

# **Drug Design**

# **Molecular docking**

**Karel Berka**

**Dpt. Physical Chemistry, RCPTM, Faculty of Science,  
Palacky University in Olomouc**

# Motto



[www.jolyon.co.uk](http://www.jolyon.co.uk)

tj. 18 hod 8 min 18 s

# Outline

- Structure-based drug design (SBDD)
  - Docking
  - Virtual screening
  - de novo design
  - Pharmacophore search
- Ligand-based drug design (LBDD)
  - Similarity matching
  - Pharmacophore search
  - QSAR

# Possibilities of Drug Design

	Known ligand	Unknown ligand
Known target structure	<p><b>Structure-based drug design (SBDD)</b></p> <p>Docking</p>	<p><b><i>De novo</i> design</b></p>
Unknown target structure	<p><b>Ligand-based drug design (LBDD)</b></p> <p><i>1 or more ligands</i></p> <ul style="list-style-type: none"><li>• Similarity search</li></ul> <p><i>Several ligands</i></p> <ul style="list-style-type: none"><li>• Pharmacophore</li></ul> <p><i>Large number of ligands (20+)</i></p> <ul style="list-style-type: none"><li>• Quantitative Structure-Activity Relationships (QSAR)</li></ul>	<p><b>CADD not possible</b> some experimental data needed</p> <p>ADMET filtering</p>

