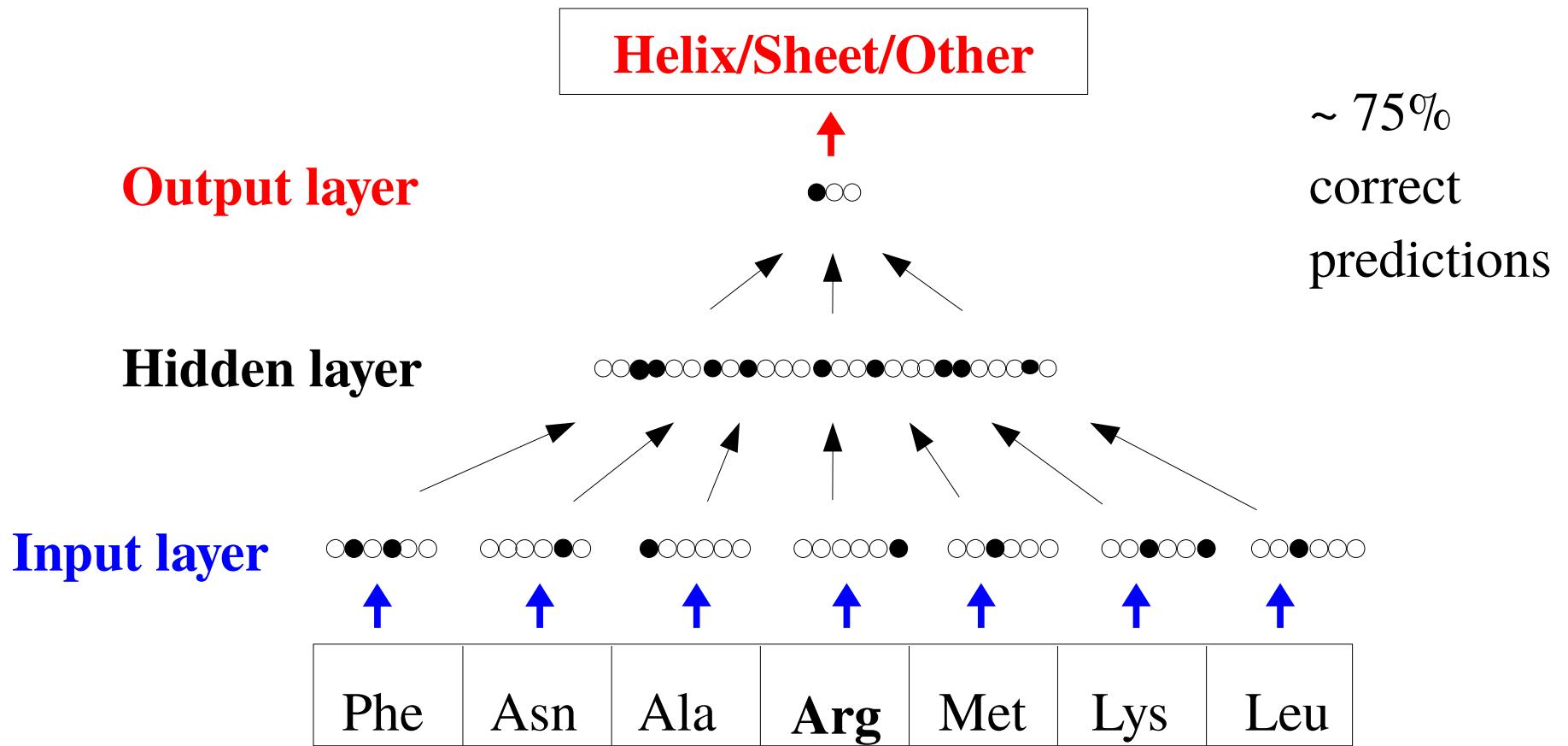
Input signals x_i

are weighted and summed: $E_j = \sum_i w_{ij} x_i$

then converted: $y_j = \sigma(E_j)$

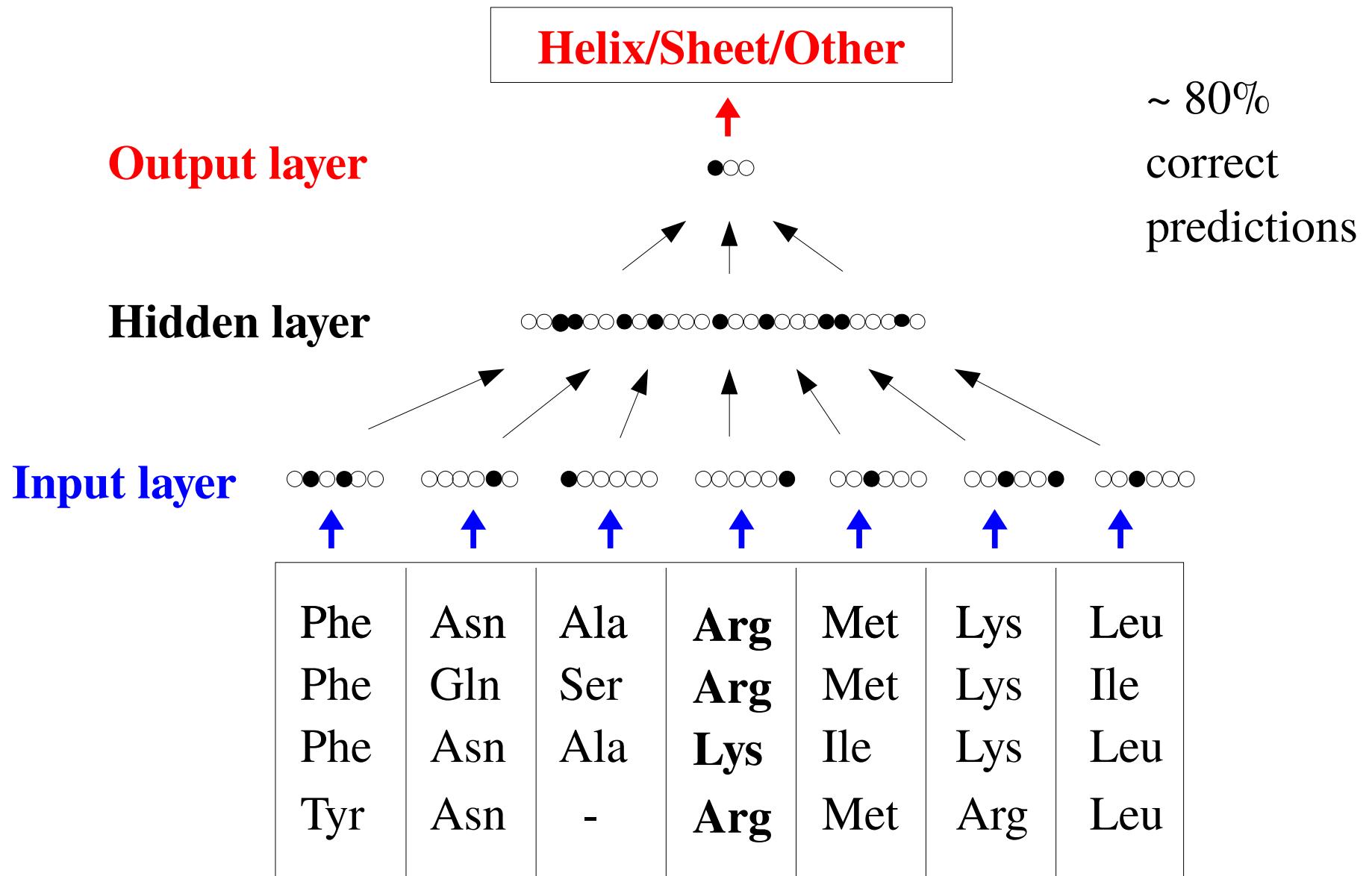


A neural network to predict secondary structure



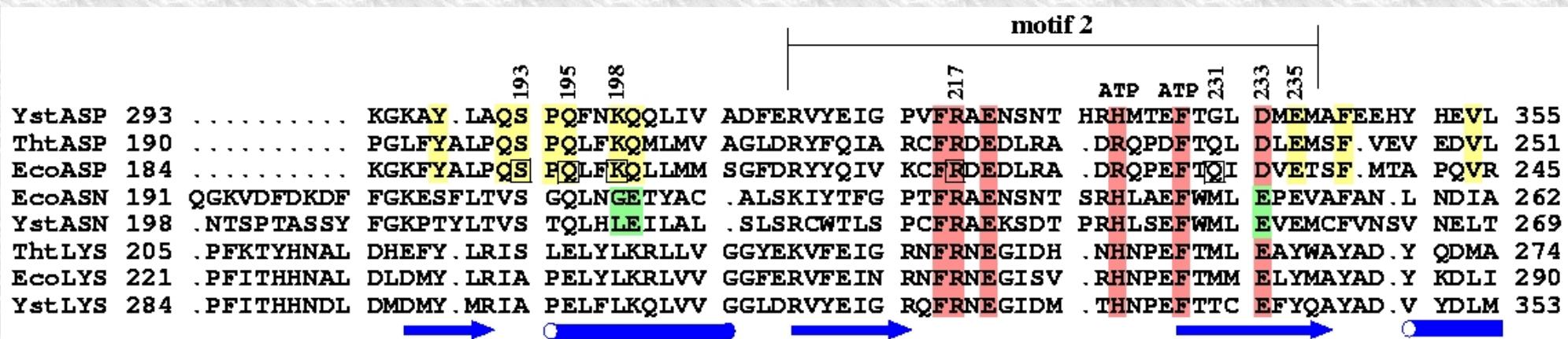
Network architecture is chosen; the connection weights w_{ij} (1212 in this example) are chosen to give correct predictions for known structures (learning phase). At this point, the network has “learned” to assign a secondary structure to any input sequence. Specifically, it predicts the state of the central residue (**Arg**).

