

Introduction to computational protein design

Thomas Simonson

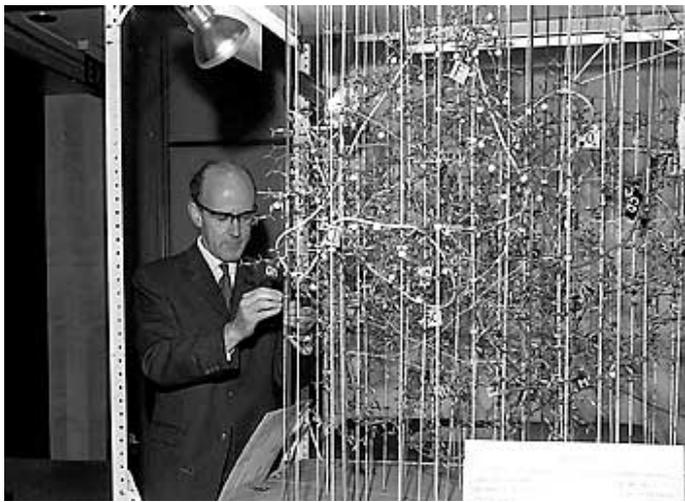
Laboratoire de Biologie Structurale de la Cellule
Ecole Polytechnique, Paris

thomas.simonson@polytechnique.fr

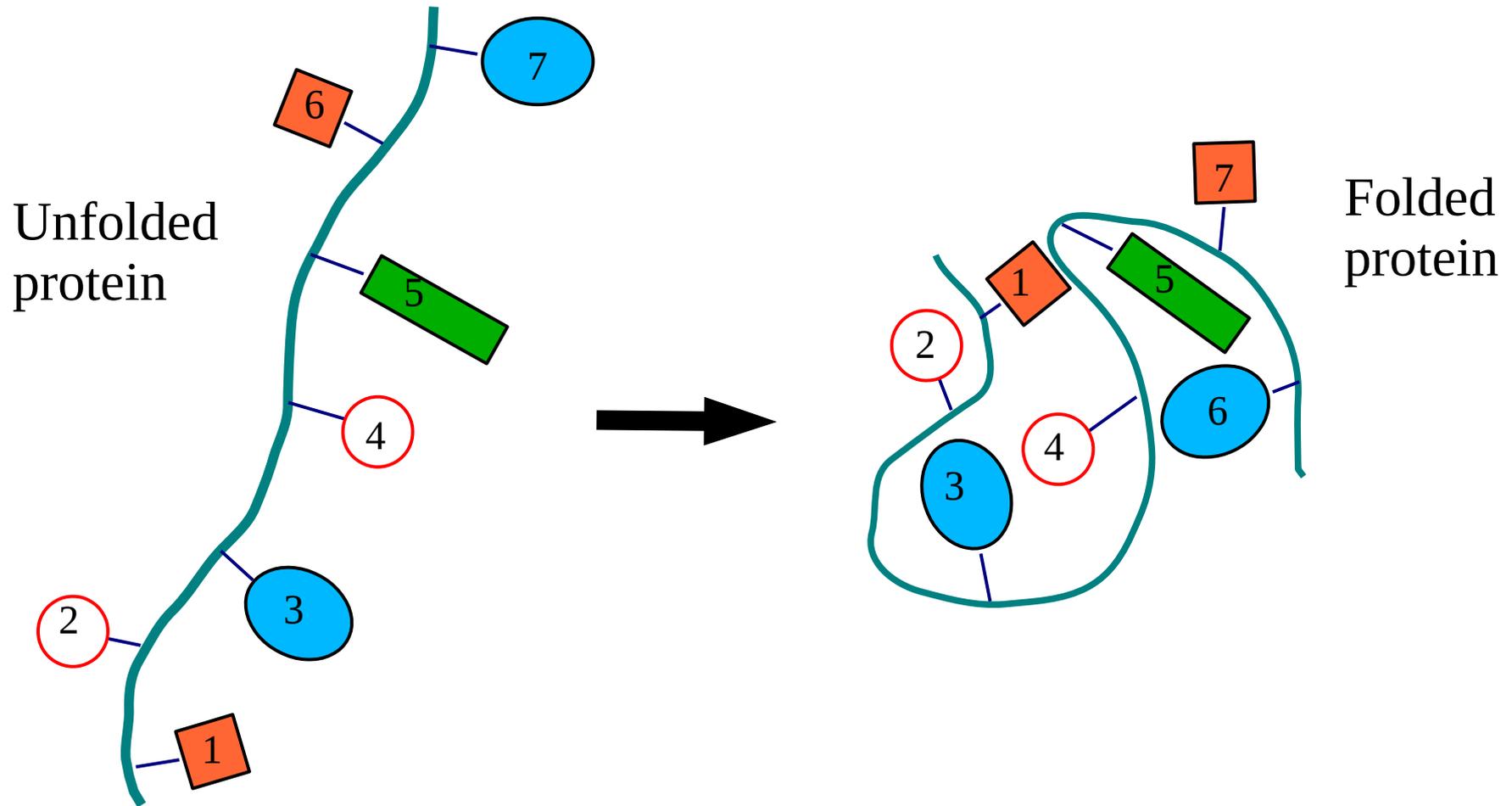
<http://biology.polytechnique.fr/biocomputing>

<https://proteus.polytechnique.fr>



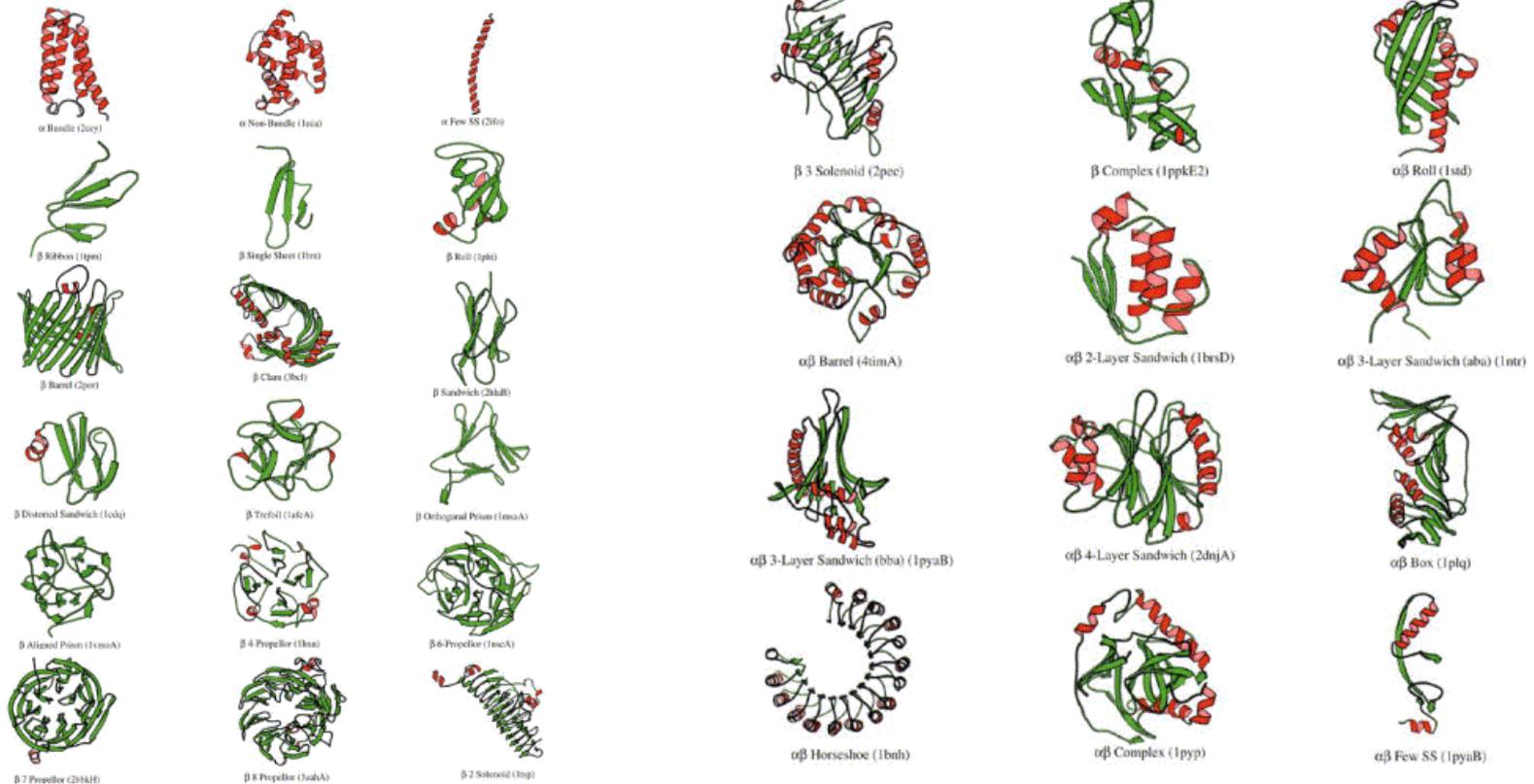


**The amino acid sequence
directs the polypeptide to
a particular «native»
structure**



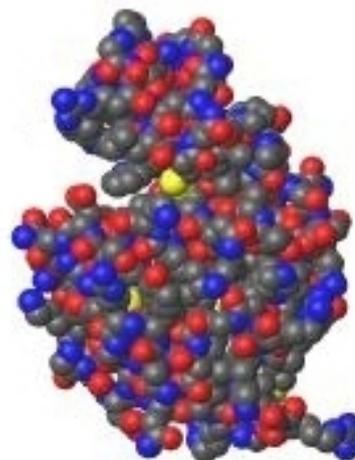


The amino acid sequence directs the polypeptide to a particular «native» structure



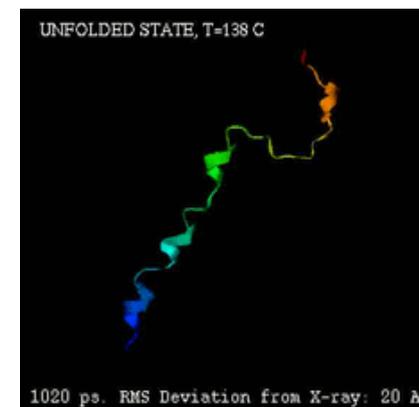
The protein folding problem: can we predict protein structures?

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



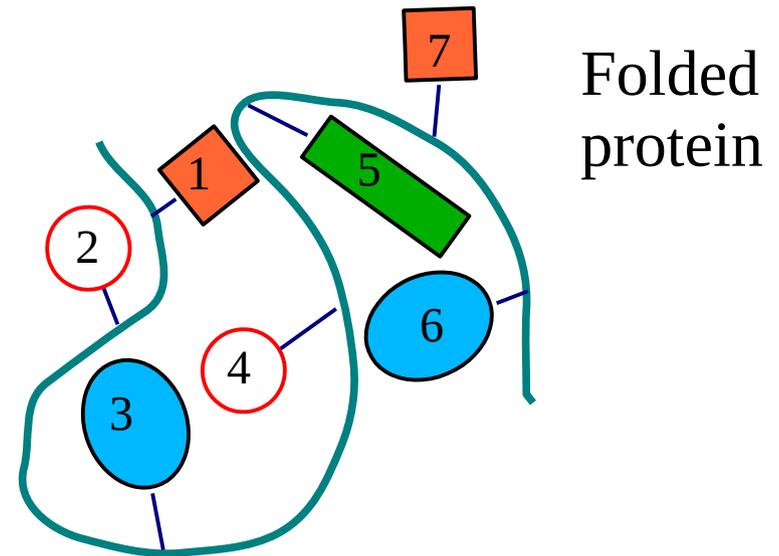
$$N_{\text{states}} \sim 10^n$$
$$n = 100-300$$

googols of possibilities...

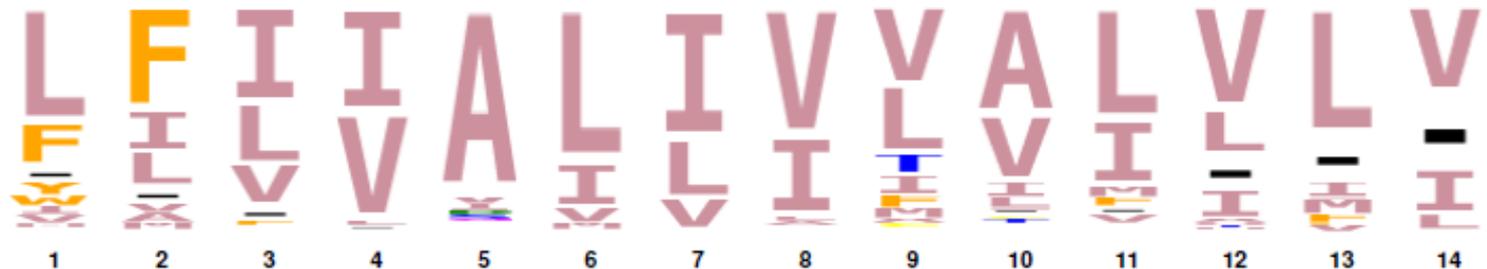


Inverse folding problem: for a given “backbone structure”, or fold, which sidechains “fit”?