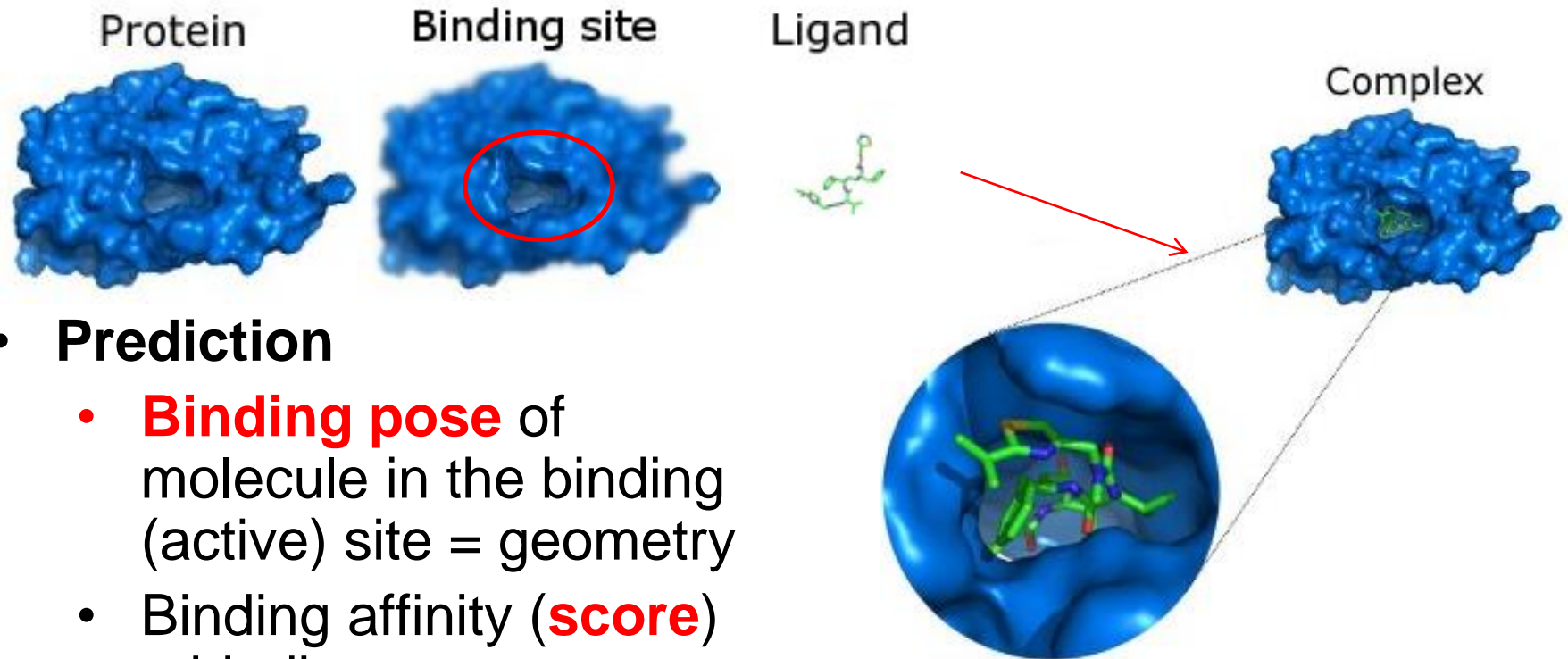
# Molecular Docking Idea

- Finding the best "fit" of ligand to receptor



# Molecular Docking

Computational method mimicking binding of ligand to receptor

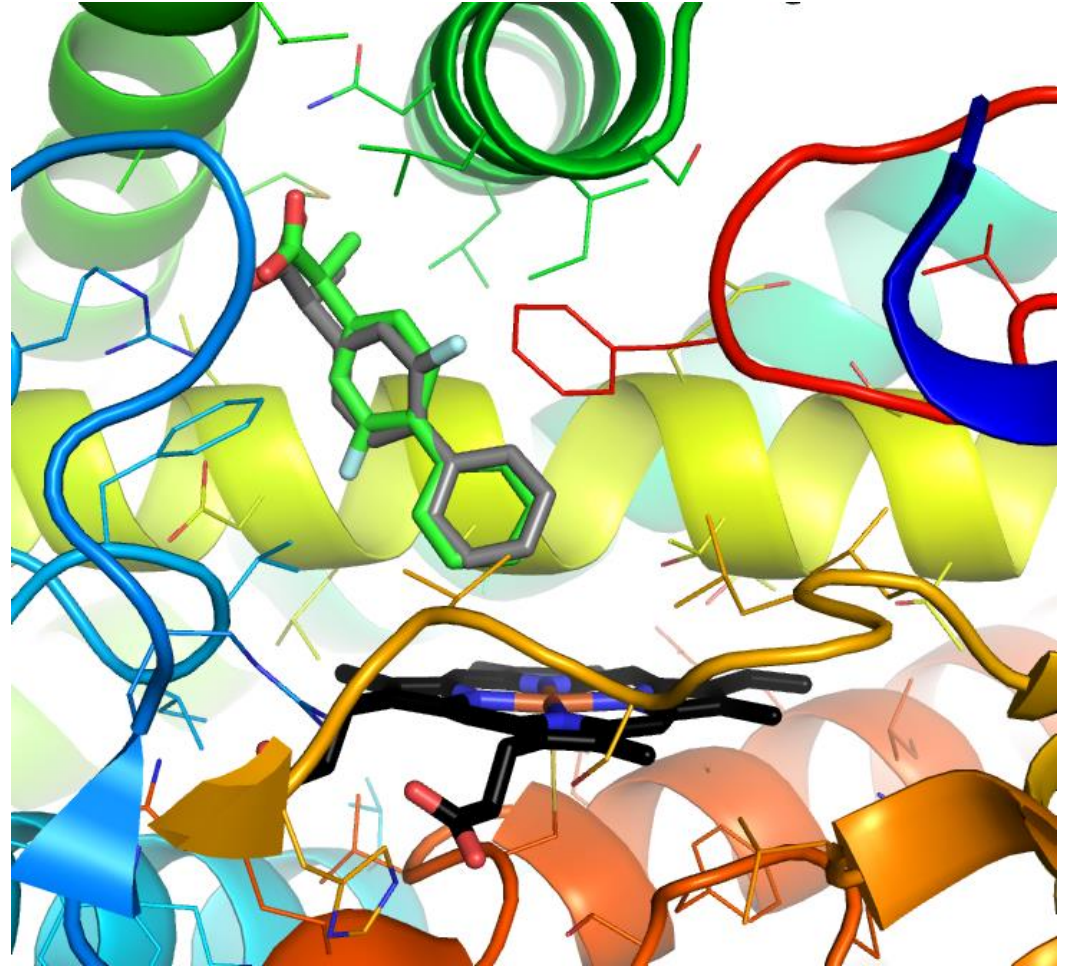


- **Prediction**

- **Binding pose** of molecule in the binding (active) site = geometry
- Binding affinity (**score**) = binding energy