A neural network to predict secondary structure

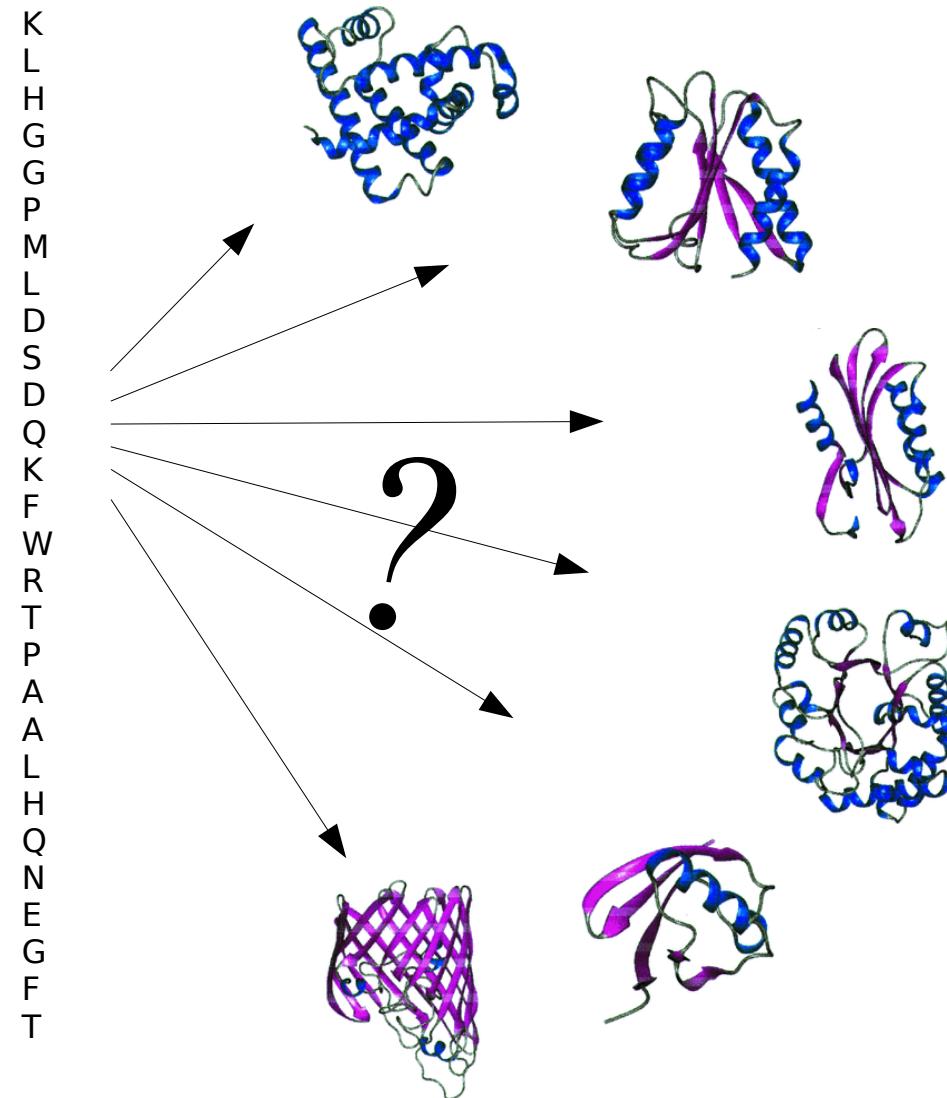


Using a multiple alignment as input

Homology modeling: an overview

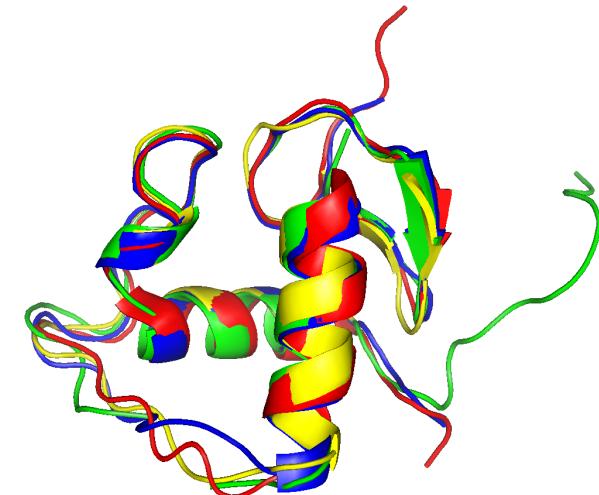


**Homology-based structure prediction: start from
the sequence of a new protein. Check if it has a
known fold through sequence comparaisons**



If so, perform “Homology modelling”

Identify and align homologues of known structure (“templates”)



Do multiple sequence alignment with many homologues

Conserved, well-structured regions: adopt mean backbone of templates

Loops, insertions:

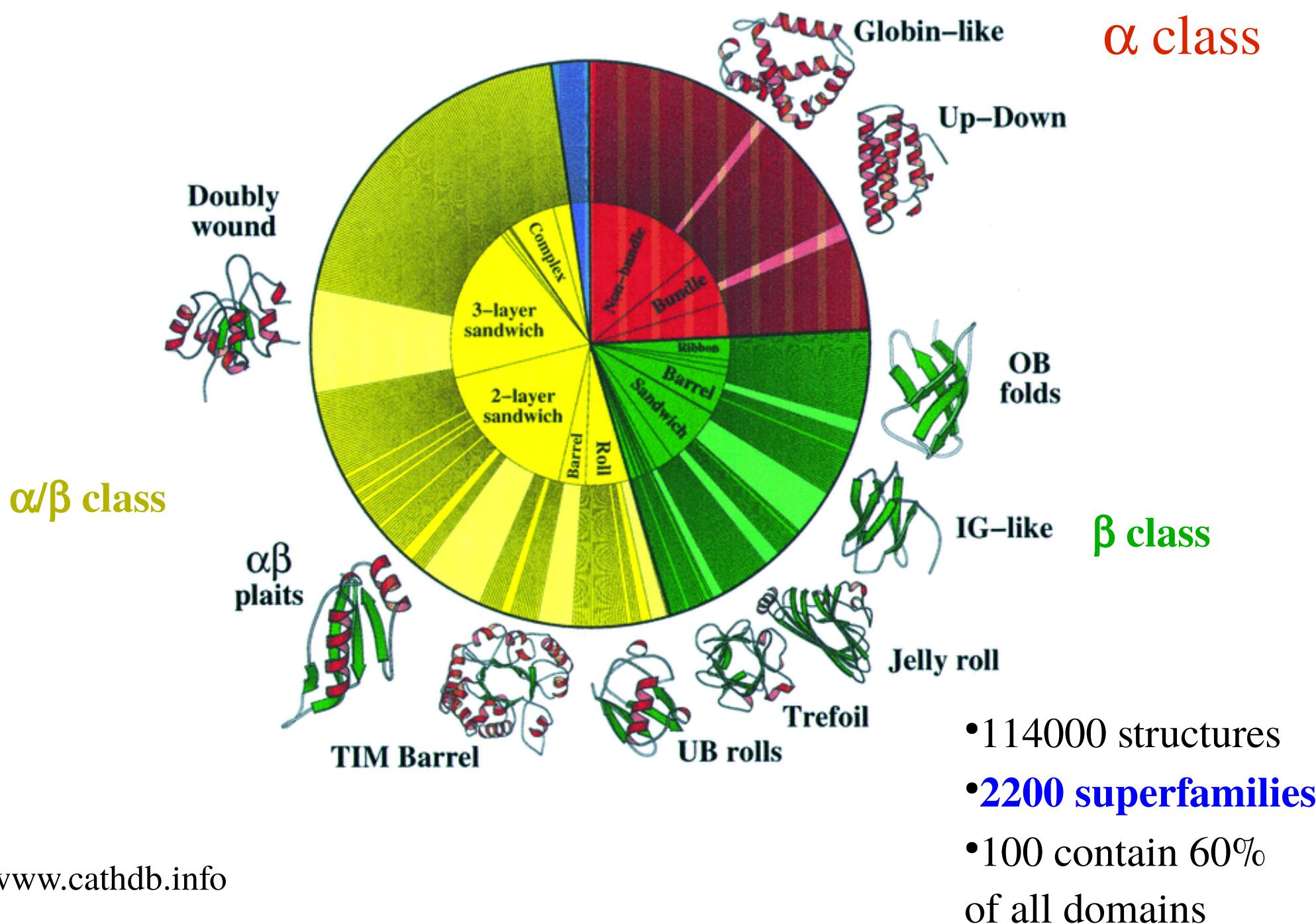
- search Protein Data Bank for similar fragments
- if not, use “ab initio” loop modelling methods

Sidechains: explore and score possible rotamers, possibly with a stochastic approach such as “Monte Carlo”

Model refinement: energy minimization, molecular dynamics

Experimental testing

Fold set is discrete and rather small



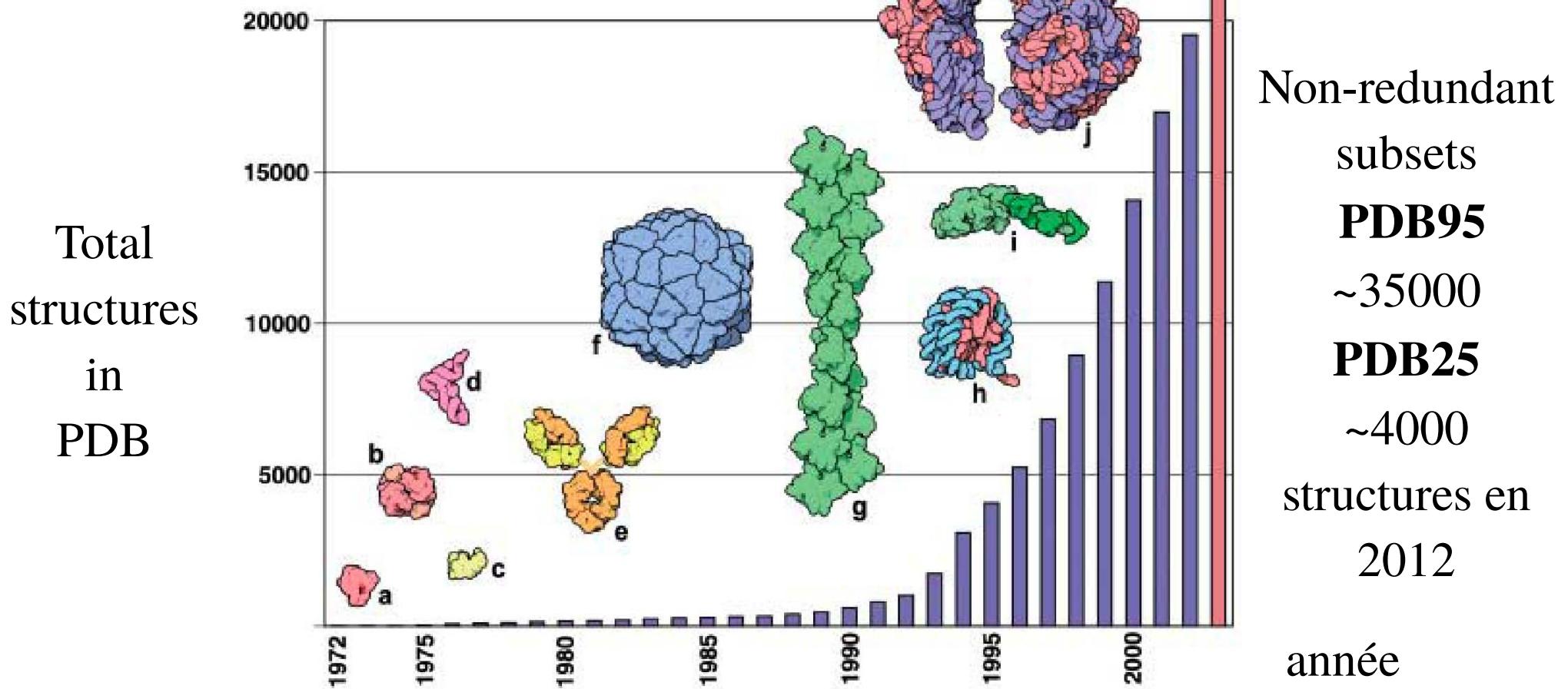
Homologues of androsterone receptor with BLAST