The size of the problem is enormous.....
100 amino acids $\Rightarrow 10^{130}$ sequences

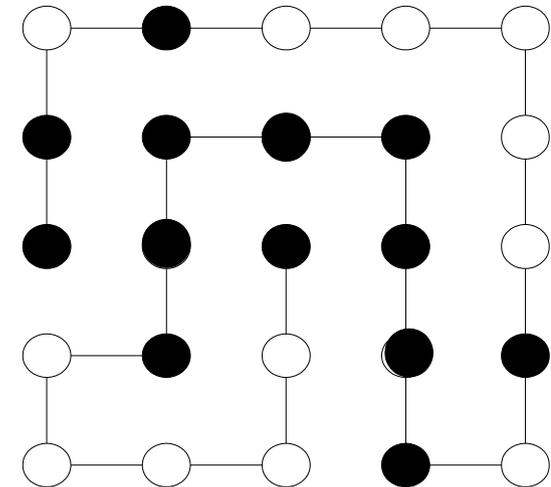
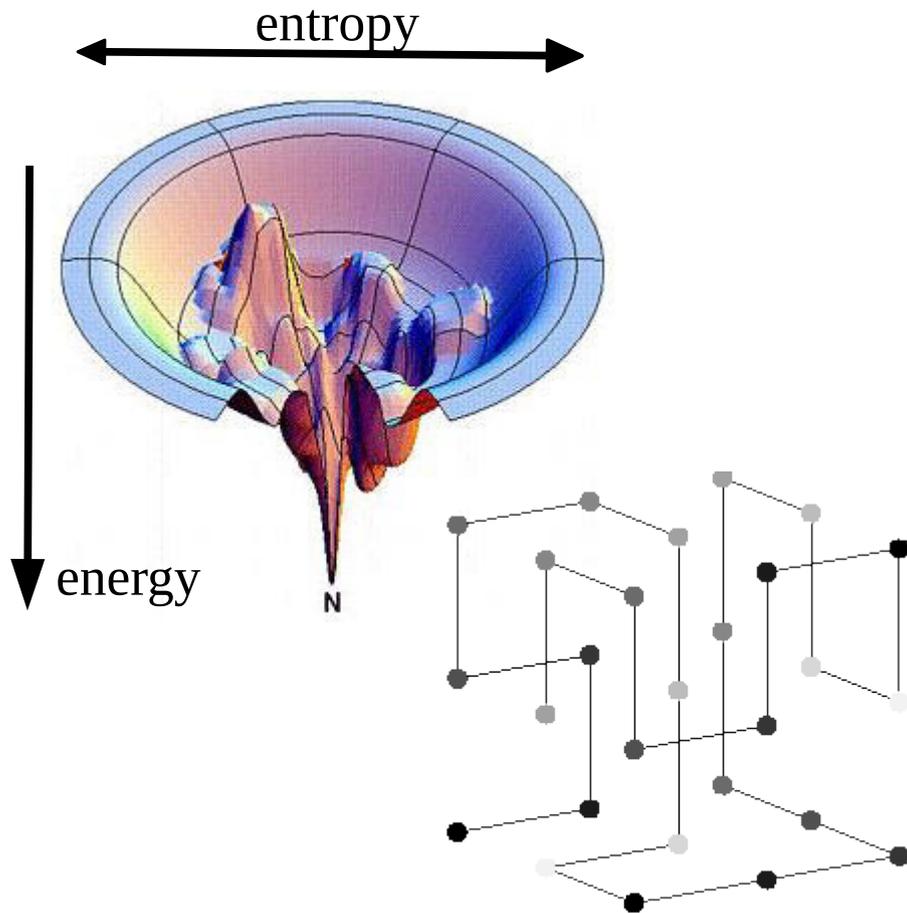


Some natural
PDZ sequences
(core positions)



>10,000 natural homologues in UniProt....

The inverse folding problem: “Easy” to solve with toy models

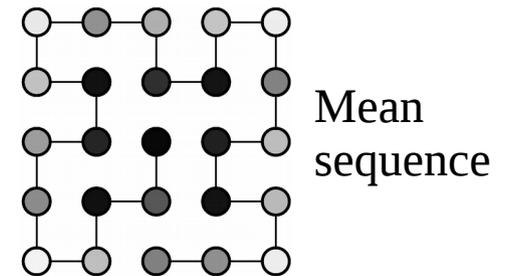
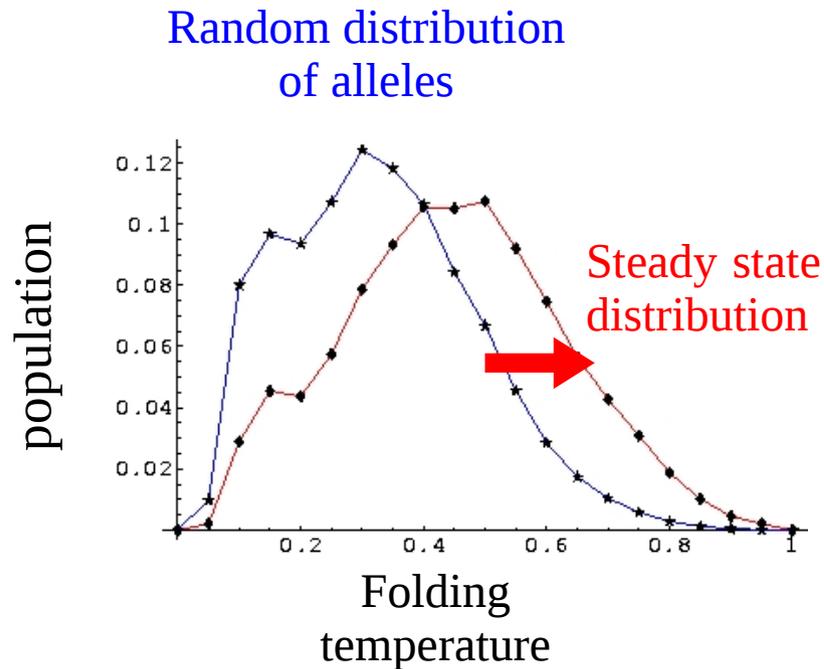


1081 conformations
33.554.432 sequences

- Designability correlates with fitness, thermophilic adaptation
- Of interest for protein engineering and biotechnology

The inverse folding problem

Easy to solve with toy models of structure and population dynamics



Noirel & TS, J Chem Phys '08

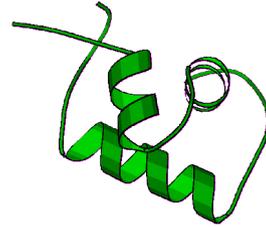
Steady state population enriched in sequences that are thermostable and tolerant of mutations (designable)

Protein design: search sequence/conformation space for preferred/functional sequences

random mutagenesis + selective pressure

Ponder, Richards (1987); Hellinga, Richards (1994); Mayo (1997) Desjarlais (1998)
Koehl, Levitt (1999); Baker, Serrano, Wodak, Handel (2000-04); others...

AHGSQNTTILIP...
DKPAIFTDLGDWV...
EKPLEVDDAAEWS...
PLIKRYWWNAQAG...
MKPVTLLDVAEYA...
GHYILKQSANCCM...
FKPIEASDIAEFV...
QKPVSLSDVGEFA...

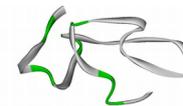
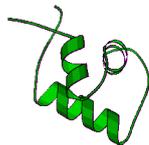


DKPAIFTDLGDWV...
EKPLEVDDAAEWS...
MKPVTLLDVAEYA...
QKPVSLSDVGEFA...

Sequence score: $\Delta G = G_{\text{folded}} - G_{\text{unfolded}}$

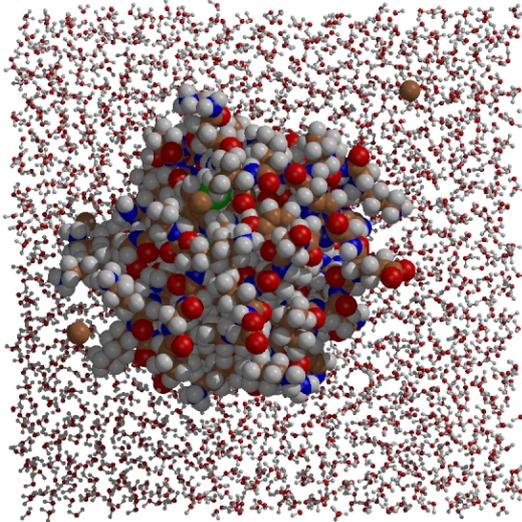
Free energy function

Structural model of both ensembles

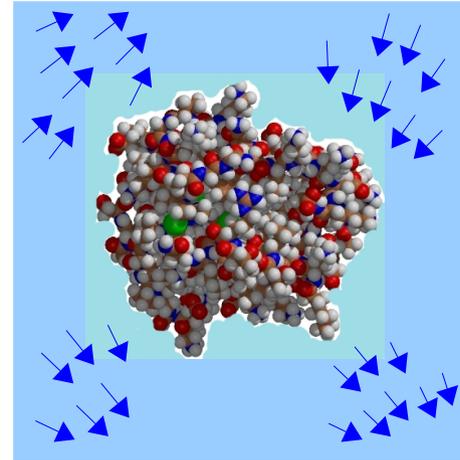


Implicit solvent: polar free energy components

Explicit solvent



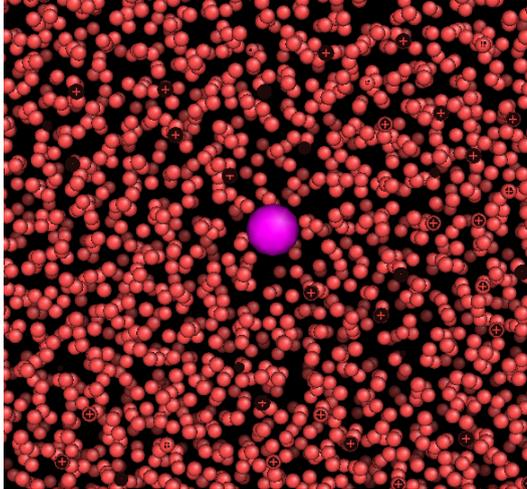
Implicit solvent



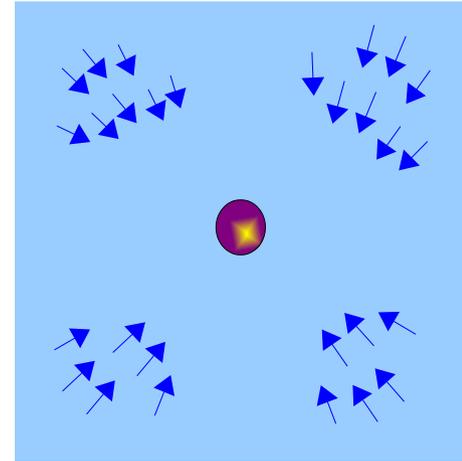
Dielectric continuum

Implicit solvent

“explicit”



“implicit”



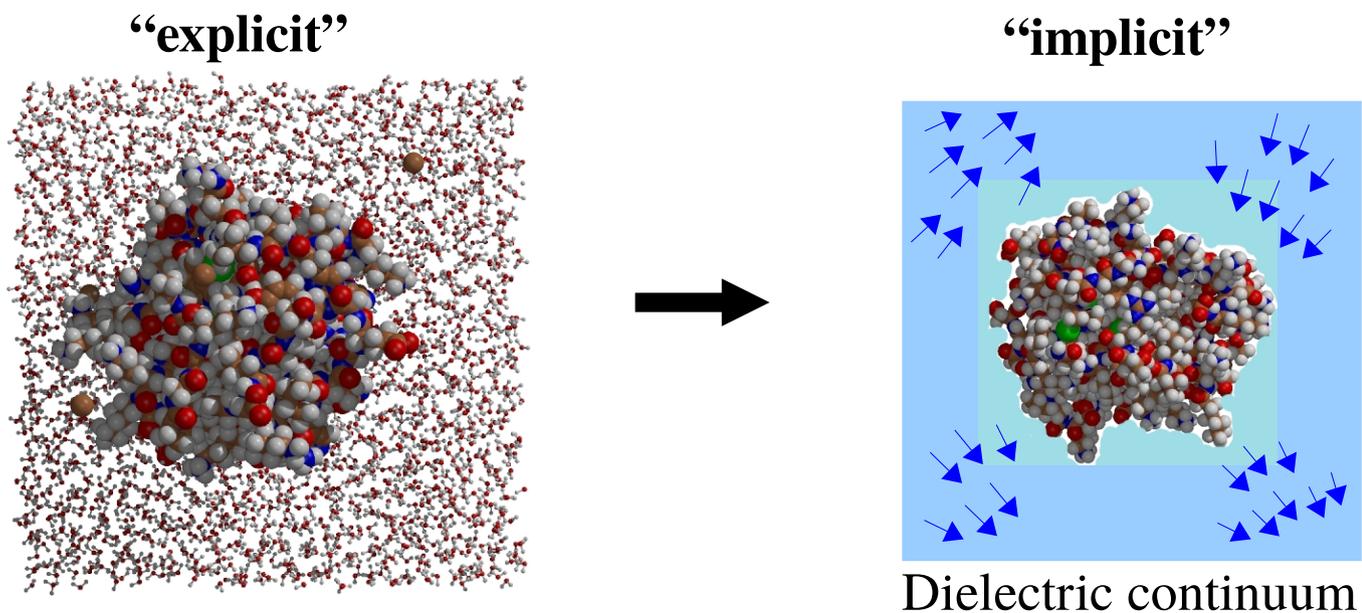
Dielectric continuum

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.....