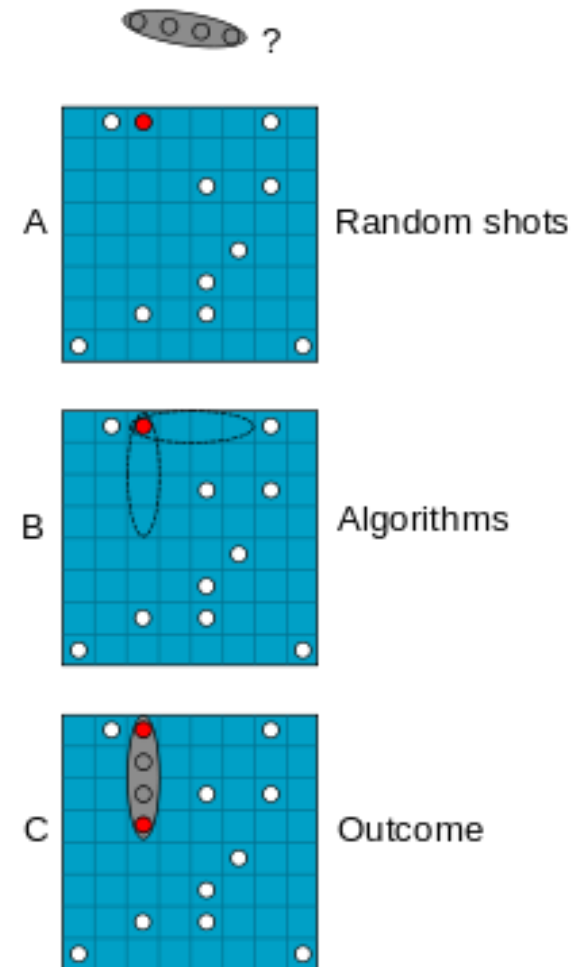
# Binding Pose

- Structural arrangement of ligand within receptor/enzyme
- Driven by molecular interactions



# Search Algorithms

- Monte Carlo
  - Random selection
  - Metropolis condition
    - (if better energy  $\rightarrow$  accept new pose; else check depend on energy difference)
- Genetic algorithms
  - Poses described by “Genes”
  - Best poses “mate” to generate offspring
  - Converge faster than MC
- Simulated heating
  - Heating – more energy – barrier crossing
  - Cooling – minima search





# Energetics

- Equilibrium binding constant

$$K_d = [P...L] / [P][L]$$

- correspond to free energy of binding:

$$\Delta G_{\text{bind}} = -RT \ln K_d$$

Free energy – combination of enthalpy and entropy

$$\Delta G_{\text{bind}} = \Delta H_{\text{bind}} - T\Delta S_{\text{bind}}$$

- $k_{\text{cat}}$ ,  $K_i$ ,  $IC_{50}$ ,  $EC_{50}$  – other values used for characterization
  - depend on concentration and affinity of substrate and concentration of protein