#	Swissprot	Hit	Description	Score (bits)	% Identity			Match Length
					E *	*		
1	P15207	ANDR_RAT	Androgen receptor.	162	1e-40	100		73
6	P19091	ANDR_MOUSE	Androgen receptor.	162	1e-40	100		73
14	Q63449	PRGR_RAT	Progesterone receptor (PR)	136	1e-32	80		72
17	P06401	PRGR_HUMAN	Progesterone receptor (PR)	136	1e-32	80		72
21	P08235	MCR_HUMAN	Mineralocorticoid receptor (MR)	136	1e-32	79		72
33	P04150	GCR_HUMAN	Glucocorticoid receptor (GR)	131	3e-31	77		72
41	Q9YH32	ESR2_ORENI	Estrogen receptor beta (ER-beta)	99	3e-21	58		72
42	Q9YH33	ESR1_ORENI	Estrogen receptor (ER-alpha)	98	4e-21	55		72
:	:	:	:	:	:	:		:
:	:	:	:	:	:	:		:
343	Q9N4Q7	NH13_CAEEL	Nuclear hormone receptor nhr-13	54	8e-08	39		66
344	Q23294	NH11_CAEEL	Nuclear hormone receptor nhr-11	54	8e-08	42		66
345	Q45460	NH54_CAEEL	Nuclear hormone receptor nhr-54	54	1e-07	37		67
346	Q09565	NH20_CAEEL	Nuclear hormone receptor nhr-20	51	7e-07	34		66
347	Q09587	NH22_CAEEL	Nuclear hormone receptor nhr-22	45	5e-05	32		66
349	P17672	E75B_DROME	Ecdysone-induced protein 75B	40	0.001	37		47
351	P20659	TRX_DROME	Trithorax protein.	31	0.74	26		49
355	P98164	LRP2_HUMAN	Lipoprotein receptor.	30	1.7	27		65

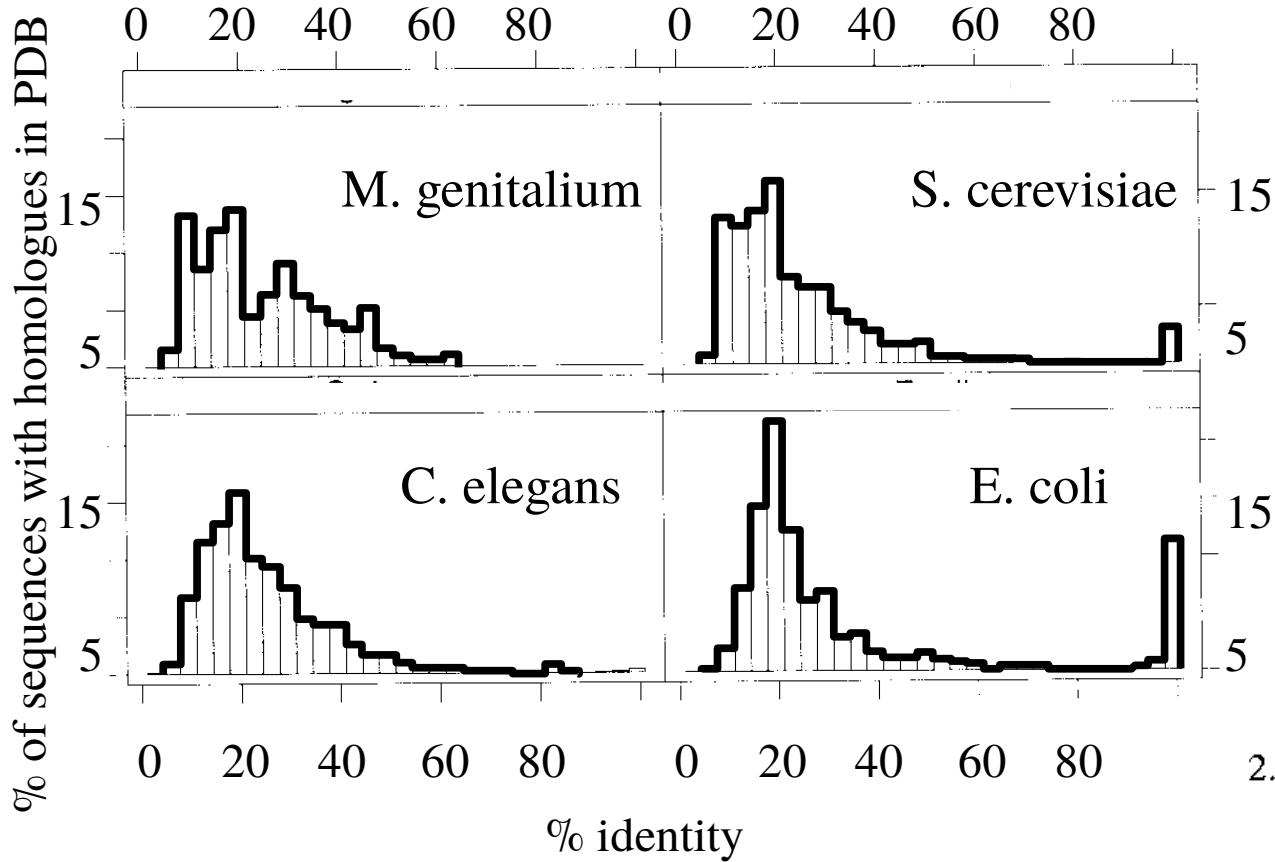
*E = expectation of number of random alignments with a higher score

Structures in “Protein Data Bank”

This chart shows the increase in the total number of structures in the PDB per year, through July 1, 2003, as well as examples of the increasing complexity of these structures. In the 1970's, the first structures available to the scientific community included proteins such as a. myoglobin,^{5,6} b. hemoglobin,^{5,6} and c. lysozyme^{7,8} and other molecules such as d. transfer RNA.^{9,12} In the 1980's, advances in experimental data collection methods allowed much larger structures to be solved, including e. antibodies,^{13,14} and f. entire viruses.¹⁵ By 2003, all aspects of structural science have advanced so that very complex and functionally significant structures could be made accessible to study, including g. actin,¹⁶ the h. nucleosome,¹⁷ i. myosin,¹⁸ j. ribosomal subunits,^{19,21} and the k. calcium pump.²² Structures pictured here were taken from PDB entries 1mbn, 2dhb, 2lyz, 4tna + 6tna, 1fc1 + 1mcp, 2stv, 1atn, 1aoi, 1dfk, 1ffk + 1fka + 1j5e, and 1iwo, respectively. Images were created by David S. Goodsell of The Scripps Research Institute, creator of the PDB Molecule of the Month series.