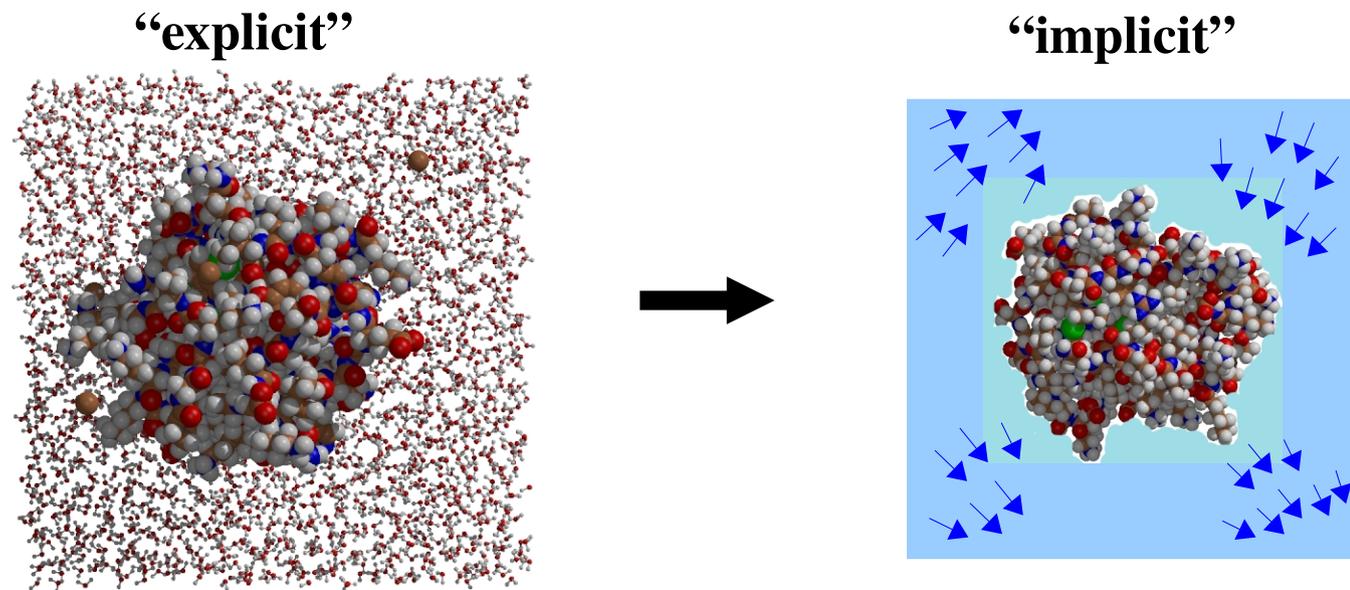
Implicit solvent



A biomolecule is more complex: heterogeneous system

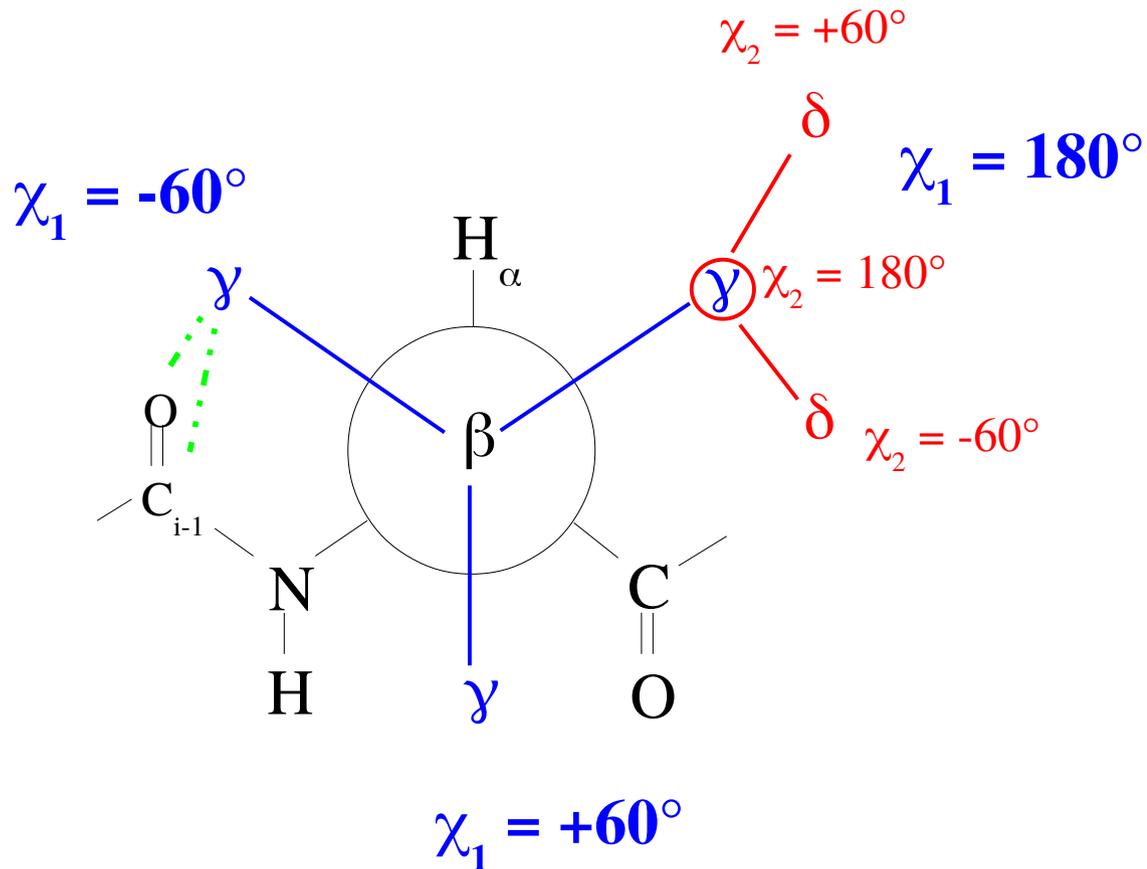
Going further: more sophisticated model

Implicit solvent: “generalized Born” model (GB)

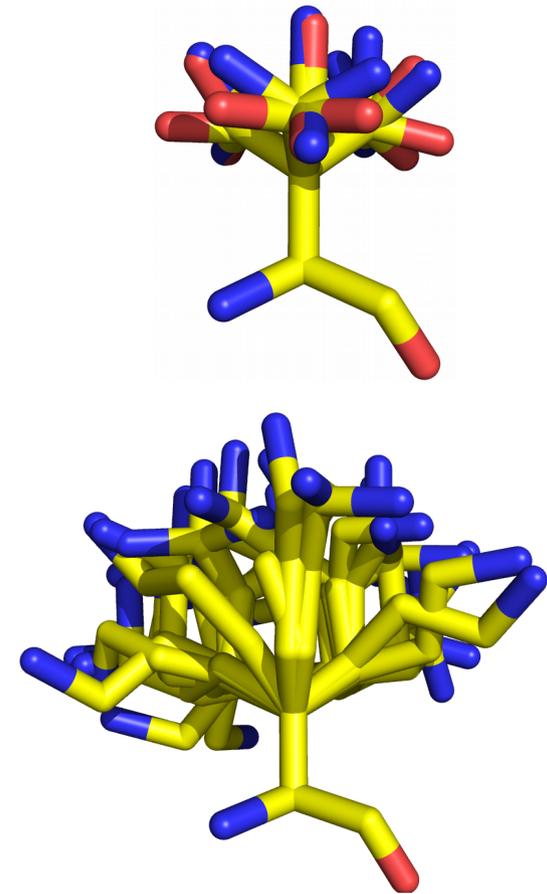


- Buried charges: $E = q q' / r$
- Exposed charges: $E = q q' / \epsilon r$
 $\epsilon = 80$
- Intermediate cases: **interpolation**

Discrete conformational space



Conformational space
=
fixed backbone
+
sidechain
“rotamers”



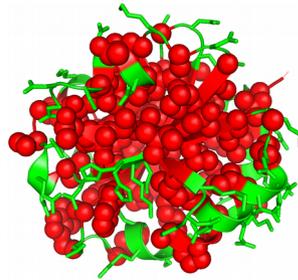


Protein structure is complex...

Side chains are dynamic and flexible...

Unfolded proteins form a molten globule that is hard to characterize...

Aqueous solvent has a rich, complex structure...



Protein structure is simple!

Side chains have just a few conformations !

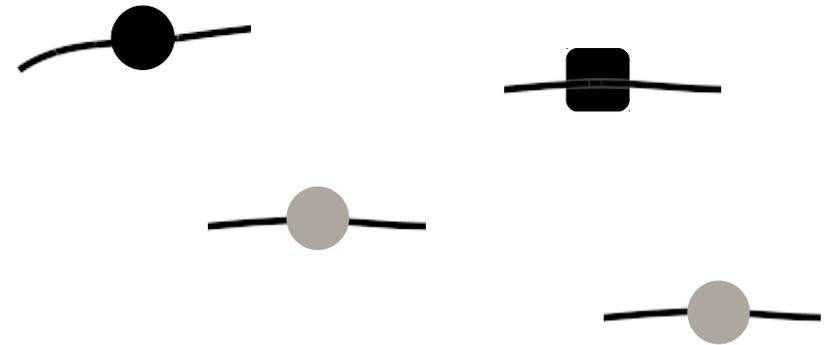
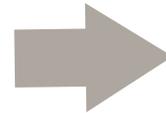
Unfolded proteins are just like an extended peptide !

Solvent is just like a simple dielectric !

The unfolded state : a simple model

extended peptide :

approximated as collection
of tripeptides :



$$E_u = \sum_i E_u(t_i)$$

- Unfolded energy only depends on composition, through type-dependent amino acid *chemical potentials* ; no structural model.
- The chemical potentials can be computed from an extended peptide model or chosen empirically (see below).
- Longer peptides (5-, 7-mers) don't help (Pokola & Handel, JMB '05)

Maximum likelihood strategy for the unfolded state parameters

Choose a set of experimental sequences \mathbf{S}

= {S, S', ...} as a reference

Boltzmann probability of sampling a

sequence is $p(S) = \exp(-\beta\Delta G)/Z$

ΔG folding free energy, Z

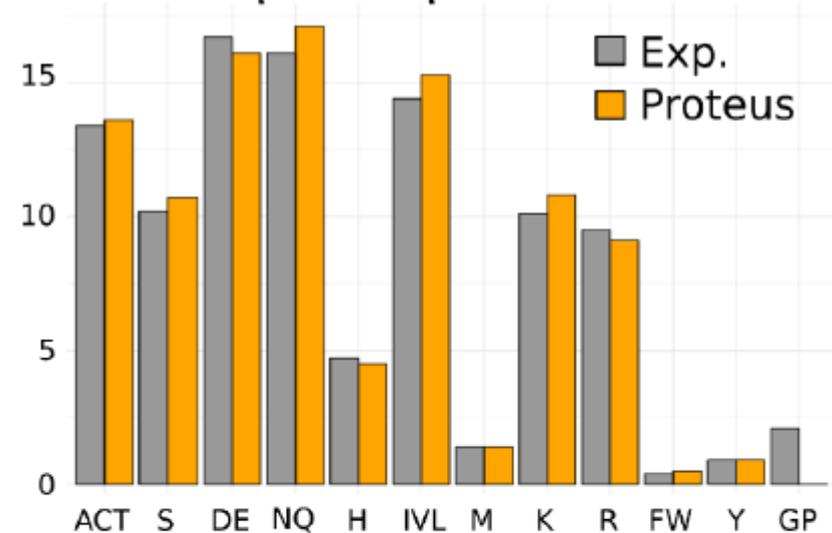
normalization constant

Choose unfolded parameters $E_u(t)$ to

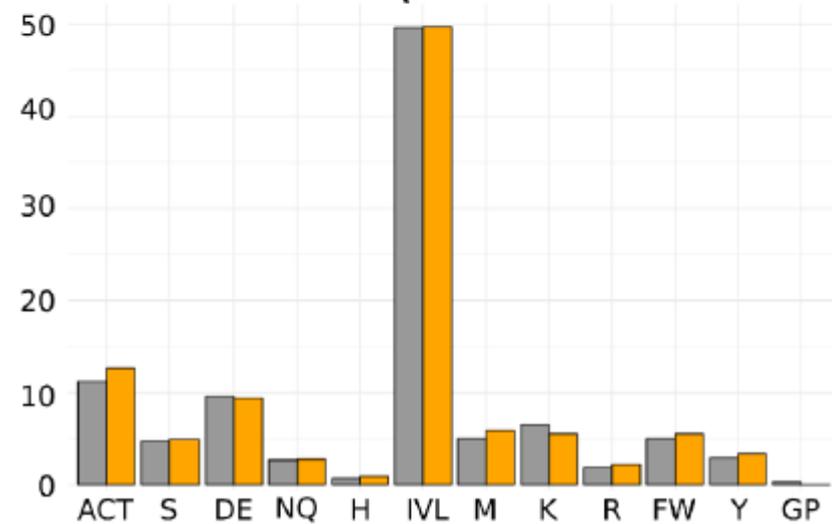
maximize probability or likelihood \mathbf{L} of \mathbf{S}

2 PDZ domains

Exposed positions



Buried positions



Maximum likelihood strategy for the unfolded state parameters

Choose a set of experimental sequences $\mathbf{S} = \{S, S', \dots\}$ as a reference

Boltzmann probability of sampling a sequence is $p(S) = \exp(-\beta\Delta G)/Z$

ΔG folding free energy, Z normalization constant

Choose unfolded parameters $E_u(t)$ to maximize probability or likelihood \mathbf{L} of \mathbf{S}

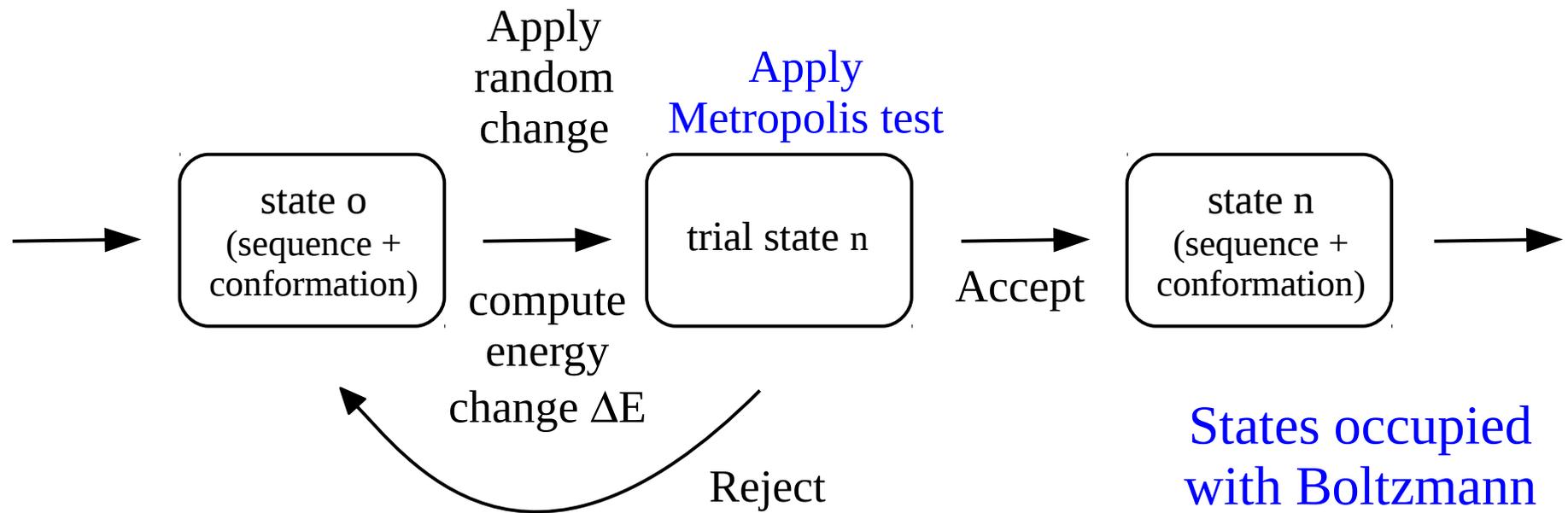
After some algebra, (exercise)