# IC<sub>50</sub>

- Concentration with 50% of inhibition activity

- Comparison of affinity between two compounds
- Cheng-Prusoff equation

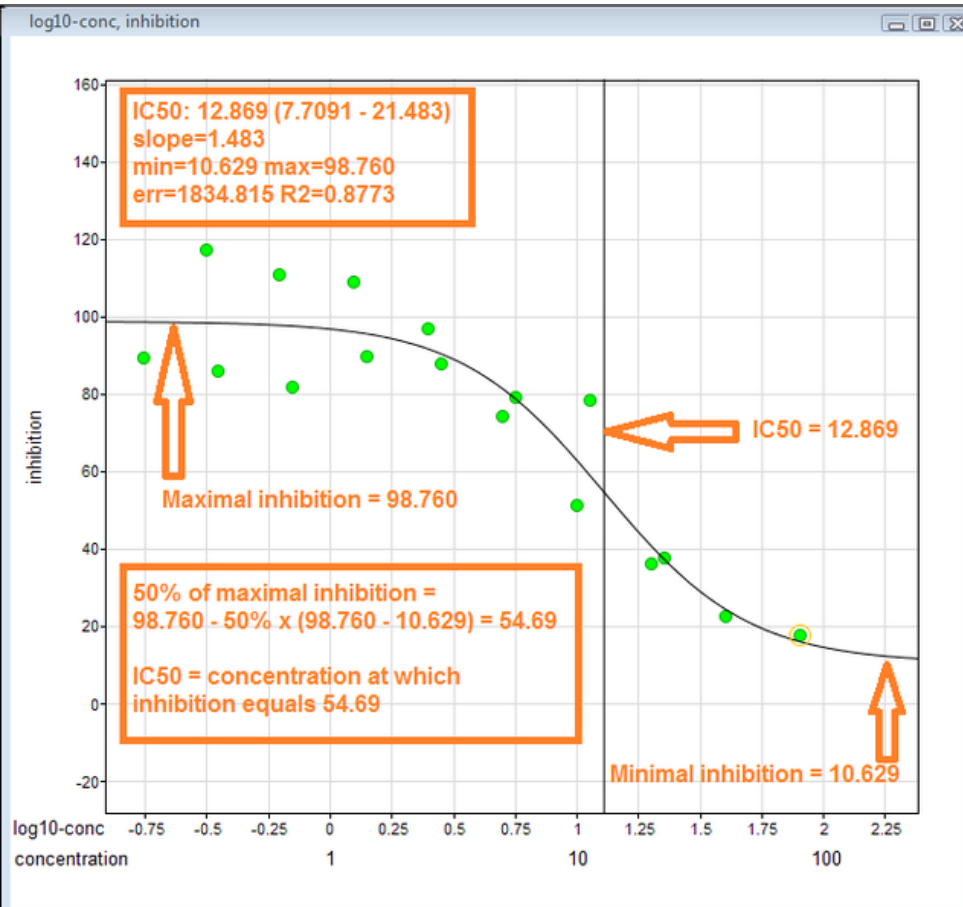
$$K_i = \frac{IC_{50}}{1 + \frac{[S]}{K_m}}$$

- Often logarithmic (mol/L)

$$pIC_{50} = -\log_{10}(IC_{50})$$

- Lower = better

pM (excellent) > nM (great) > μM (common) > mM (unusable)



Visual demonstration of how to derive IC<sub>50</sub> value: Arrange data with inhibition on vertical axis and log(concentration) on horizontal axis; then identify max and min inhibition; then the IC<sub>50</sub> is the concentration at which the curve passes through the 50% inhibition level. (wikipedia)

# Molecular Interactions

## Enthalpy:

- Electrostatics  
(partial charges)
- van der Waals  
(dispersion and repulsion)
- Hydrogen bonding  
(directionality)
- Desolvatation  
(cavitation energy)

## Entropy

- Conformation selection  
(flexibility)
- Solvation  
(hydrophobic effect)

# Scoring Function

- Binding affinity approximation

$$\Delta G_{bind} = \Delta G_{solvent} + \Delta G_{conf} + \Delta G_{int} + \Delta G_{rot} + \Delta G_{t/r} + \Delta G_{vib}$$

- It should be:
  - Quick
  - Score the right pose the best
- Parameterized against known binding poses and affinities
- Types:
  - Force-field (DOCK, Autodock, GoldScore)
  - Empirical (Glide, ChemScore)
  - Knowledge-based (DrugScore)

# Scoring Function

1. Score individual binding poses during search – **objective function**
  2. Identification of lowest (best) binding energy
  3. Sort **binding free energies** between individual ligands – selection of the best ligand
- Not necessarily the same for all points
- First part is most computationally intensive – needs to be quickest
  - Sorting should be the finest