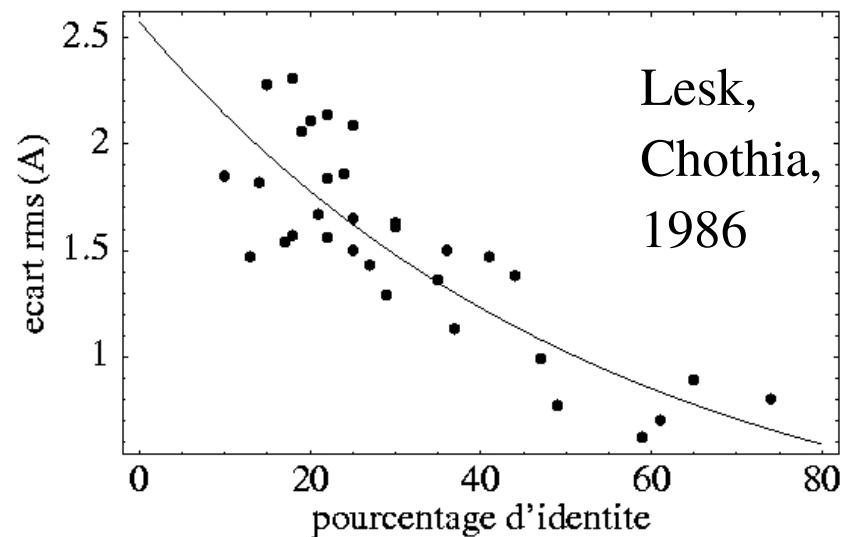


Homology modelling: a template is not always available

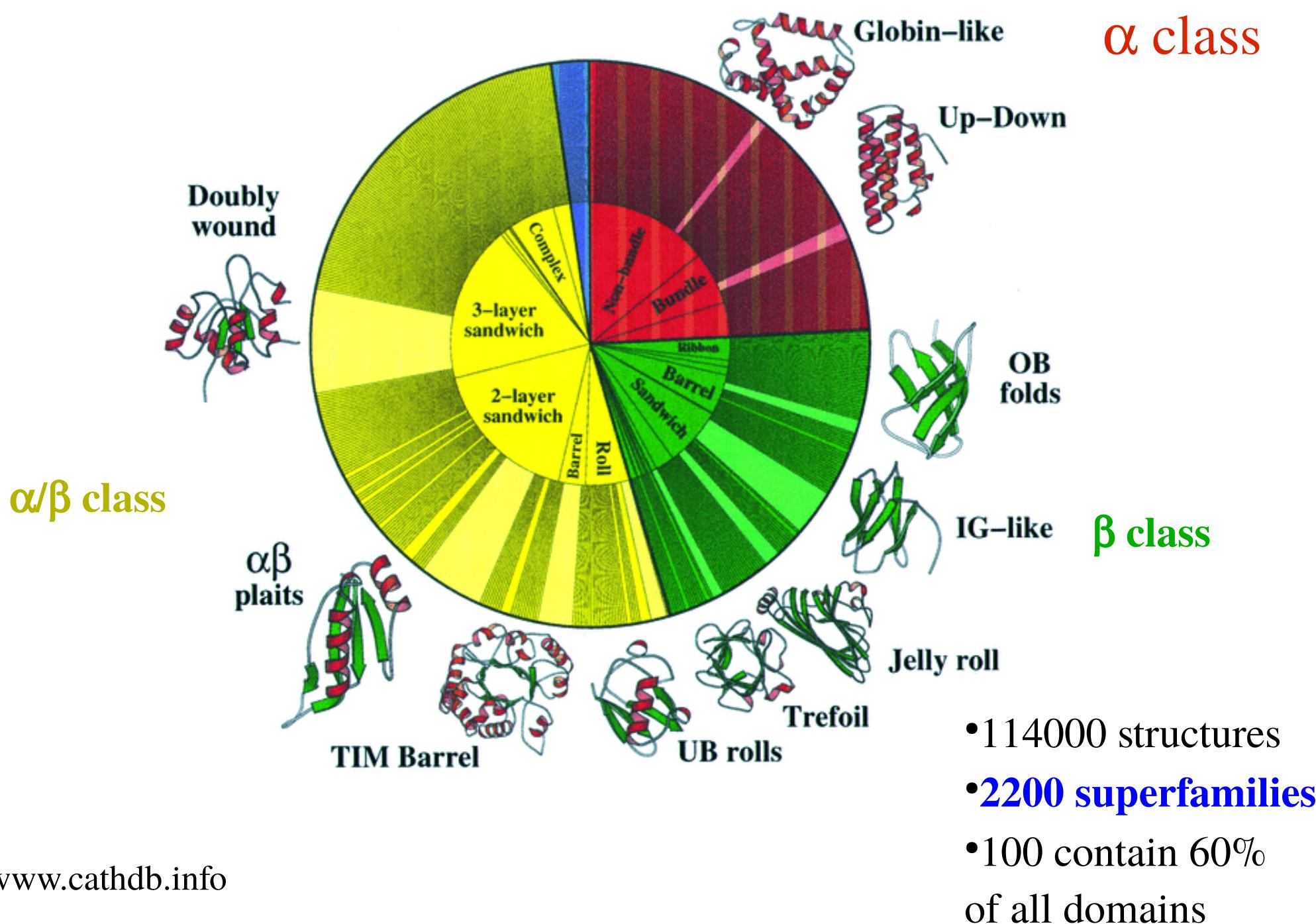


Distribution of
sequence identities
in Protein Data Bank

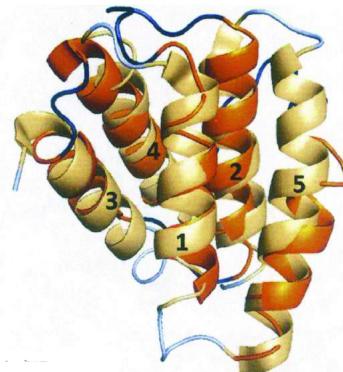
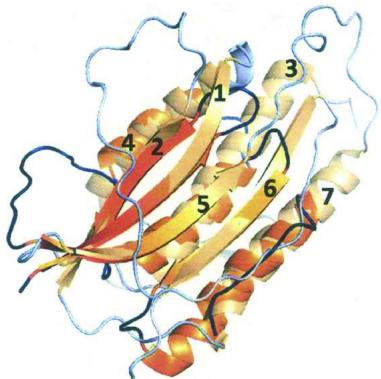
A template should have > 25% sequence
identity with the modelled protein



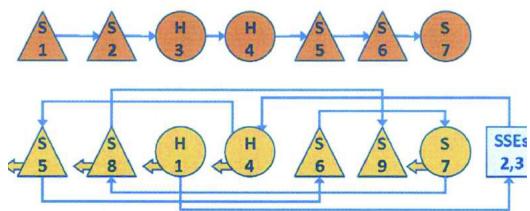
Fold set is discrete and rather small



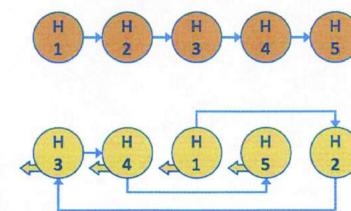
“New” folds aren’t always new....



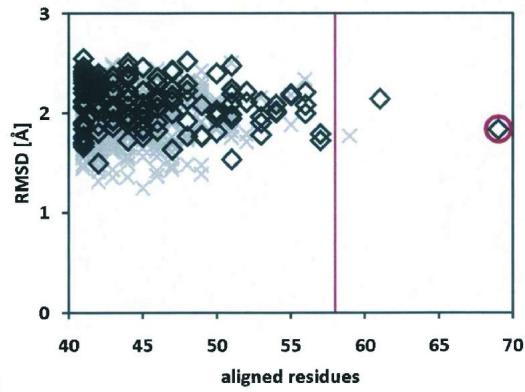
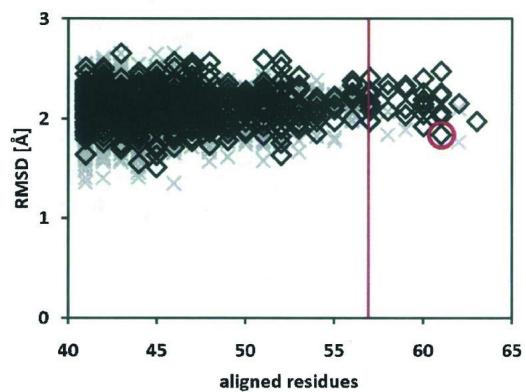
New fold



Similar fold



Alignment
with known
domains



“topologic”
view

Homology of protein:protein interactions: 35% threshold

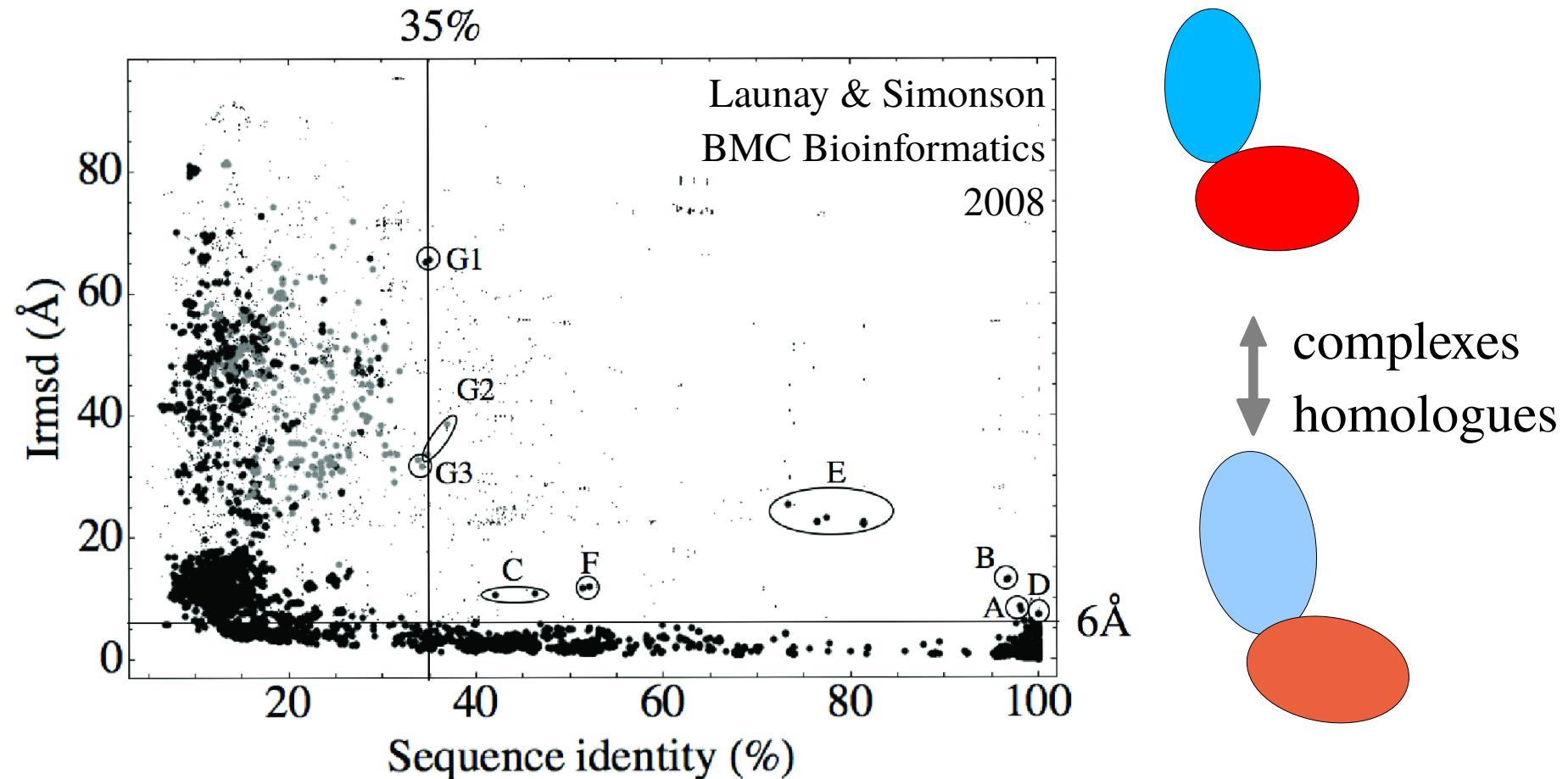
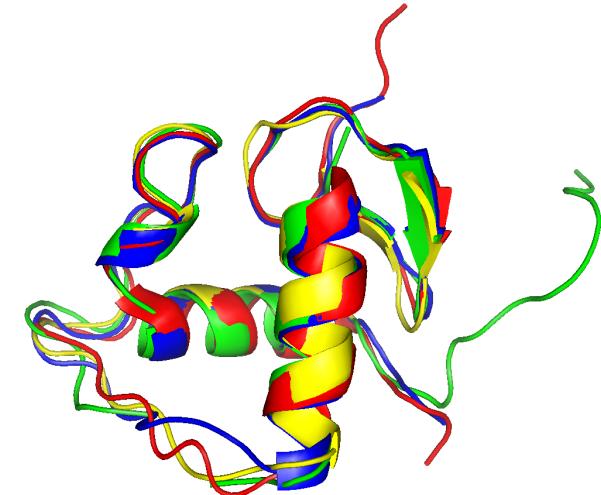


Figure 1

The relationship between sequence and structural similarity. 743 complexes from 40 interacting superfamily groups (ISGs) were analyzed. All pairs within each ISG were compared, for a total of 9630 pairwise comparisons. Small points correspond to comparisons involving at least one complex with either a small interface (buried area $< 600 \text{ \AA}^2$) or a weak association energy ($E_{int} > -10 \text{ kcal/mol}$; see text). Points labelled A-G are discussed in the text. The horizontal line corresponds to $I_{rmsd} = 6 \text{ \AA}$; the vertical line corresponds to a 35% sequence identity. Gray points correspond to comparisons where the MATRAS structural alignment provided fewer than 80% of the equivalent residues used for the I_{rmsd} calculation. All the gray points lie below the 35% similarity threshold.