$\partial\mathbf{L}/\partial E_u(t) = 0 \iff$ experimental frequency $n_{\text{exp}}(t) =$ computational one $n_{\text{comp}}(t)$

Grad \mathbf{L} along $E_u(t) \propto n_{\text{exp}}(t) - n_{\text{comp}}(t)$

Search for maximum likelihood parameters using gradient method of your choice

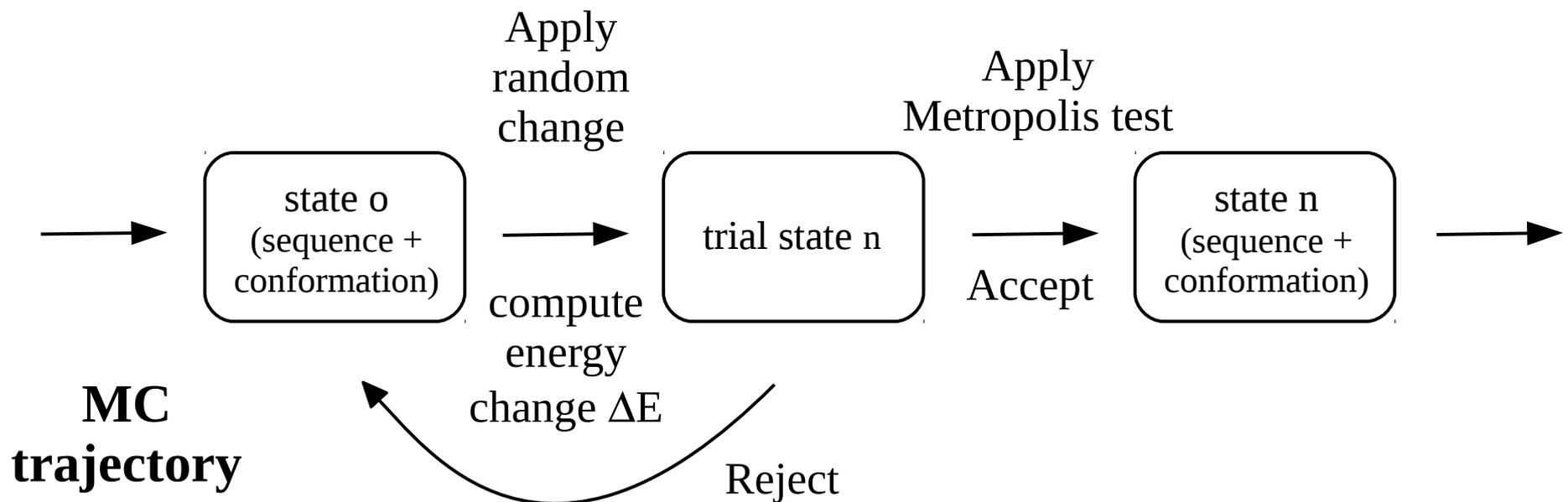
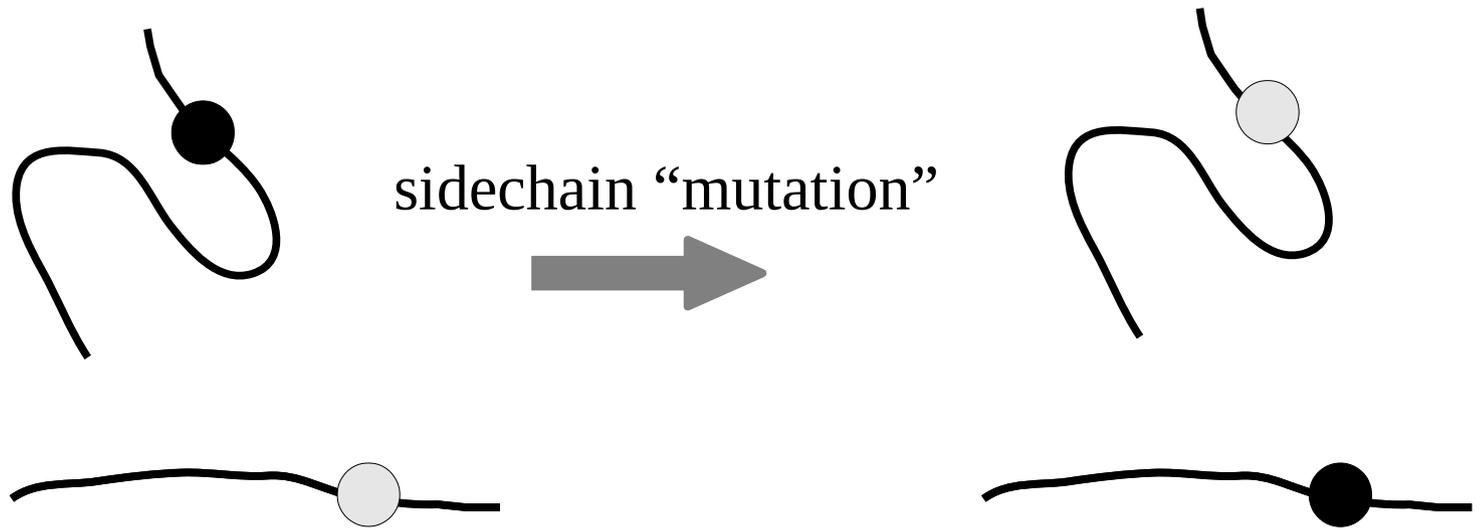
Eg: run MC with $\{E_u(t)\}$; adjust E_u values based on $n_{\text{exp}}(t) - n_{\text{comp}}(t)$; repeat



**MC trajectory
in sequence and
structure space**

States occupied
with Boltzmann
probability
 $\exp(-\beta E)$

**MC move in
“sequence
space”**

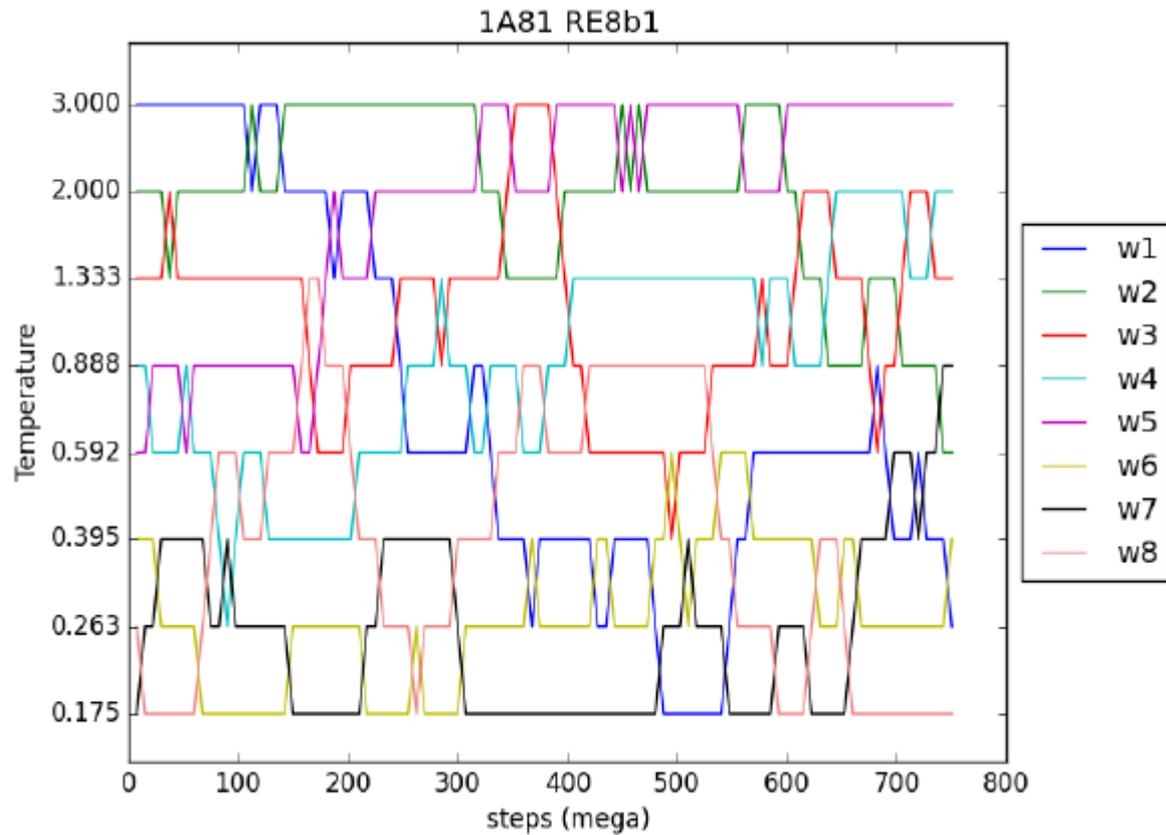


The simulation mimics a collection of sequences, present at equal concentrations, at equilibrium, distributed between folded/unfolded states according to Boltzmann statistics

Replica Exchange Monte Carlo

Multiple simulations or « replicas » are done in // at different temperatures; periodic exchanges of conformations between them are attempted following a Metropolis test: $\text{acc}(\text{swap}_{ij}) = \text{Min}[1, \exp[(\beta_i - \beta_j)(\Delta E_i - \Delta E_j)]]$

Thermal
energy kT
in kcal/mol

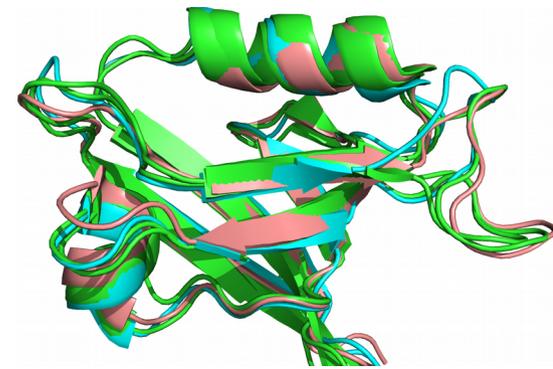


Replicas

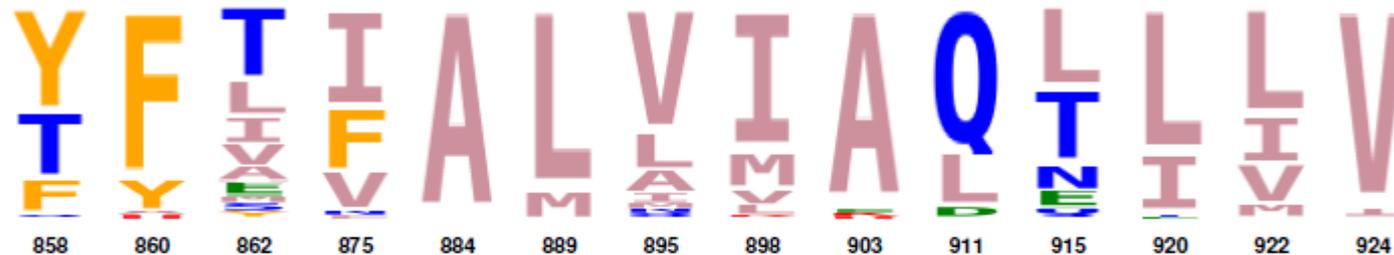
Monte Carlo steps (millions)

Enhanced sampling ; trivial cost thanks to OpenMP parallelization

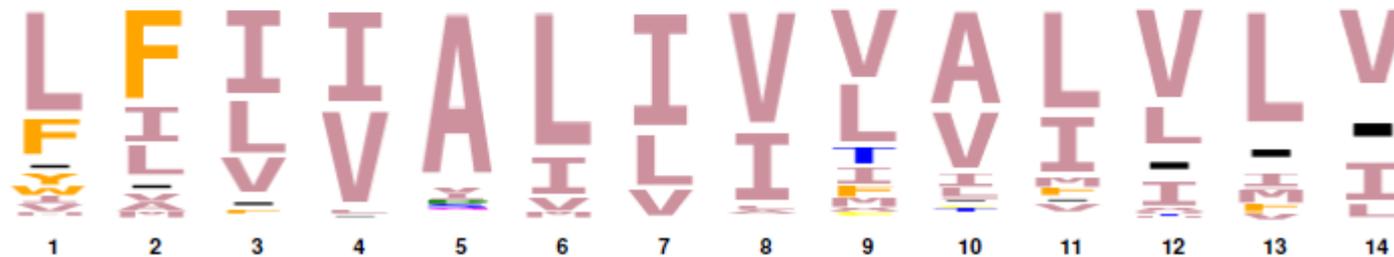
Experimental and designed PDZ sequences are similar (core positions)