# Scoring Function Types

- Force-field – based on molecular mechanical force-fields
  - Physical model - Interaction terms (elstatic, vdW,...)
  - Goldscore, DOCK, Autodock
- QM-based – based on quantum chemical calculations
  - PM6-DH
- Empirical – parameterized against exp. binding affinities ( $K_d$ ,  $IC_{50}$ )
  - Arbitrary terms (H-bonds, hydrophobic contacts)
  - ChemScore, PLP, Glide SP/XP
- Knowledge-based – based on protein-ligand complexes
  - Boltzmann hypothesis
    - typical binding motives -> stronger binding
  - PMF, DrugScore, ASP

# Force-field Scoring Functions

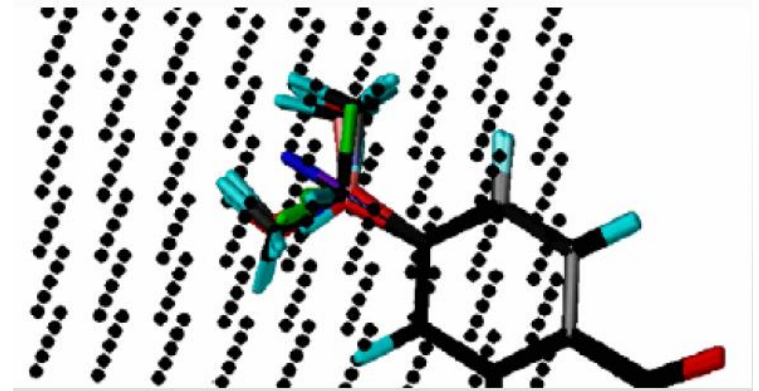
- **Physical interaction terms**

$$E = E_{\text{bond}} + E_{\text{angle}} + E_{\text{dih}} + E_{\text{coulomb}} + E_{\text{vdw}} + E_{\text{solv}}$$

- Often only **intermolecular** terms ( $E_{\text{coul}} + E_{\text{vdw}} + E_{\text{solv}}$ )
- **Intramolecular** are usually changed to rigid (bonds, angles) or screened by some value (dihedrals by 5 deg)

- **Grid** – time-saving

- Protein is divided into grid and interactions are pre-calculated at each point
- Ligands interaction is evaluated by multiplication of grid potential with ligand atoms
- Table search is quicker than full energy evaluation
- Receptor is usually one, while there is a series of ligands