Perform “Homology modelling”

Identify and align homologues of known structure (“templates”)



Do multiple sequence alignment with many homologues

Conserved, well-structured regions: adopt mean backbone of templates

Loops, insertions:

- search Protein Data Bank for similar fragments
- if not, use “ab initio” loop modelling methods

Sidechains: explore and score possible rotamers, possibly with a stochastic approach such as “Monte Carlo”

Model refinement: energy minimization, molecular dynamics

Experimental testing

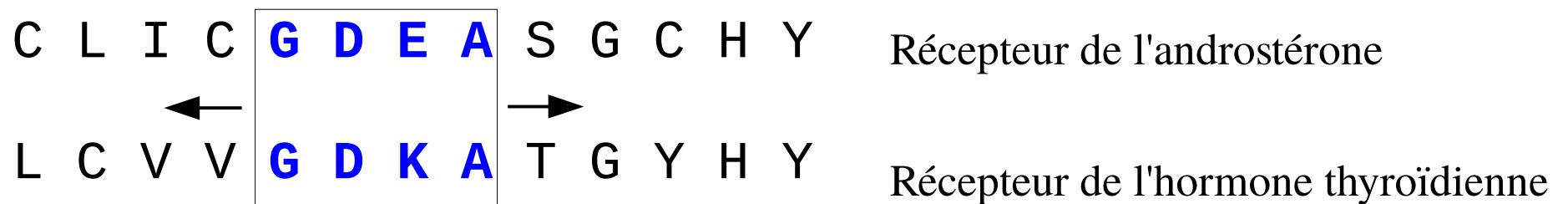
Homologues of androsterone receptor with BLAST

#	Swissprot	Hit	Description	Score (bits)	% Match		
					E	* Identity	Length
1	P15207	ANDR_RAT	Androgen receptor.	162	1e-40	100	73
6	P19091	ANDR_MOUSE	Androgen receptor.	162	1e-40	100	73
14	Q63449	PRGR_RAT	Progesterone receptor (PR)	136	1e-32	80	72
17	P06401	PRGR_HUMAN	Progesterone receptor (PR)	136	1e-32	80	72
21	P08235	MCR_HUMAN	Mineralocorticoid receptor (MR)	136	1e-32	79	72
33	P04150	GCR_HUMAN	Glucocorticoid receptor (GR)	131	3e-31	77	72
41	Q9YH32	ESR2_ORENI	Estrogen receptor beta (ER-beta)	99	3e-21	58	72
42	Q9YH33	ESR1_ORENI	Estrogen receptor (ER-alpha)	98	4e-21	55	72
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
343	Q9N4Q7	NH13_CAEEL	Nuclear hormone receptor nhr-13	54	8e-08	39	66
344	Q23294	NH11_CAEEL	Nuclear hormone receptor nhr-11	54	8e-08	42	66
345	Q45460	NH54_CAEEL	Nuclear hormone receptor nhr-54	54	1e-07	37	67
346	Q09565	NH20_CAEEL	Nuclear hormone receptor nhr-20	51	7e-07	34	66
347	Q09587	NH22_CAEEL	Nuclear hormone receptor nhr-22	45	5e-05	32	66
349	P17672	E75B_DROME	Ecdysone-induced protein 75B	40	0.001	37	47
351	P20659	TRX_DROME	Trithorax protein.	31	0.74	26	49
355	P98164	LRP2_HUMAN	Lipoprotein receptor.	30	1.7	27	65

*E = expectation of number of random alignments with a higher score

BLAST: *Basic Local Alignment Search Tool*

- Find homologous tetrapeptides
- Extend as long as similarity > threshold



Blast output

Query = Human androgen receptor
Hit = Human Estrogen receptor beta

28	CLICGDEASGCHYGALT CGSCKVFFKRAAE GKQKYLCASRN DCTIDKFRRKNCPSCRLRK	87	Query
149	CACSDYASGYHYGVWSCEGCKAFFKRSIQGHNDYICPATNQCTIDKNRRKSCQACRLRK	208	Q92731-6
88	CYEAGMTL GARKLKKLGNLKLQEEGEASSTTSPTEETTQ-----KLTSHIEGYECQPI	141	Query
209	CYE GM + ++ G ++ + A + + ++ ++ + +	268	Q92731-6
142	FLNVLEAIEPGVVVCAGHDNNQPDSFAALLSSNLGERQLVHVVKWAKALPGFRNLHVDD	201	Query
269	L +LEA P V+ + + P + A+++ SL +L +++LVH++ WAK +PGF L + D	326	Q92731-6
202	VLTLEAEPPHV LIS--RPSAPFTEASMMMSLTKLADKELVHMISWAKKIPGFVELSLFD	255	Query
327	QMAVIQYSWMGLMVFAMGWRSFTNVNSRMLYFAPDLVFNEYRMHKSRMYSQCVR-----	377	Q92731-6
256	Q+ +++ WM +++ + WRS + L FAPDLV + R +CV	283	Query
378	QVRLLESCWMEVLMMGLMWRSIDHPGK--LIFAPDLVLD-----RDEGKCVEGILEIF	407	Q92731-6

Similarité
en gris

Job information

Query sequence

>P10275

etc

Date of job execution Dec 1, 2010

Running time 23.3 seconds

Program blastp (BLASTP 2.2.23 [Feb-03-2010])

Database uniprotkb (Protein) generated for BLAST on Nov 2, 2010

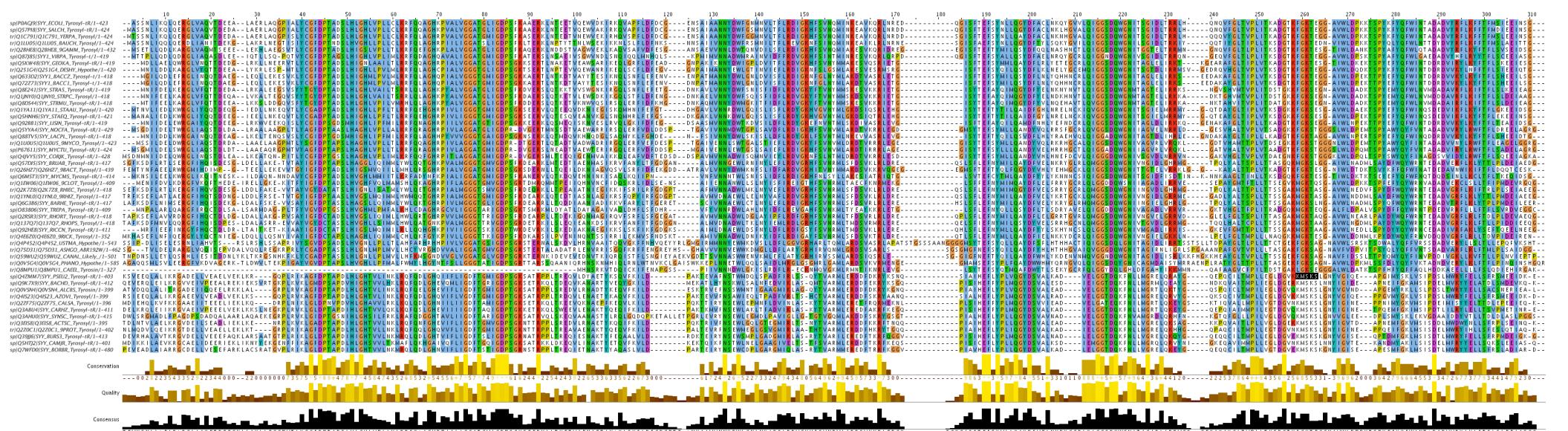
Sequences 12,898,884 sequences consisting of 4,176,319,342 letters

Matrix blosum62 Threshold 0.001 Gapped true

Heuristic alignment method: 3 phases

- " Align sequence pairs
- " Construct “guide” tree
- " Progressive sequence-profile alignment

Alignment of 30 tyrosyl-tRNA synthetases, catalytic domain



Alignment libraries

<http://pfam.sanger.ac.uk> Pfam 24.0 (October 2009, 11912 families)



HOME | SEARCH | BROWSE | FTP | HELP | ABOUT



Family: zf-C4 (PF00105)

51 architectures 3525 sequences 1 interaction 401 species 96 structures

Summary
Domain organisation
Alignments
HMM logo
Trees
Curation & models
Species
Interactions
Structures
Jump to...
<input type="button" value="enter ID/acc"/> <input type="button" value="Go"/>

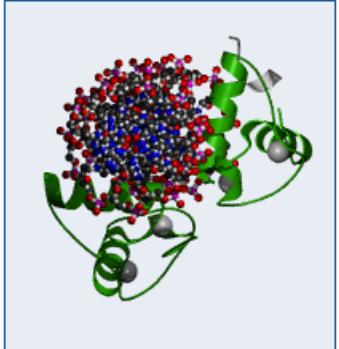
Summary

Zinc finger, C4 type (two domains) [Add annotation](#)

In nearly all cases, this is the DNA binding domain of a nuclear hormone receptor. The alignment contains two Zinc finger domains that are too dissimilar to be aligned with each other.