Rosetta

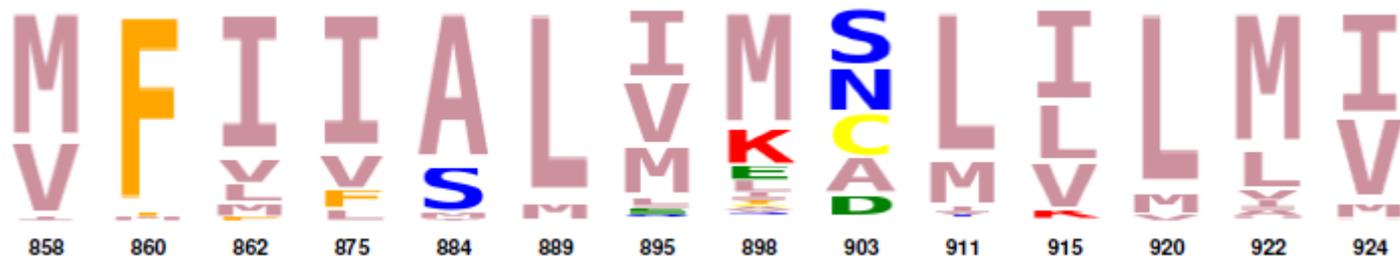


Pfam



Proteus

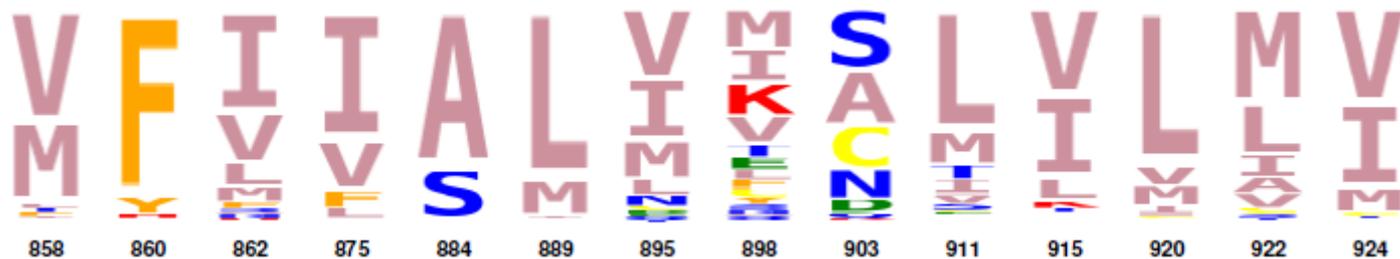
kT=0.6



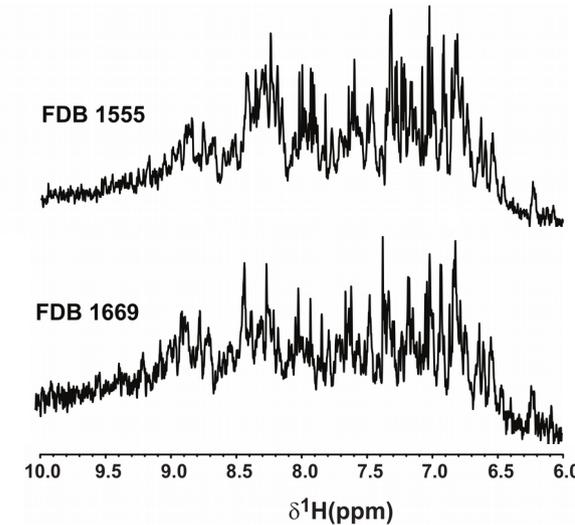
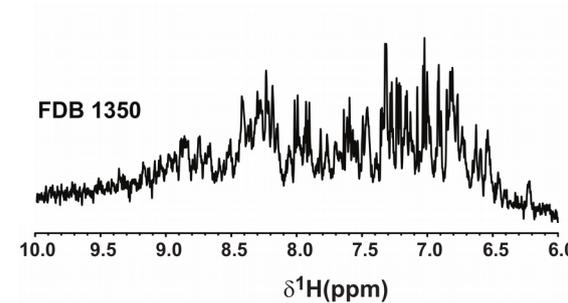
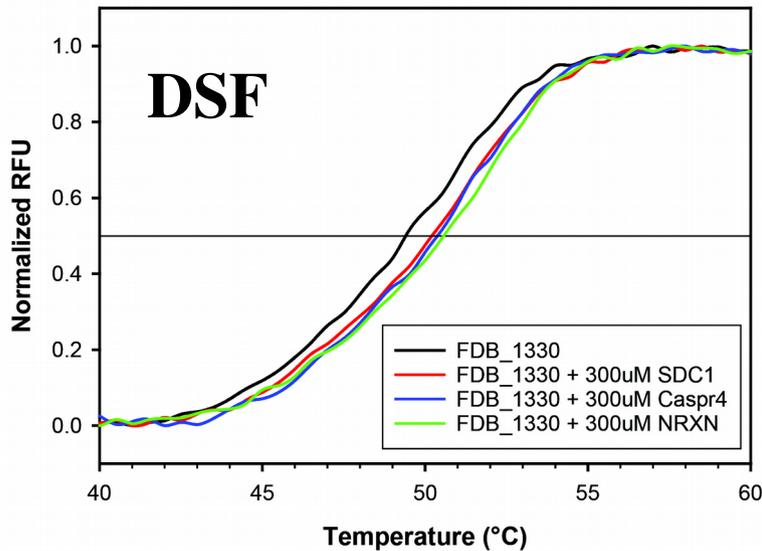
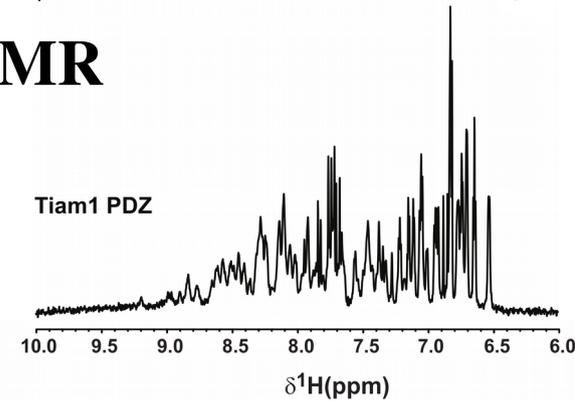
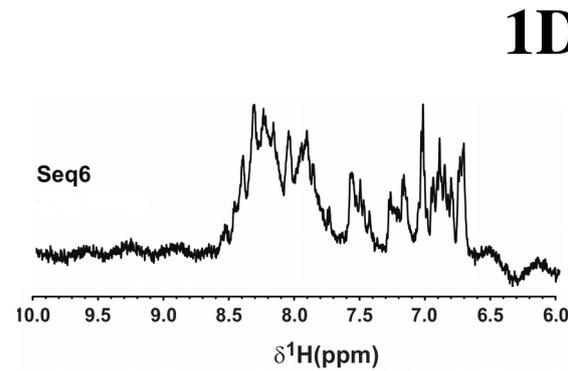
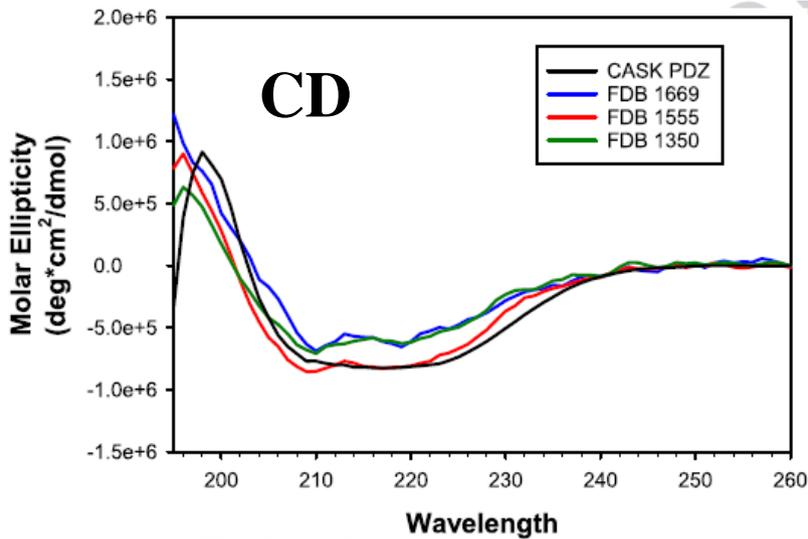
Proteus

kT=0.9

(kcal/mol)



3 out of 3 designed PDZ sequences tested adopt the correct fold and 2 bind the correct ligand

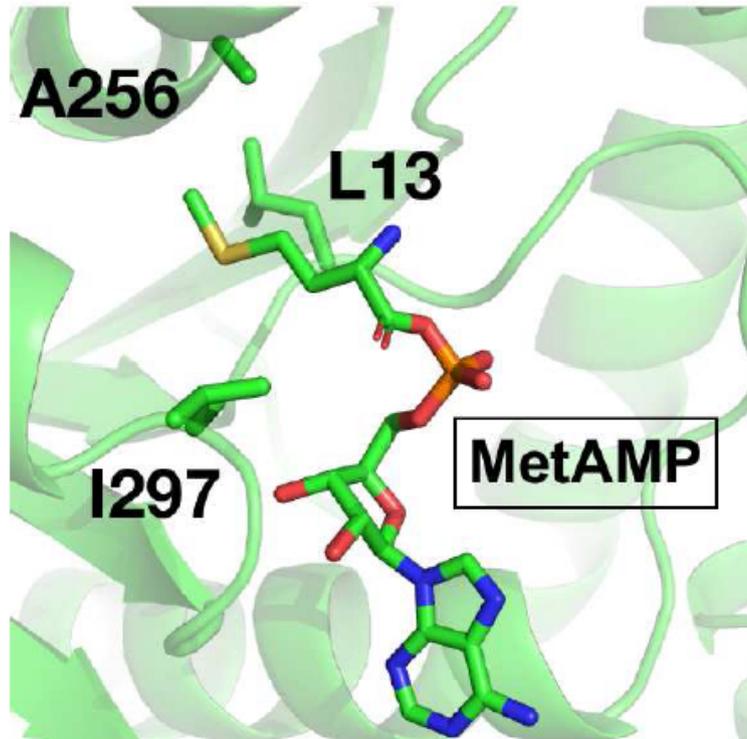


E Fuentes, U Iowa

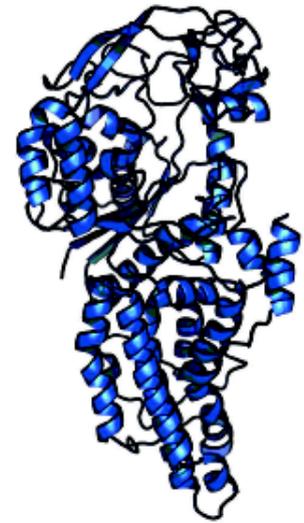
First successful complete protein redesign with a physics-based energy
Opou et al, Scientific Reports, 2020

Computational enzyme design

To design new enzymes, or redesign old ones, choose a few amino acids in the active site and explore mutations, in search of activity.