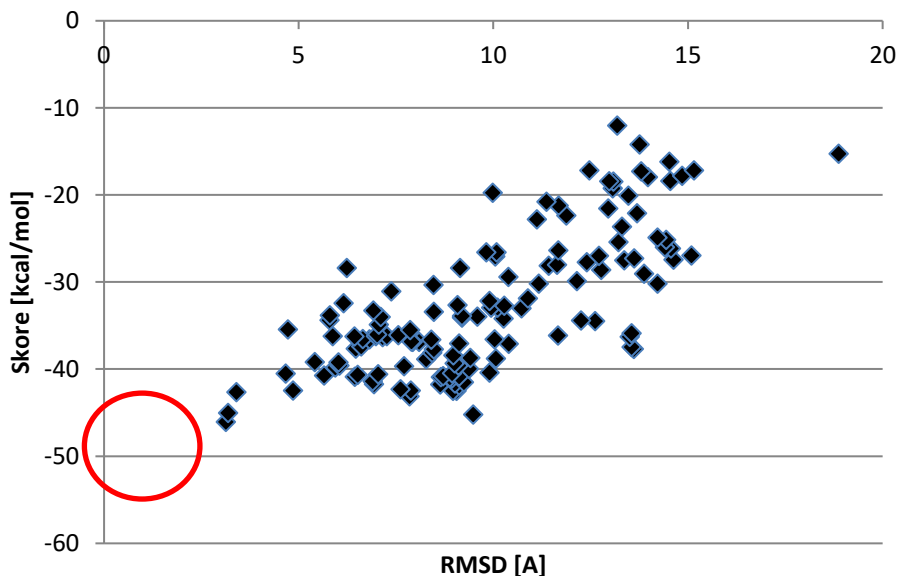


# Scoring Function Problems Example

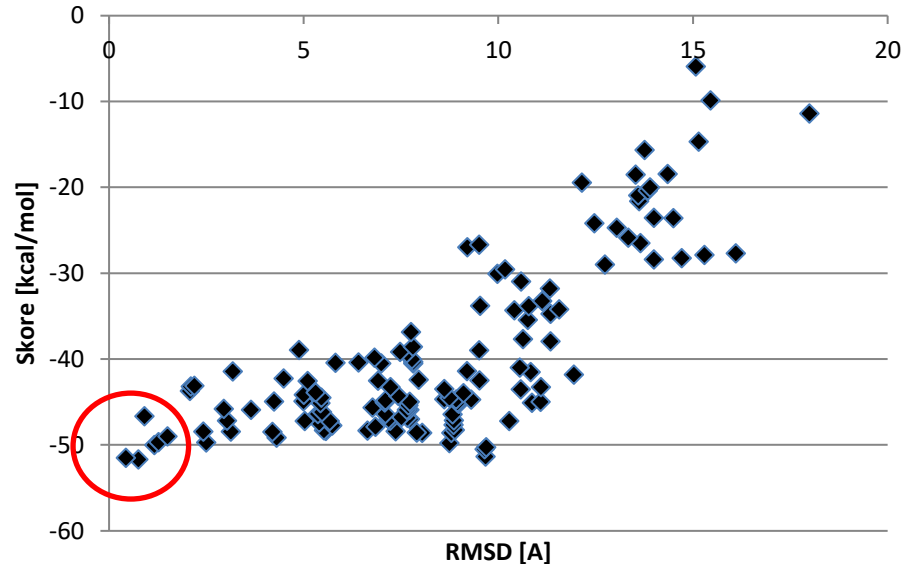


- Problems:
  - Repulsion
  - Electrostatics

original DOCK 6.6



DOCK 6.6  
with exponential repulsion





# QM based Scoring Function

- Based on quantum chemical calculations
- PM6-DH2

$$\Delta G'_w = \Delta H_w - T\Delta S_w + \Delta E_{\text{def}}(\text{I}) + \Delta\Delta G_w(\text{I}).$$

- $\Delta H_w$  - interaction enthalpy
- $-T\Delta S_w$  - interaction entropy
- $\Delta E_{\text{def}}$  - correction for inhibitor deformation
- $\Delta\Delta G_w$  - correction for inhibitor hydration

# Empirical scoring function

- Decomposition of binding energy into pre-defined “chemical terms”
- Specific interactions taken explicitly
  - H-bonding,  $\pi$ - $\pi$  stacking, ...

Linear form of terms is usually used (albeit unphysical)

$$DG_{bind} = DG_{solvent} + DG_{conf} + DG_{rot} + DG_t + DG_r + DG_{vib}$$

# Böhm's empirical scoring function

- linear summation of individual binding terms

- **Bohm's scoring function**

$$\Delta G_{bind} = \Delta G_0 + \Delta G_{hb} \sum_{h-bonds} f(\Delta R, \Delta \alpha) + \Delta G_{ionic} \sum_{ionic\ interactions} f(\Delta R, \Delta \alpha)$$
  - H-bonding, ion interaction, lipophilic interactions and conformational terms

- **Hydrogen bonding and ionic interactions**

$$+ \Delta G_{lipo} |A_{lipo}| + \Delta G_{rot} NROT$$
  - Depend on geometrical interaction – large deviations are penalized (ideal distance R, ideal angle  $\alpha$ ).