InterPro entry IPR001628

Steroid or nuclear hormone receptors constitute an important superfamily of transcription regulators that are involved in widely diverse physiological functions, including control of embryonic development, cell differentiation and homeostasis. The receptors function as dimeric molecules in nuclei to regulate the transcription of target genes in a ligand-responsive manner. Nuclear hormone receptors consist of a highly conserved DNA-binding domain that recognises specific sequences, connected via a linker region to a C-terminal ligand-binding domain (). In addition, certain nuclear hormone receptors have an N-terminal modulatory domain (). The DNA-binding domain can elicit either an activating or repressing effect by binding to specific regions of the DNA known as hormone-response elements [PUBMED:15242341](#), [PUBMED:15242339](#). These response elements position the receptors, and the complexes recruited by them, close to the genes of which transcription is affected. The DNA-binding domains of nuclear receptors consist of two zinc-nucleated modules and a C-terminal extension, where residues in the first zinc module determine the specificity of the DNA recognition and residues in the second zinc module are involved in dimerisation. The DNA-binding domain is furthermore involved in several other functions including nuclear localisation, and interaction with transcription factors and co-activators [PUBMED:15242320](#).



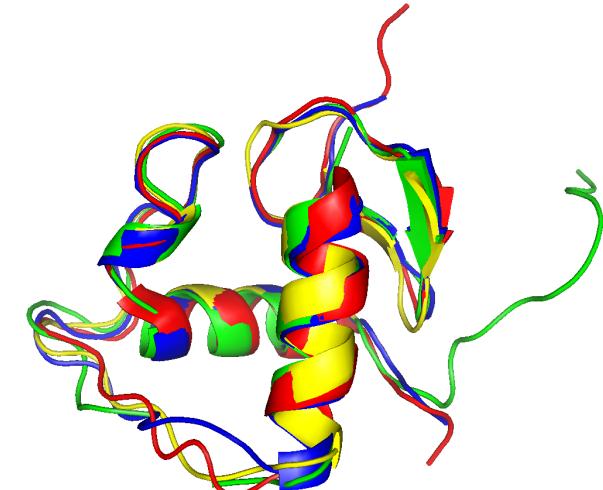
Example structure
[PDB entry 1r0o](#): Crystal Structure of the Heterodimeric Ecdysone Receptor DNA-binding Complex
[View a different structure](#): 1r0o | ▲

Famille zf-C4 → “expert” alignment “expert” of 26 sequences; “automatic” of 3525 séquences

Balibase: 217 expert alignments benchmark set to tester methods/programs

<http://www-bio3d-igbmc.u-strasbg.fr/balibase> Thompson et al (1999) Bioinformatics, 57, 87-88

Perform “Homology modelling”



Identify and align homologues of known structure (“templates”)

Do multiple sequence alignment with many homologues

Structural alignment with available homologues

Loops, insertions:

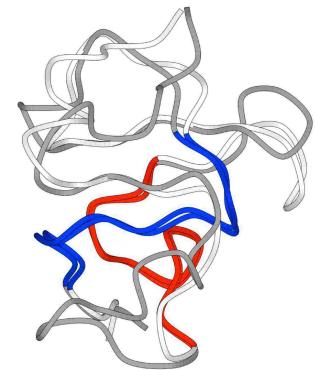
- search Protein Data Bank for similar fragments
- if not, use “ab initio” loop modelling methods

Sidechains: explore and score possible rotamers, possibly with a stochastic approach such as “Monte Carlo”

Model refinement: energy minimization, molecular dynamics

Experimental testing

Remarks on structural alignments

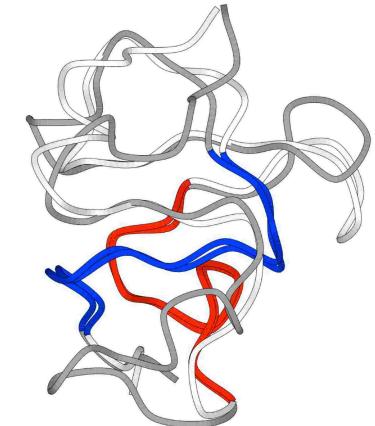


- 1) Differ from sequence alignments: an amino acid is not coded by a letter
- 2) A “trivial” method is possible: align sequences then superimpose homologous C_α atoms
- 3) We can align complex objects, like alignments... Code an amino acid by a vector of properties: type, secondary structure, burial, ...
- 4) In general, various heuristic approaches are used.
- 5) For multiple alignments, use a progressive method.

Remarks on structural alignments

STRUCTAL: **a)** align sequences; **b)** superimpose homologous C α atoms;
c) compute scoring matrix; **d)** repeat a-c.

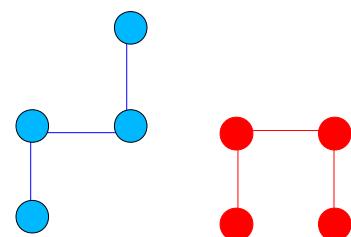
Score M(i,j) = -distance between i, j (Levitt, 1993)



SSAP: describe each amino acid by a vector of properties;
Align with Needleman-Wunsch (Taylor & Orengo, 1989)

CE: find two similar octapeptide; extend with gaps. (Bourne, 1998)

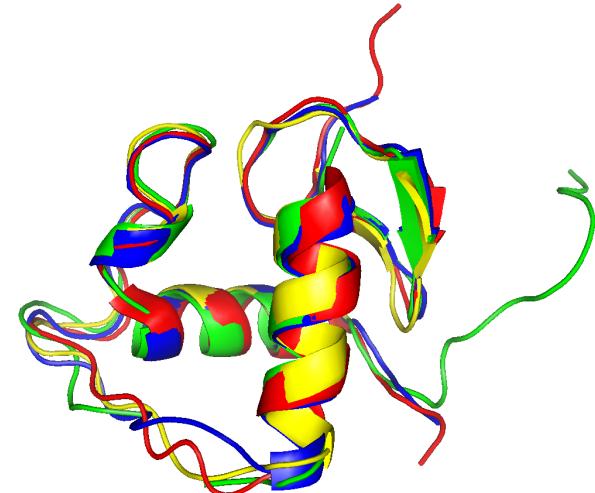
DALI: similar to CE (Holm & Sander, 1993).



Large-scale test: Kolodny et al (2005) J Mol Biol, 346:1173

Over 80 methods in Wikipedia...

Perform “Homology modelling”



Identify and align homologues of known structure (“templates”)

Do multiple sequence alignment with many homologues

Structural alignment with available homologues

Loops modeling methods

Side chains: explore and score possible rotamers

Model refinement: energy minimization, molecular dynamics

Experimental testing

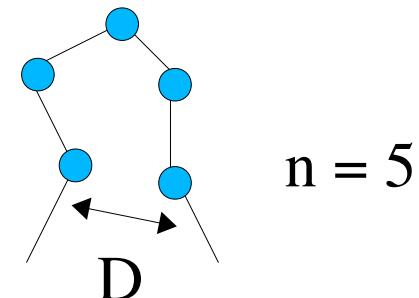
Loop modeling with an experimental conformation library

Wojcik et al (1999) J Mol Biol, 289:1469

13,563 loops of 3-8 amino acids from the PDB

Ranked by:

- length n
- distance D between ends
- "height" and "width"



n = 5

To model a new loop X of length n:

- estimate D
- find loops in library compatible with n, D;
- find those with sequence similar to X
- refine with molecular dynamics

What are the main hypotheses of the method?

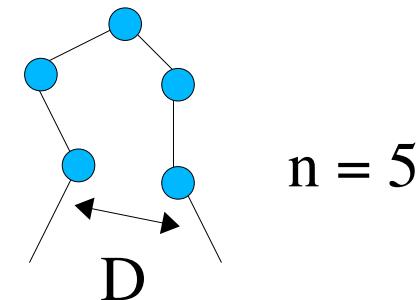
Loop modeling with an experimental conformation library

Wojcik et al (1999) J Mol Biol, 289:1469

13,563 loops of 3-8 amino acids from the PDB

Ranked by:

- length n
- distance D between ends
- "height" and "width"



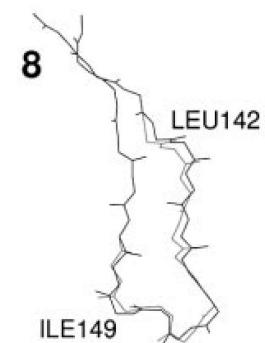
n = 5

To model a new loop X of length n:

- estimate D
- find loops in library compatible with n, D;
- find those with sequence similar to X
- refine with molecular dynamics

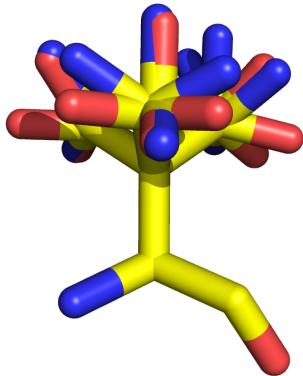
Quality of
results

	Longueur de boucle	3	4	5	6	7	8
	Nombre de cas testés	92	81	54	61	37	37
	Rmsd des C _α (A)	1.0	1.2	1.6	2.1	2.7	2.3
	Avec Modeller:			0.6		1.2	



Side chain placement by simulated annealing

Side chains explore rotamers.

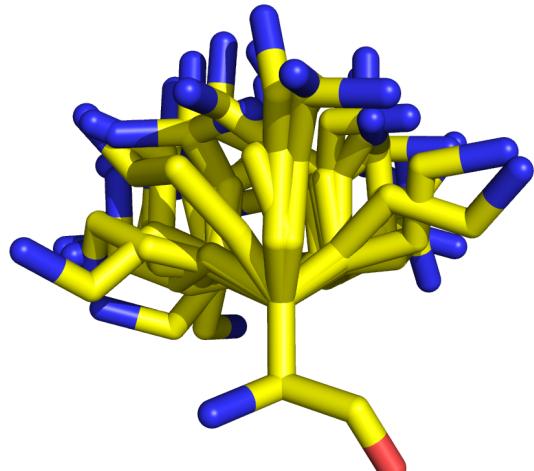


Explore with Monte Carlo:

Many elementary steps (10^6).