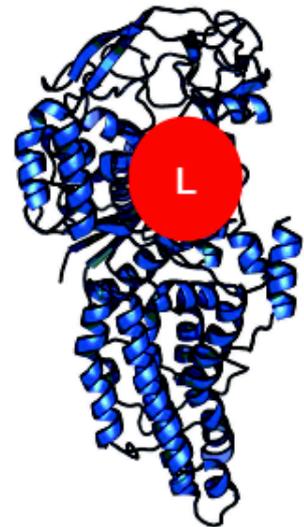


APO

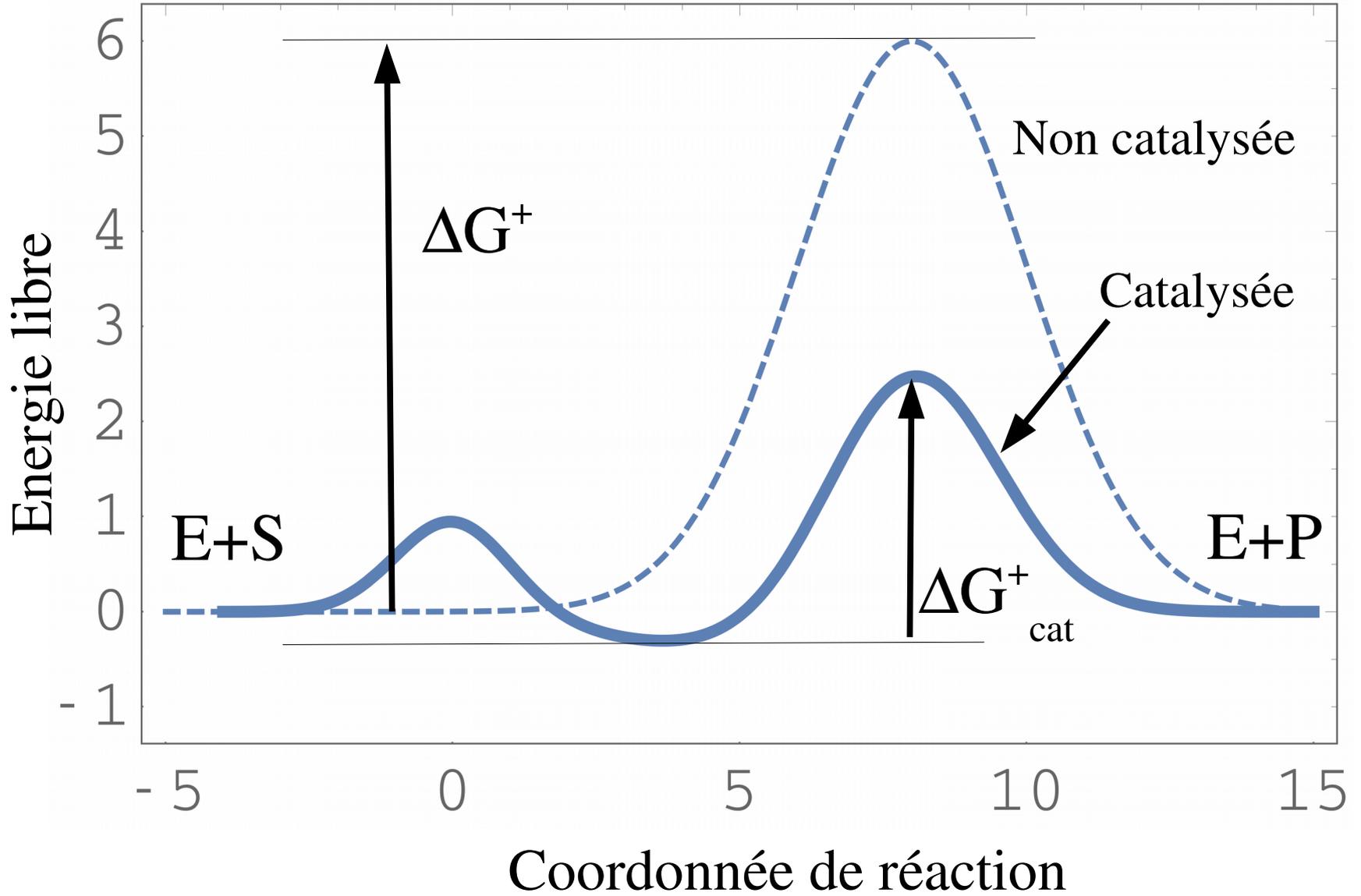


Affinity ↓

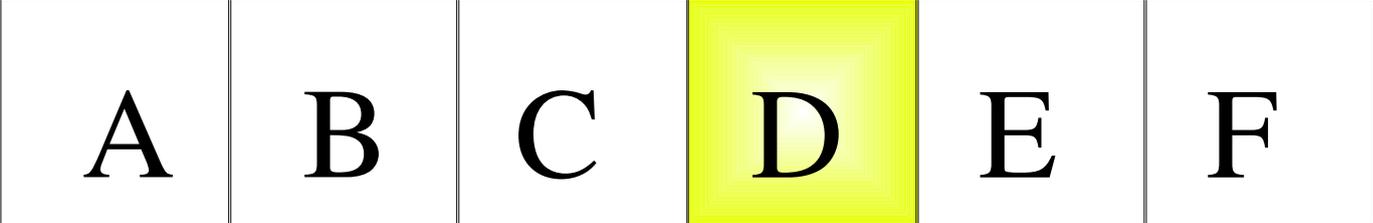
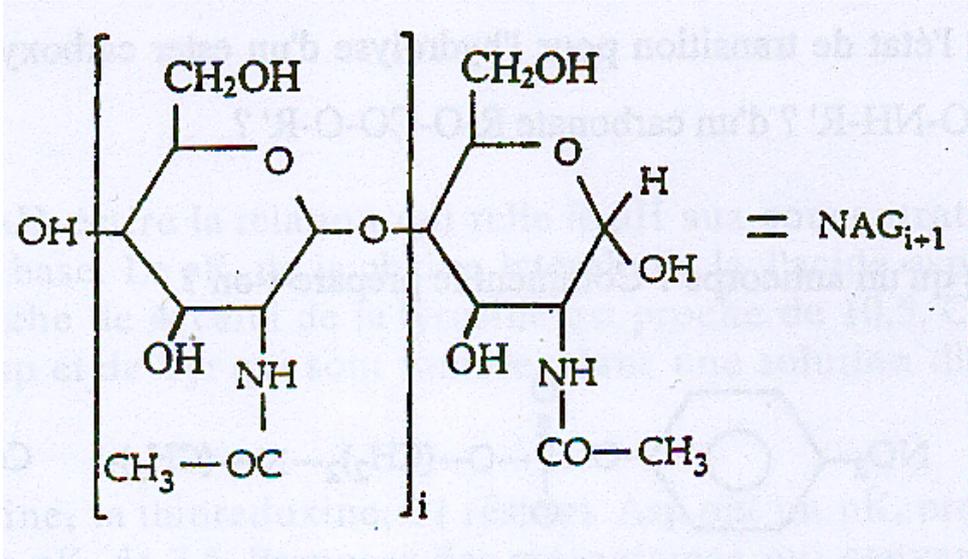
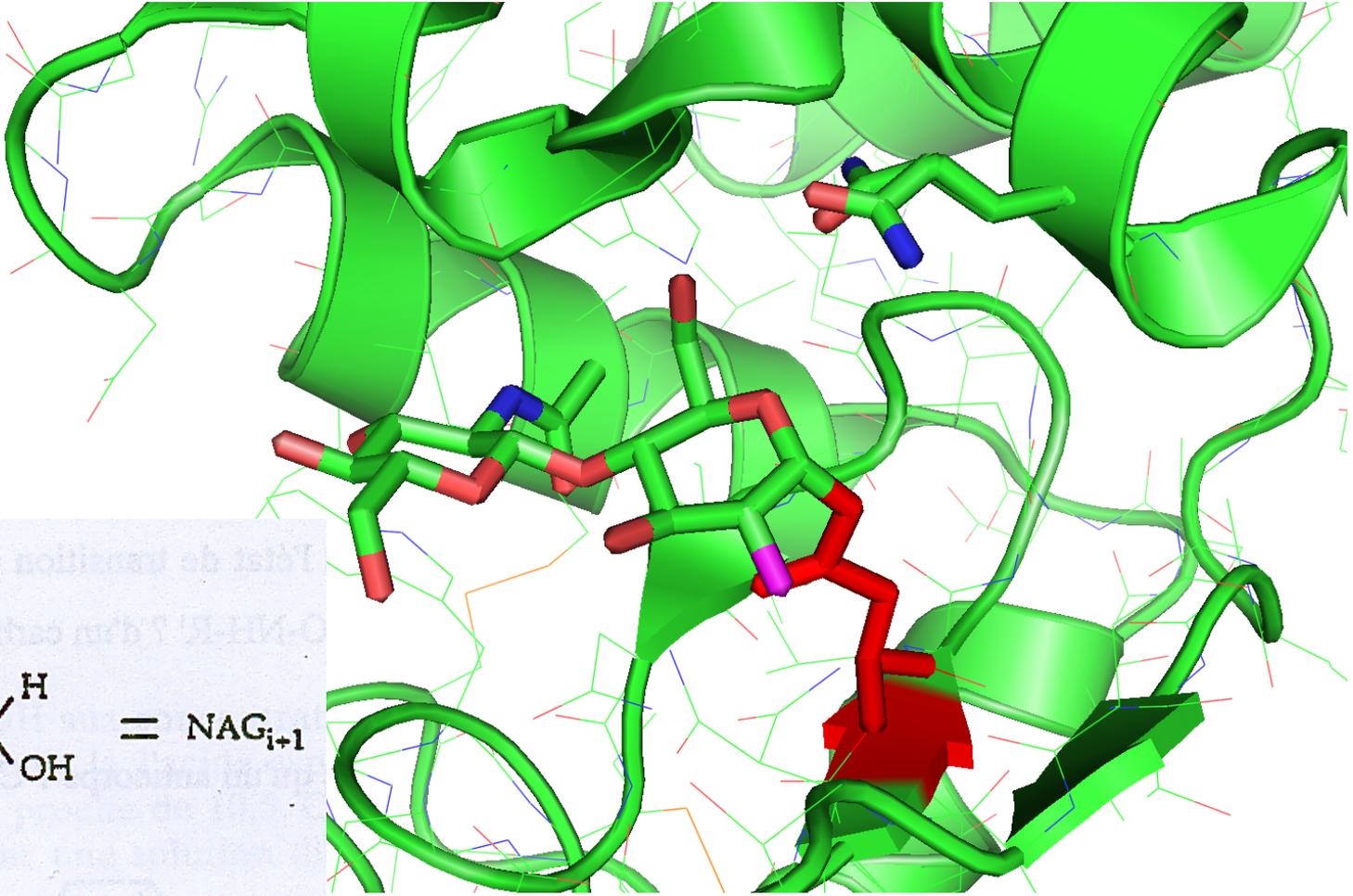
HOLO



Design criteria : protein stability, substrate binding, low reaction barrier
Novel substrate, new reaction.

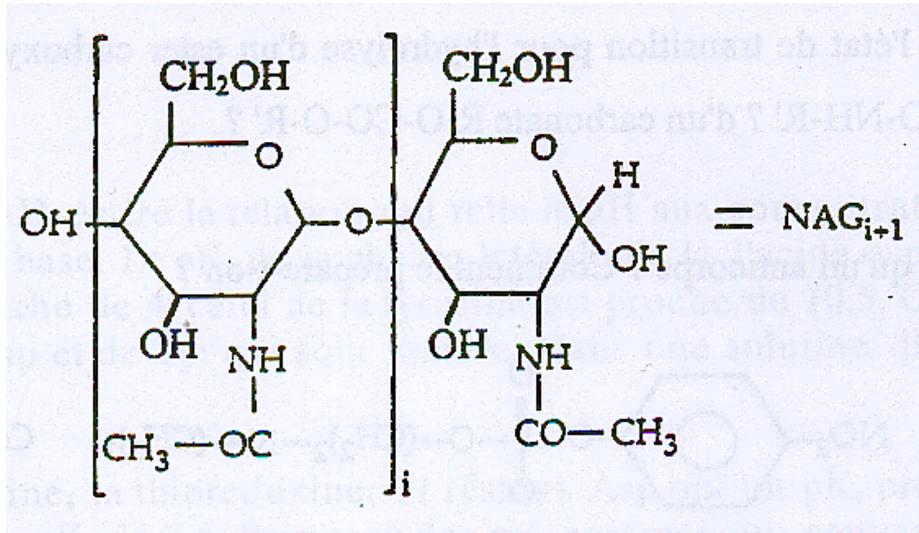


Lysozyme active site



Cuts between sites D and F

N-acetyl-glucosamine (NAG)