- **Lipophilic term**
  - Proportional to lipophilic surface contact between protein and ligand ( $A_{lipo}$ )

- **Conformational entropic term**
  - penalization for freezing of internal rotations of ligand - entropy
  - Proportional to number of rotationable bonds of ligand (NROT)

- $\Delta G$  values of individual terms are constants obtained by linear regression on experimental binding data on 45 protein–ligand complexes

# Chemscore

- Original Chemscore function for binding free energies

$$\Delta G_{binding} = \Delta G_o + \Delta G_{hbond} S_{hbond} + \Delta G_{metal} S_{metal} \\ + \Delta G_{lipo} S_{lipo} + \Delta G_{rot} H_{rot}$$

- $S_{hbond}$  – hydrogen bonding
- $S_{lipo}$  – lipophilic interactions
- $S_{metal}$  – acceptor-metal interactions
- $H_{rot}$  – loss of conformational entropy on ligand binding

J. Comput. Aided Mol. Des. 11, 425-445, 1997

# Chemscore

## Chemscore for docking

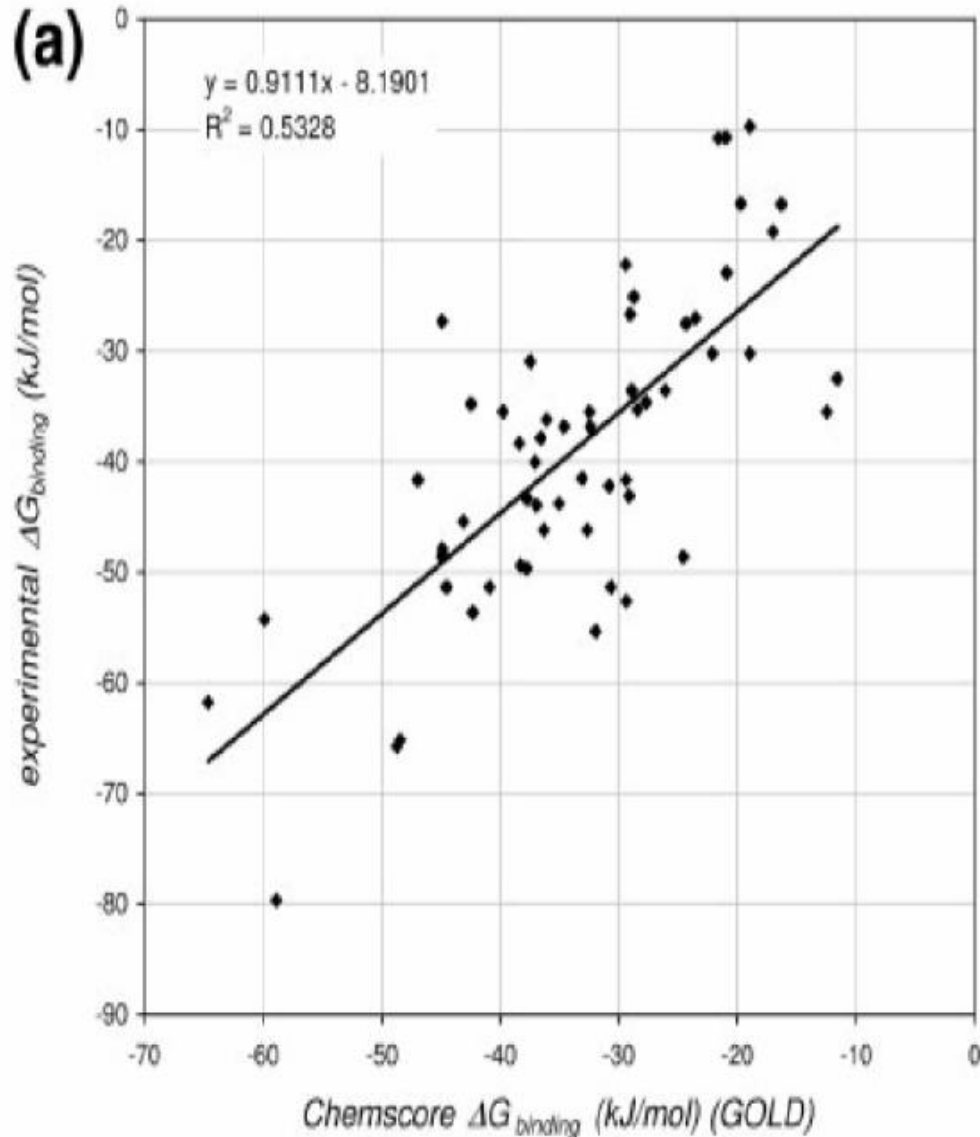
- Add more terms – clash, ligand internal, protein-ligand covalent