At each step:

- pick an amino acid randomly
- pick a new rotamer randomly
- compare new energy E_n to previous energy E_p
- if $E_n < E_p$, keep the new rotamer
- if $E_n > E_p$, pick a random number between 0 and 1
 - if $\exp[-(E_n - E_p)/RT] > x$, keep the new rotamer
 - else keep the previous one



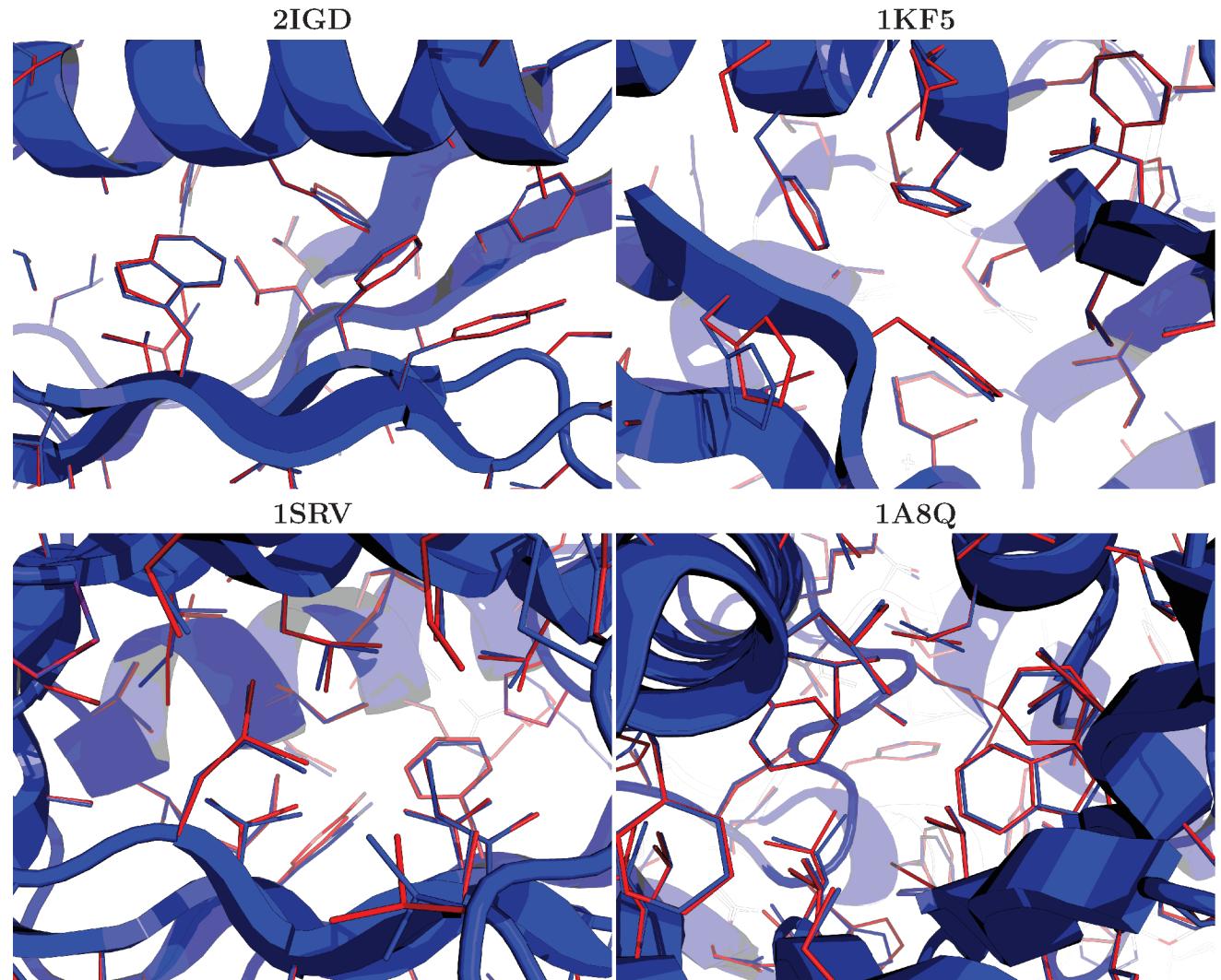
Start at high temperature T and lower T progressively.

Remarkably, the method explores conformations with probability $p(\text{Conf}) = A \exp(-E_{\text{Conf}}/RT)$, the experimental, Boltzmann distribution.

Prediction exercise: delete, then reconstruct side chains; main chain known exactly.

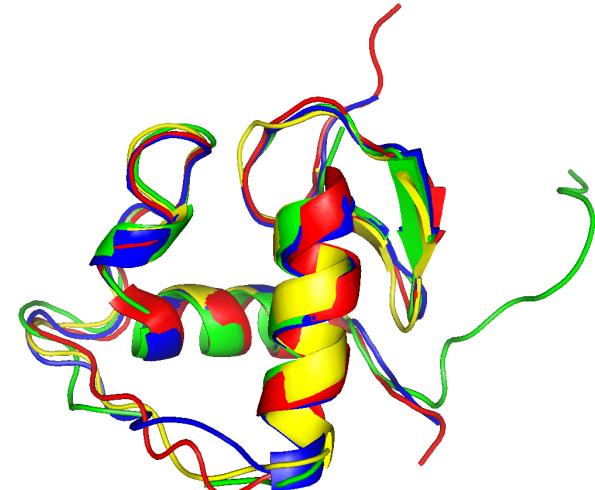
86% chi1 correct
73% chi1 et chi2

With perfect
backbone, side chain
predictions are good
but not perfect.



Gaillard, Panel, Simonson (2016) Proteins

Perform “Homology modelling”



Identify and align homologues of known structure (“templates”)

Do multiple sequence alignment with many homologues

Structural alignment with available homologues

Loops modeling methods

Side chains: explore and score possible rotamers

Model refinement: energy minimization, molecular dynamics

Experimental testing

Performance of selected methods

	% correct predictions	
	chi1	chi2+chi1
Gaillard et al, 2016	86%	73%
Koehl & Delarue, 1994	72%	62%
Yang et al, 2002	80%	66%
Dunbrack et al, 1999	80%	?
Abagyan et al, 1993	80%	67%
Mendes et al, 1999	87%	(improved rotamer library)

In these tests, the side chains were placed on the known backbone...

In homology modeling, the backbone is not exact...

Homology modeling. Krieger E, Nabuurs SB, Vriend G.
In Structural Bioinformatics; editors PE Bourne, H Weissig; Wiley, 2003

Homology modelling in biology and medecine. R Dunbrack.
In Bioinformatics: from genomes to drugs; editor T Lengauer; Wiley, 2002; Vol. 1.

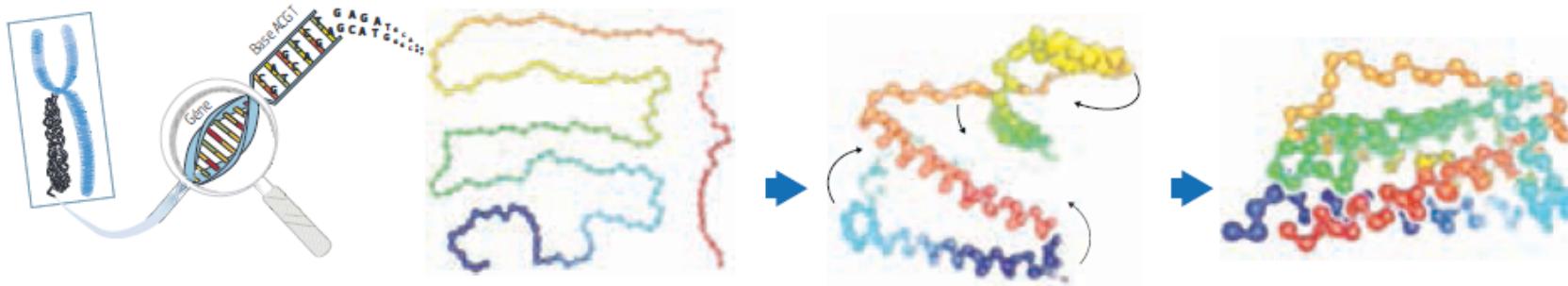
Structure comparison and alignment. Bourne & Shindyalov
In Structural Bioinformatics; editors PE Bourne, H Weissig; Wiley, 2003

Proteins: Structure, Function, and Bioinformatics; Volume 77, Issue S9, pages 1-228.
Numéro spécial décrivant la compétition CASP8:

“Critical Assessment of Structure Prediction”
64 cibles à prédire; 159 méthodes/équipes

Le secret du pliage des protéines dévoilé

Une intelligence artificielle de Google arrive à prédire la structure 3D d'à peu près n'importe quelle molécule codée par le génome humain. Une avancée majeure qui ouvre un champ des possibles considérable.



① Chaque chromosome contient une multitude de bases (A, C, G ou T). Combinées entre elles, ces bases forment le **code génétique**.

② Le code génétique permet la constitution de longues **chaînes d'acides aminés**...

③ ... qui se replient naturellement sur elles-mêmes pour constituer la **protéine** et lui donner sa **forme spécifique**.

④ Grâce à l'intelligence artificielle Alpha Fold, les chercheurs peuvent désormais **décoder ce « schéma de pliage »** et donc prédire la structure 3D de la protéine.

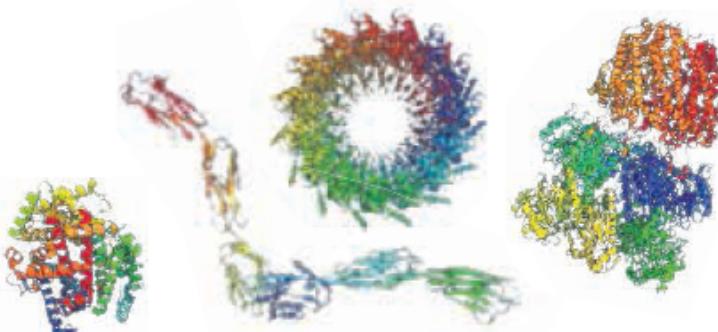
Infographie LE FIGARO

STEPHANY GARDIER @S_Gardier

BIOLOGIE «Incrroyable», «révolutionnaire», «bluffant»... Quand ils évoquent la plateforme AlphaFold, dévoilée il y a quelques jours par DeepMind (une filiale d'Alphabet, maison mère de Google), les spécialistes ne manquent pas de superlatifs. Cet outil d'intelligence artificielle est capable de prédire la structure de la quasi-totalité des protéines codées par le génome humain. De quoi ouvrir de nouvelles perspectives notamment en médecine, la plupart des produits thérapeutiques ayant pour cible des protéines. Après la publication du génome humain en 2003, la mise à disposition de la quasi-totalité de ce «protéome» humain pourrait marquer un nouveau tournant pour la recherche biomédicale.