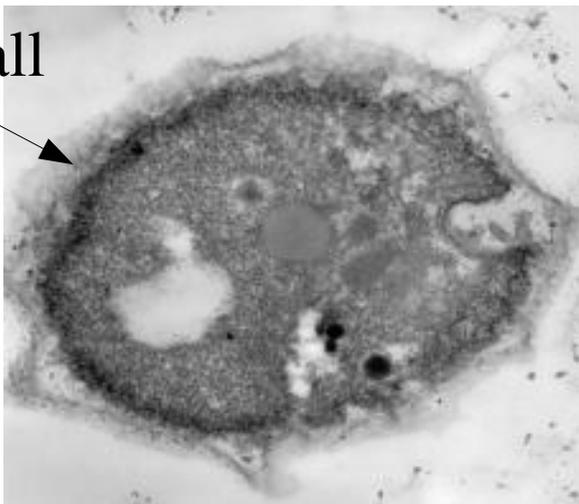
Lysozyme problem

Starch, glycogen
energy storage

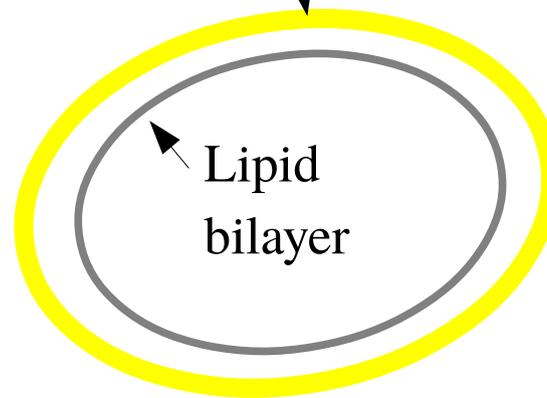
cellulose, chitin
structure



Cell wall

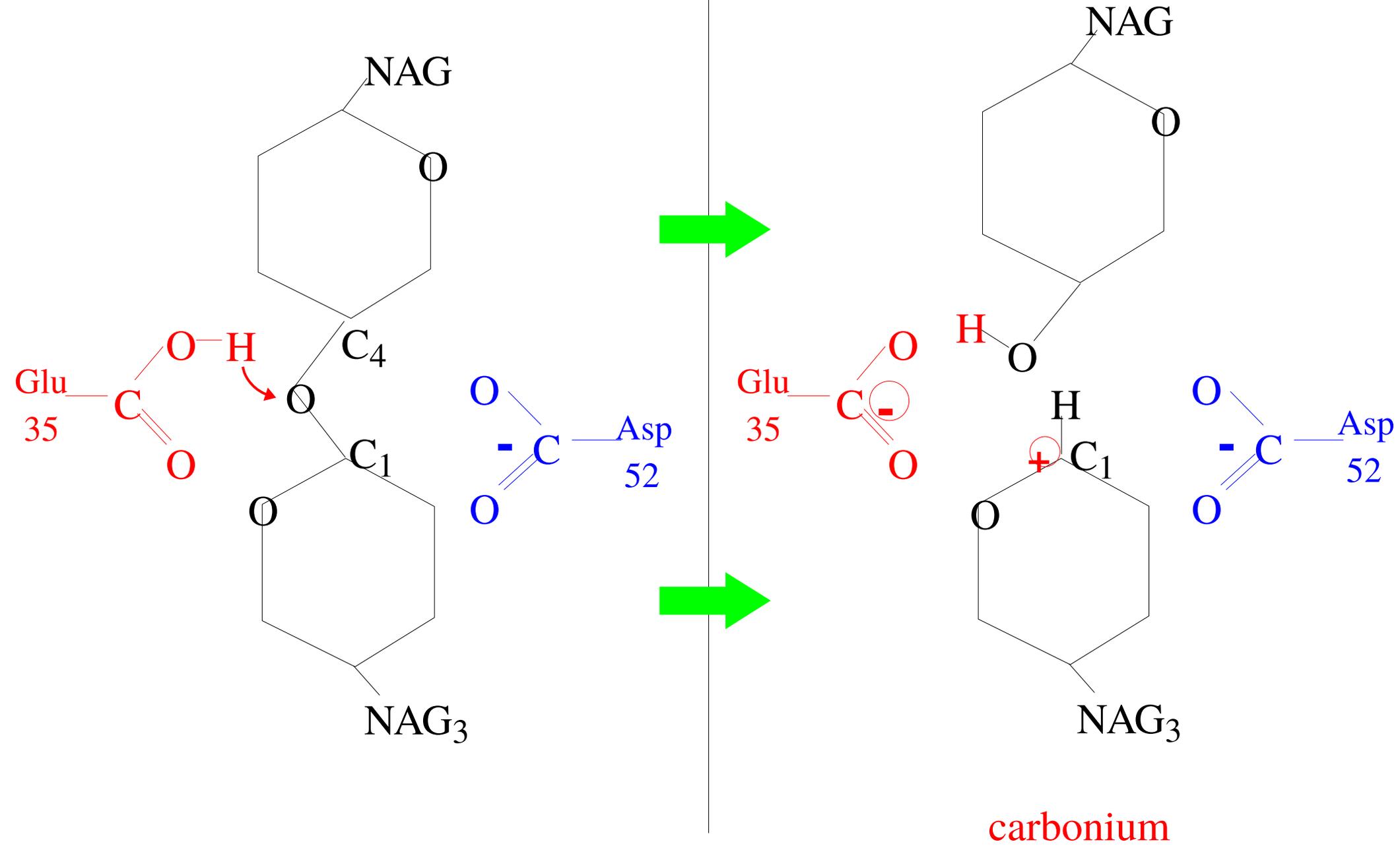


Polysaccharide
outer wall



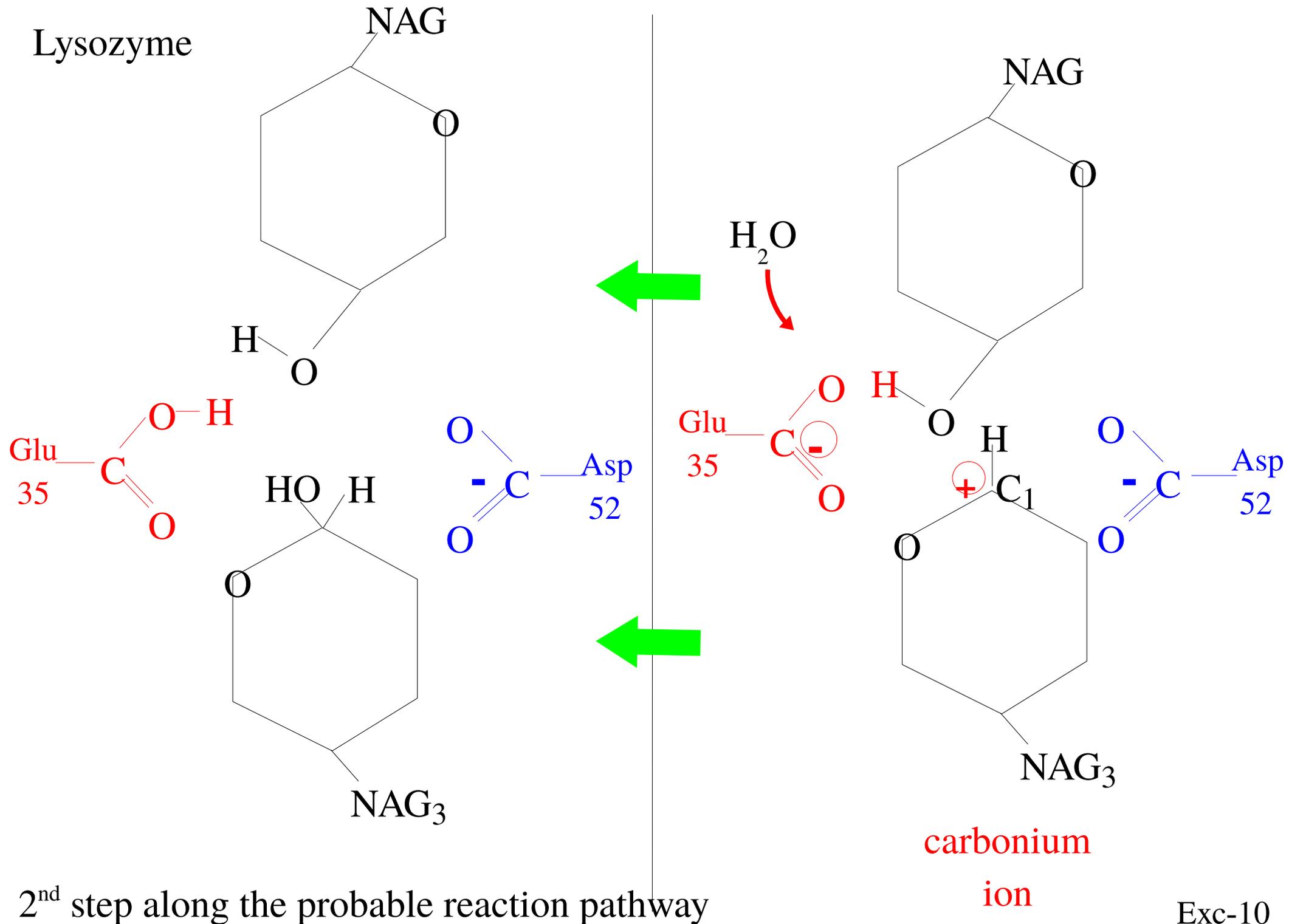
1 μm

Lysozyme



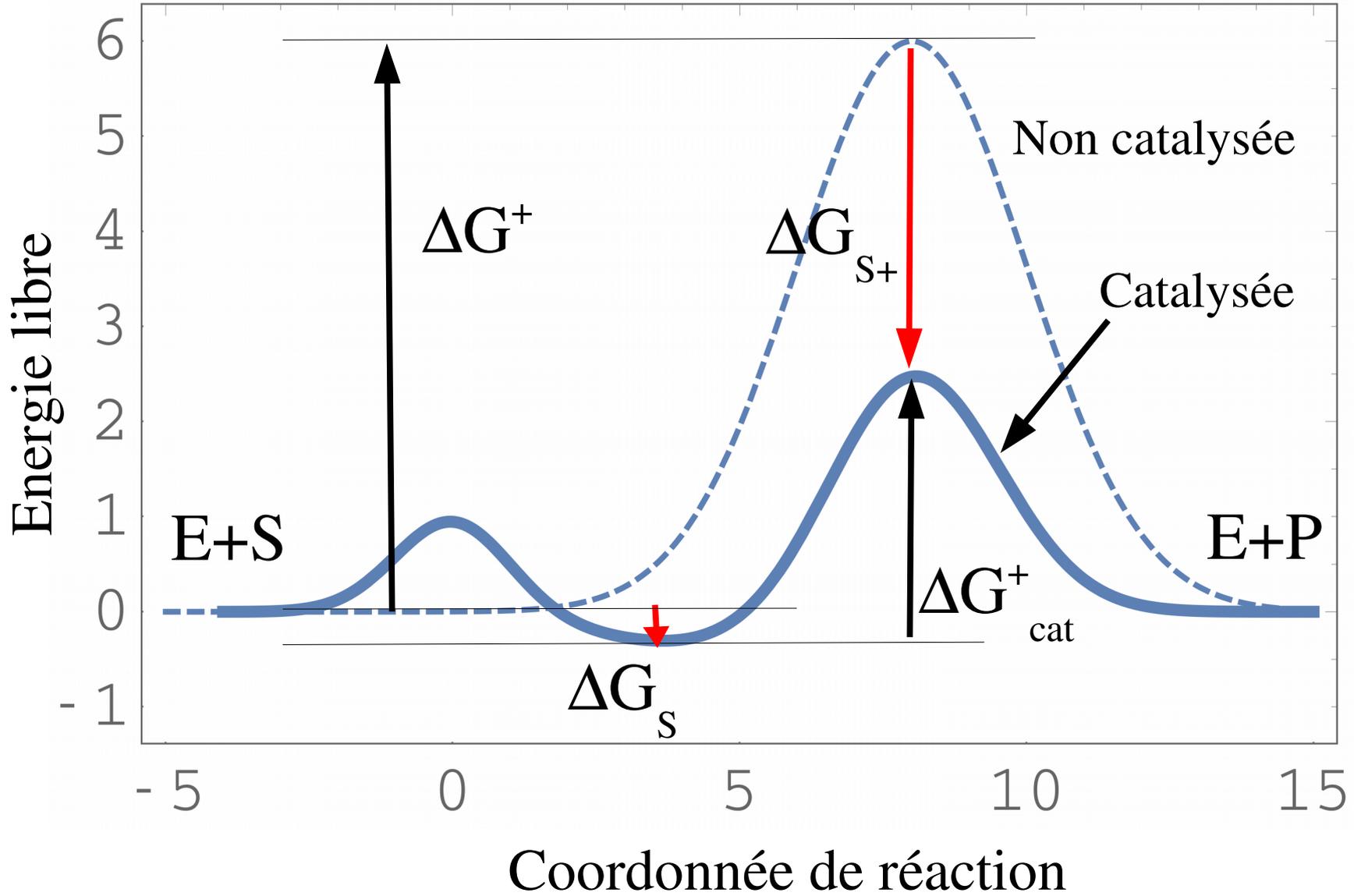
1st step along the probable reaction pathway

Lysozyme

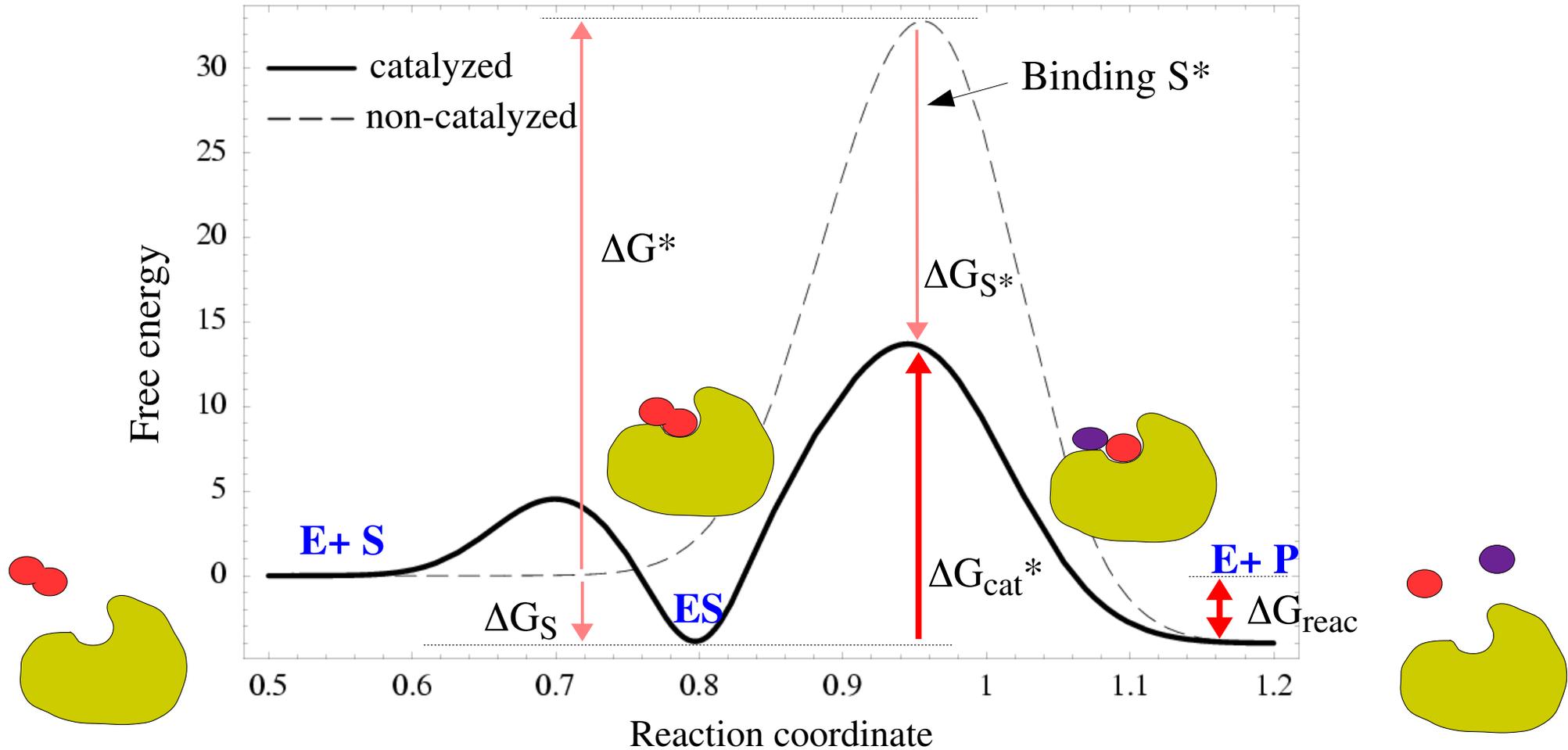


2nd step along the probable reaction pathway

Exc-10



Reaction coordinate and reaction pathway of an enzyme.
 Activation barrier for the reactions with and without the enzyme.