Une grande diversité de structures possibles

Une fois repliées sur elles-mêmes, les chaînes d'acides aminés peuvent former toutes sortes de **protéines**, qui constituent la base du vivant.



«Pouvoir prédire la structure d'une protéine offre un gain de temps considérable»

THOMAS SIMONSON, DIRECTEUR DE RECHERCHE SPÉCIALISÉ EN BIOLOGIE DES PROTÉINES À L'ÉCOLE POLYTECHNIQUE (PALAISEAU)

cellules grâce à une machinerie (elle-même faite de protéines) capable de traduire le code génétique en acides aminés. Ces petites briques élémentaires sont assemblées de manière linéaire un peu comme les perles d'un collier. Sous l'effet de forces qui s'exercent entre ces perles, la protéine va se replier pour former une structure en trois dimensions. Et c'est le résultat de cet origami de précision qui détermine les fonctionnalités de la protéine, dont la forme globale joue un rôle crucial dans son fonctionnement. «Pouvoir prédire la structure d'une protéine offre un gain de temps considérable, en permettant de cibler les domaines qui ont un intérêt particulier et sur lesquels il faut réaliser des vérifications expérimentales», souligne Thomas Simonson, directeur de re-

cherche spécialisé en biologie des protéines à l'École polytechnique (Palaiseau), qui imagine déjà comment AlphaFold va accélérer le rythme des travaux de son laboratoire.

Rapidité, simplicité et fiabilité sont trois des atouts d'AlphaFold, qui a supplanté tous ses concurrents lors de la dernière édition du CASP, la compétition internationale biseannuelle de prédiction de structure de protéines. «Ce qui est encore plus surprenant, c'est que DeepMind ne s'est attaqué à cette thématique du repliement des protéines que récemment», souligne Thomas Simonson. «On les a vus la première fois à la 13^e édition en 2019, et deux ans plus tard, ils sont les premiers, loin devant les autres!» Développés majoritairement par des laboratoires académiques, les autres programmes n'ont pas fait le poids face à la puissance de calcul et aux investissements du nouveau venu. Interrogé sur le coût de ce projet, Demis Hassabis, le fondateur de DeepMind, n'a pas donné de chiffres, mais a souligné les énormes efforts consentis pour ce qu'il a présenté comme «le plus gros investissement» de sa société, qui a monopolisé une équipe de 30 personnes durant cinq ans. «Il est évident qu'aucun laboratoire académique ne peut déclencher autant de ressources

pour tenter de coder un programme», commente Cécile Breyton.

Comme sa consoeur, Alexis Verger, chercheur à l'Institut des sciences biologiques du CNRS, salue la décision de DeepMind d'avoir opté pour un accès libre à AlphaFold et d'avoir publié le code source du programme. «Mais n'oublions pas que, sans le travail de tous les biologistes structuraux des cinquante dernières années, AlphaFold n'existerait pas, car ce sont leurs résultats qui ont servi à "entraîner" l'intelligence artificielle du programme. Et c'est grâce à nos futures données expérimentales que l'IA d'AlphaFold deviendra encore plus performante. C'est un excellent exemple du cercle vertueux des sciences ouvertes.»

AlphaFold contient les structures des protéines de la bactérie de la tuberculose et de *Plasmodium falciparum*, responsable du paludisme

AlphaFold devrait intéresser de nombreux laboratoires car, en plus des protéines humaines, le banc de

Highly accurate protein structure prediction with AlphaFold

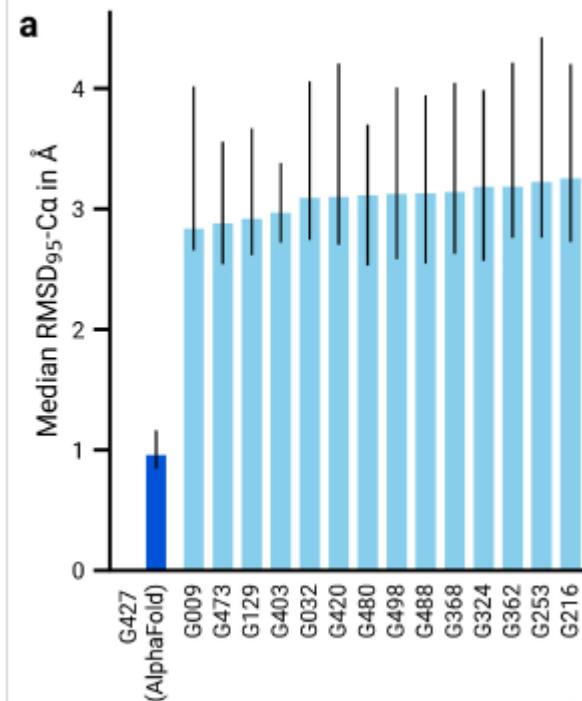
<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

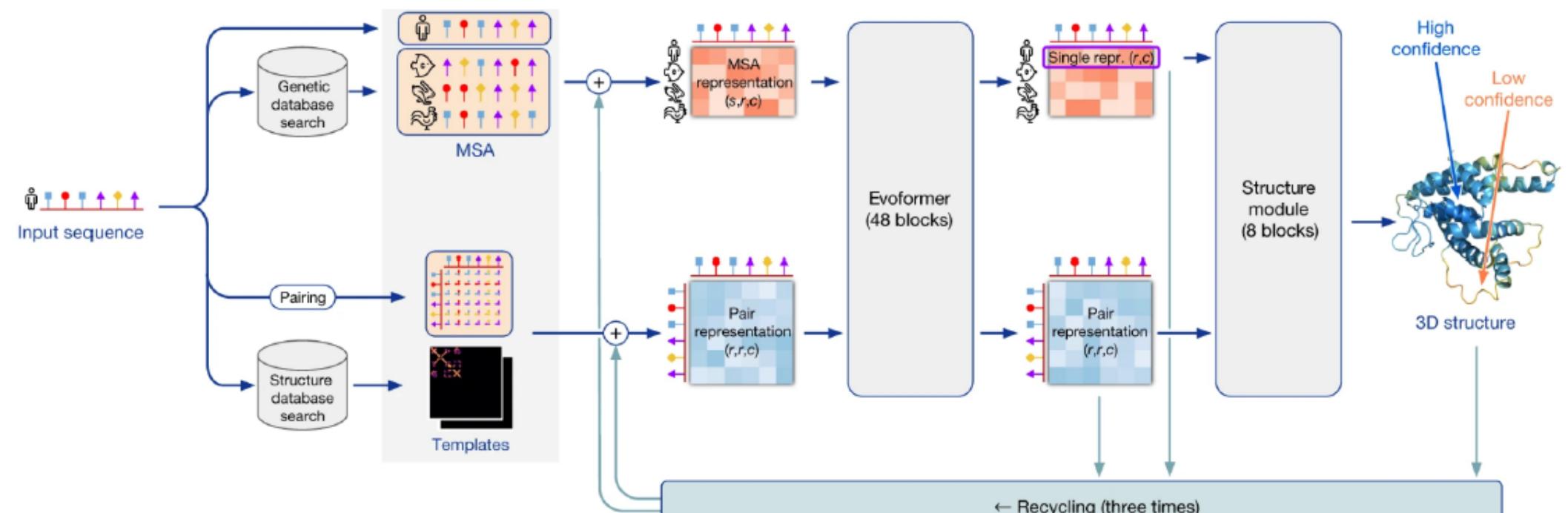
Published online: 15 July 2021

John Jumper^{1,4}, Richard Evans^{1,4}, Alexander Pritzel^{1,4}, Tim Green^{1,4}, Michael Fligurnov^{1,4}, Olaf Ronneberger^{1,4}, Kathryn Tunyasuvunakool^{1,4}, Russ Bates^{1,4}, Augustin Žídek^{1,4}, Anna Potapenko^{1,4}, Alex Bridgland^{1,4}, Clemens Meyer^{1,4}, Simon A. A. Kohl^{1,4}, Andrew J. Ballard^{1,4}, Andrew Cowle^{1,4}, Bernardino Romera-Paredes^{1,4}, Stanislav Nikolov^{1,4}, Rishabh Jain^{1,4}, Jonas Adler¹, Trevor Back¹, Stig Petersen¹, David Relman¹, Ellen Clancy¹, Michal Zelinski¹, Martin Stelmegger^{2,3}, Michałina Pacholska¹, Tamas Berghammer¹, Sebastian Bodensteiner¹, David Silver¹, Oriol Vinyals¹, Andrew W. Senior¹, Koray Kavukcuoglu¹, Pushmeet Kohli¹ & Demis Hassabis^{1,4}



Proteins are essential to life, and understanding their structure can facilitate a mechanistic understanding of their function. Through an enormous experimental effort^{1–4}, the structures of around 100,000 unique proteins have been determined⁵, but this represents a small fraction of the billions of known protein sequences^{6,7}. Structural coverage is bottlenecked by the months to years of painstaking effort required to determine a single protein structure. Accurate computational approaches are needed to address this gap and to enable large-scale structural bioinformatics. Predicting the 3-D structure that a protein will adopt solely on its amino acid sequence, the structure prediction component of the ‘protein folding problem’⁸, has been an important open research problem for more than 50 years⁹. Despite recent progress^{10–14}, existing methods fall far short of atomic accuracy, especially when no homologous structure is available. Here we provide the first computational method that can regularly predict protein structures with atomic accuracy even where no similar structure is known. We validated an entirely redesigned version of our neural network-based model, AlphaFold, in the challenging 14th Critical Assessment of protein Structure Prediction (CASP14)¹⁵, demonstrating accuracy competitive with experiment in a majority of cases and greatly outperforming other methods. Underpinning the latest version of AlphaFold is a novel machine learning approach that incorporates physical and biological knowledge about protein structure, leveraging multi-sequence alignments, into the design of the deep learning algorithm.

AlphaFold: the general structure



search in databases

- for MSA construction
- template search

construction of the initial pair representation

simultaneous optimization of the MSA representation and the pair representation thanks to a neural network

Determination of 3D position for each atom thanks to a neural network