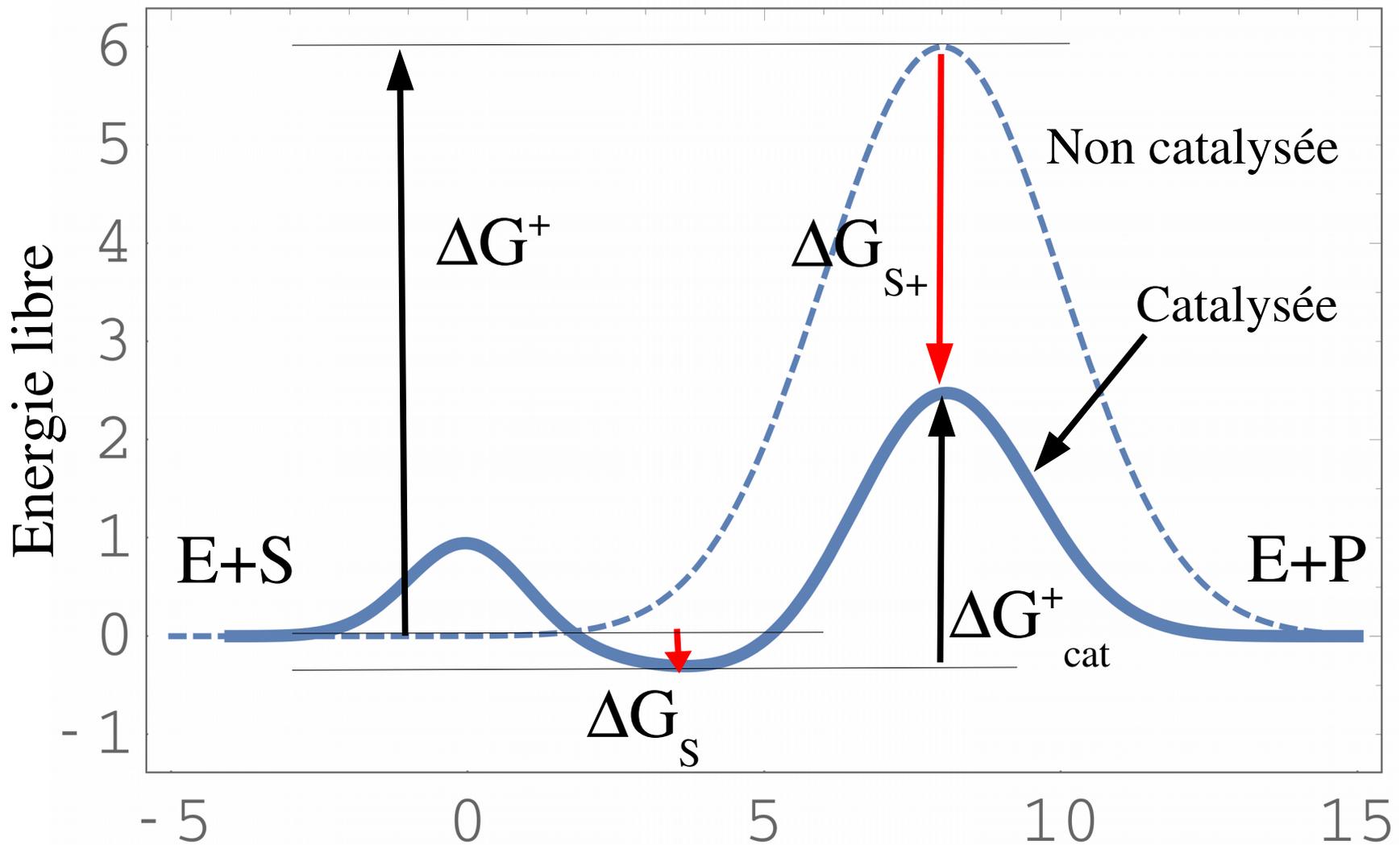


$$\Delta G_{cat}^* - \Delta G^* = \Delta G_{S^*} - \Delta G_S$$

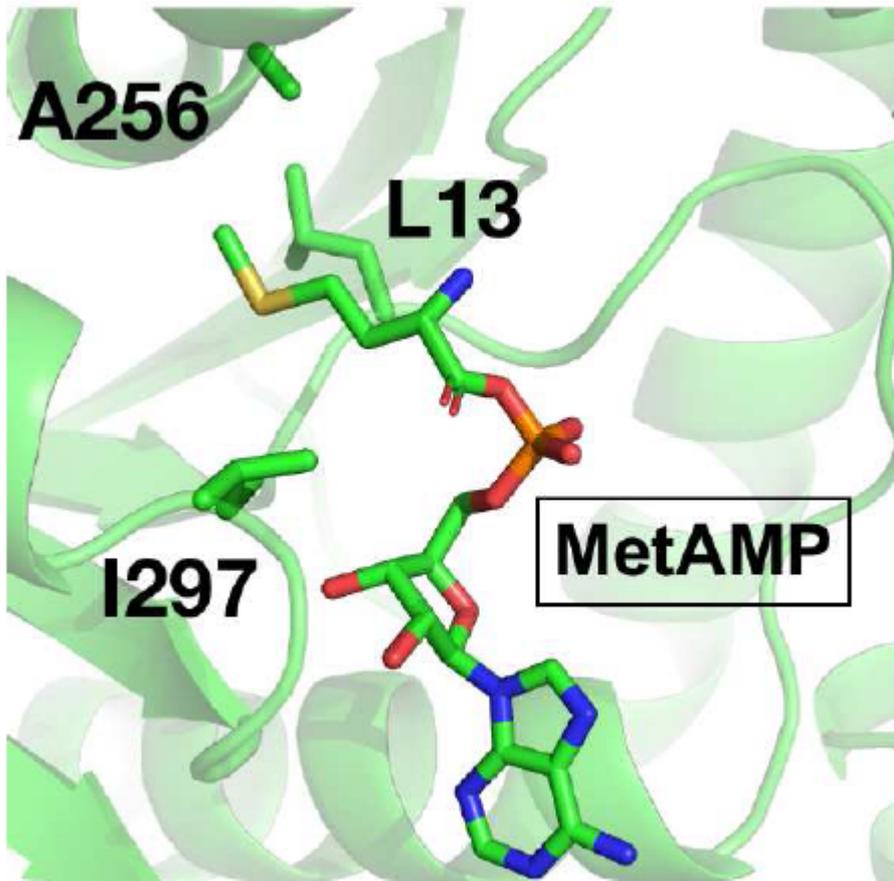
$$\Delta G^+ - \Delta G_S = \Delta G_{cat}^+ - \Delta G_{S^+}$$

$$\Delta G^+ - \Delta G_{cat}^+ = \Delta G_S - \Delta G_{S^+}$$

Pouvoir catalytique \Leftrightarrow Liaison de $S^+ >$ liaison de S



Computational enzyme design



Substrate complex

Design criterion : transition state binding.

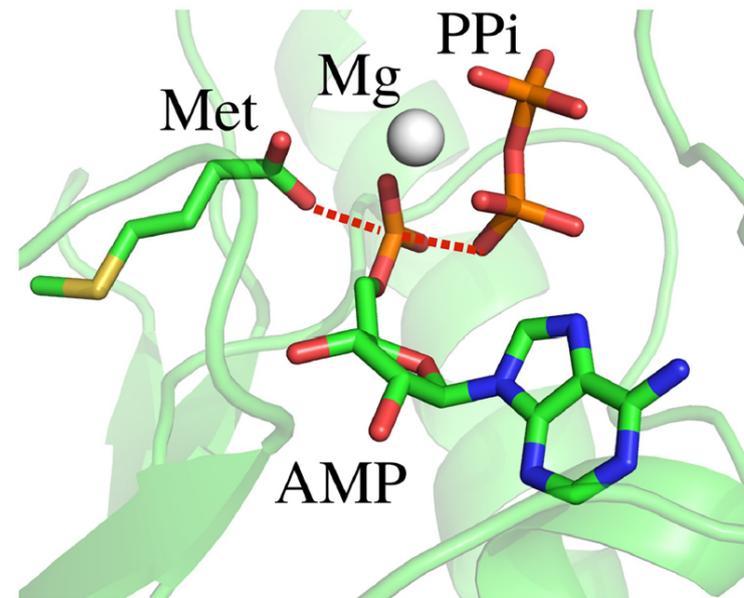
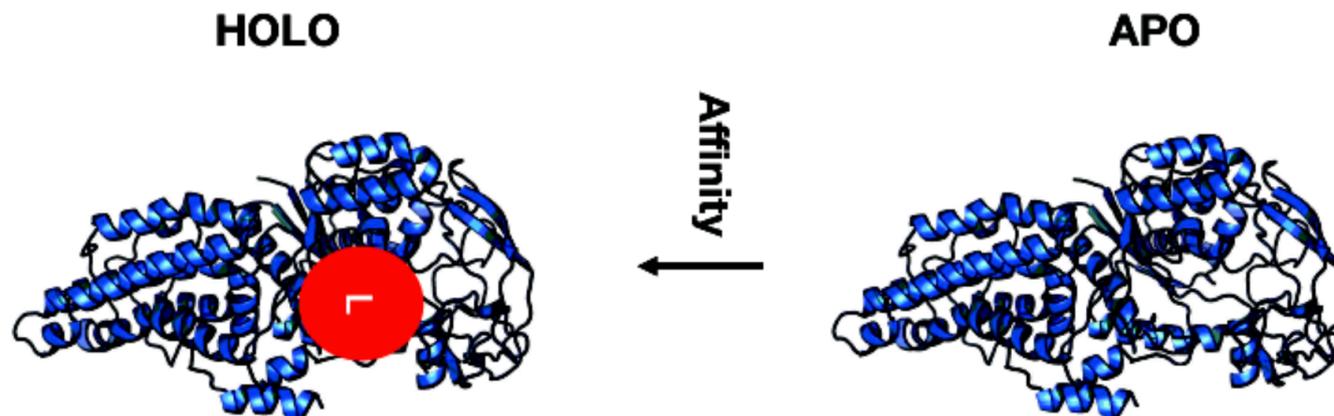
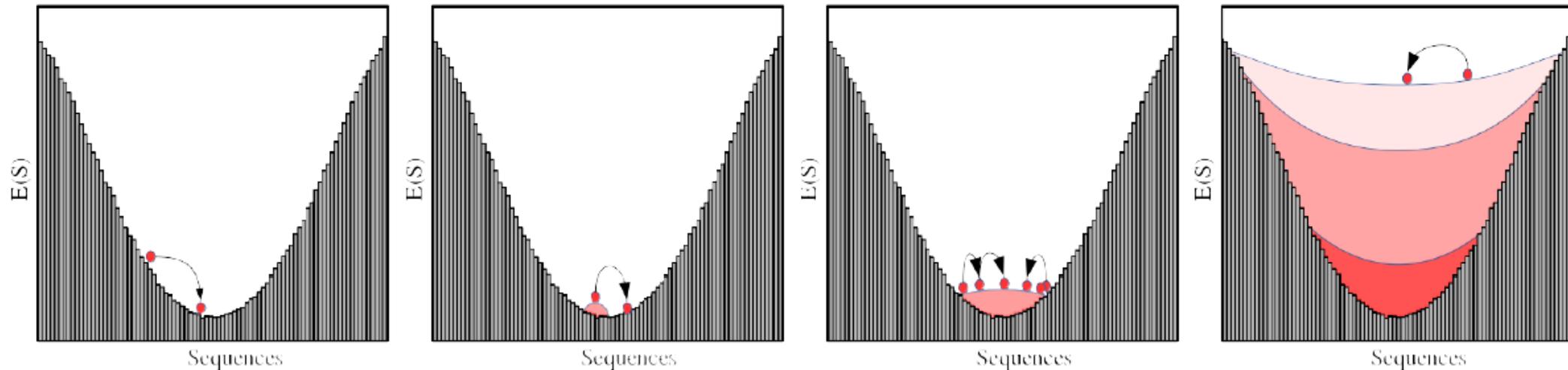


Fig 1. MetRS transition state for MetAMP formation. Closeup of the ligands.

Transition state complex

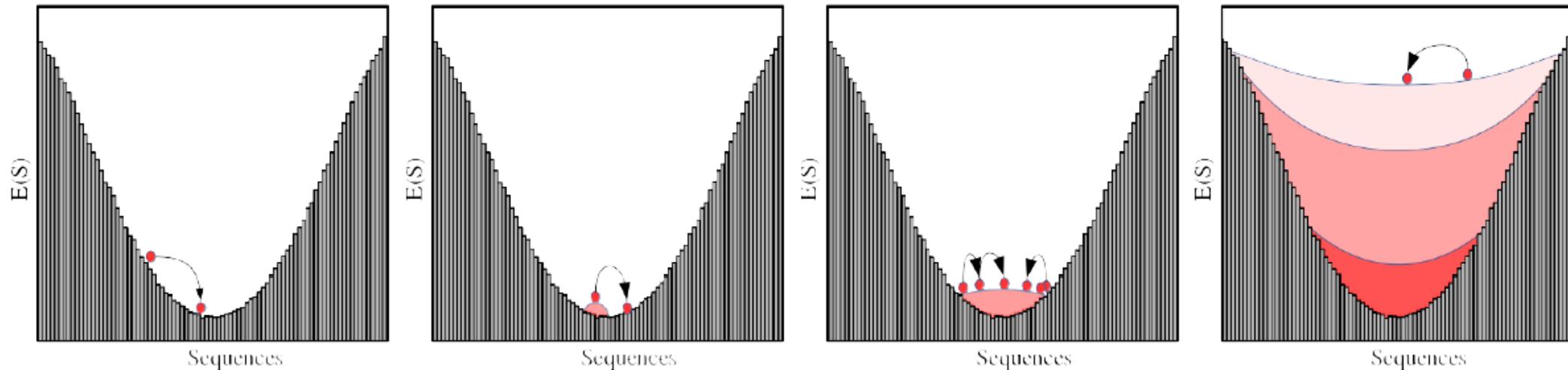
Step 1 : Adaptive landscape flattening of the unbound state

Simulate unbound state ; when a sequence S is sampled, add a small penalty to the energy function. After a while, the energy surface has been flattened.



Step 1 : Adaptive landscape flattening of the unbound state

Simulate unbound state ; when a sequence S is sampled, add a small penalty to the energy function. After a while, the energy surface has been flattened. The penalty function is $-E_{\text{apo}}(S)$.



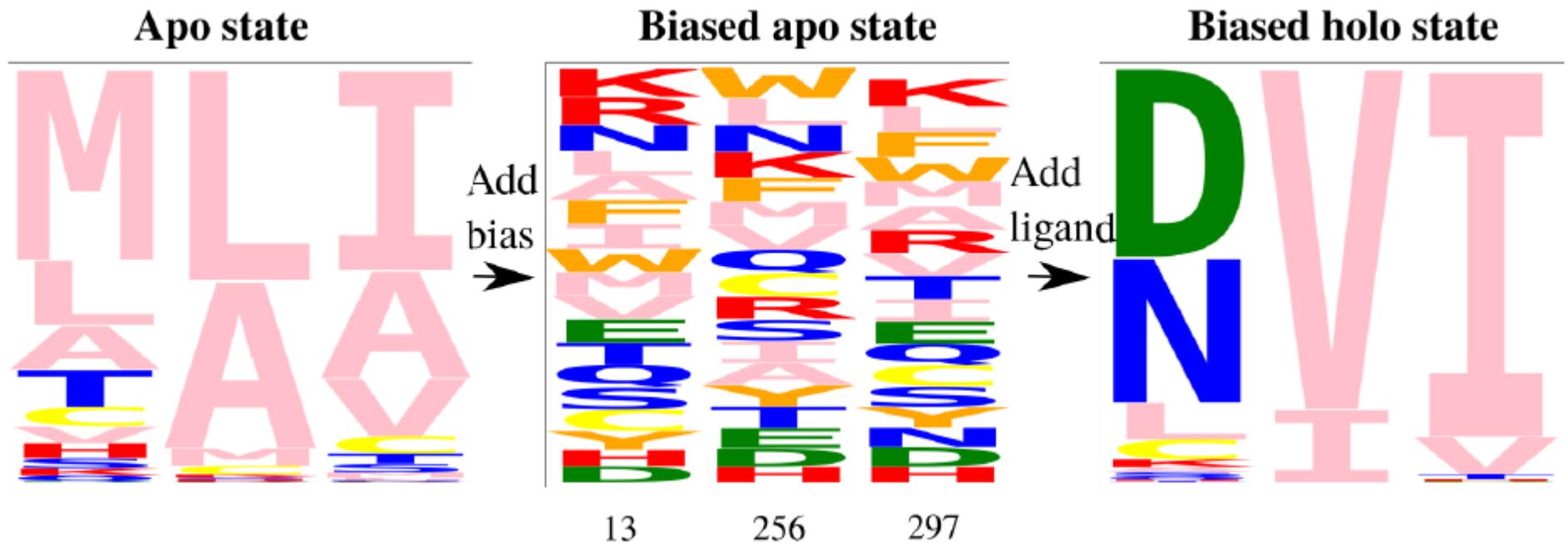
Step 2 : simulate the bound state, including the **penalty function**

The simulation is controlled by $E_{\text{holo}}(S) - E_{\text{apo}}(S)$, the binding energy

Step 1 : Adaptive landscape flattening of the unbound state

Step 2 : simulate the bound state, including the **penalty function**

The simulation is controlled by $E_{\text{holo}}(S) - E_{\text{apo}}(S)$, the binding energy



Methionyl-tRNA synthetase redesigned to accept nonnatural amino acids as substrates

A tool for synthetic biology and biotechnology

Inverse folding problem

Monte Carlo simulation, biased simulation techniques

Likelihood maximization = unsupervised machine learning of unfolded parameters

Adaptive landscape flattening : another unsupervised learning technique

An example targeting genetic code expansion (underway)

Enzyme catalysis !