$$\Delta G'_{binding} = \Delta G_{binding} + E_{clash} + E_{int} + E_{cov}$$

- Complex functional forms – look them up!
- Parameters carefully rederived

Proteins 52, 609-623, 2003

# Chemscore Accuracy



Correlation coefficient –  $r$

$r^2 < -1, 0, 1 >$

-1 – anticorrelation

0 – no correlation

1 – full correlation

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

# Empirical Scoring Functions Problems

- Heavy dependence on training set
- Can have missing interaction terms
  - metal-ion
- Parameterized on success
  - Use of molecules that bind in parameterization => artificial binding of molecules that otherwise would not bind
  - => Use of **decoys** – molecules, which are of similar size as those really binding but not binding

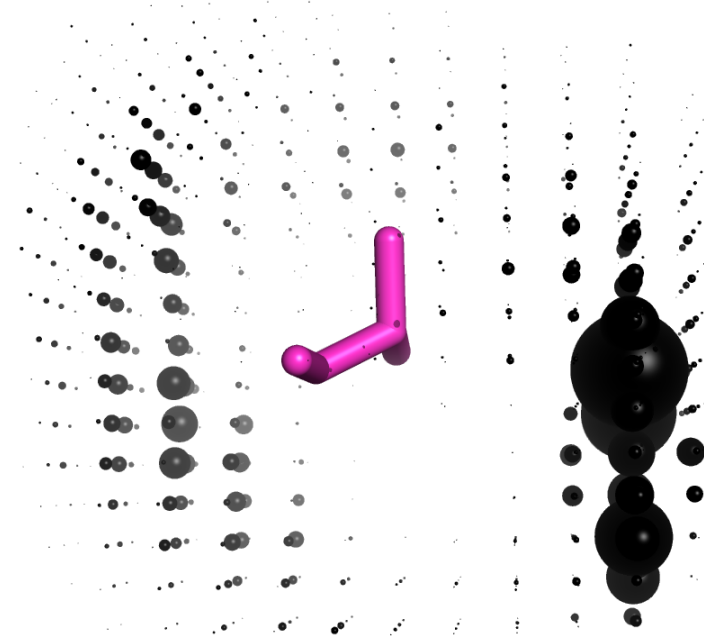
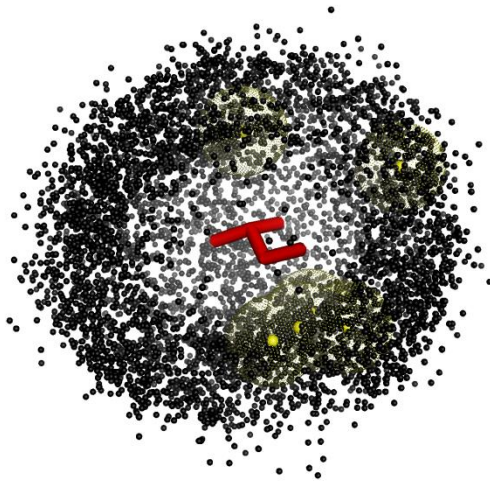
# Knowledge-based Function

Correlation of structural data from ligand/protein complexes with free energy of binding

- Use a rigorous statistical mechanical result:

$$A = -kT \ln g(r)$$

- This equation holds for an ensemble of particles at equilibrium (in gas)
- not necessarily proteins





# Drugscore

## DRUGSCORE

$$\Delta W_{i,j}(r) = W_{i,j}(r) - W(r) = -\ln \frac{g_{i,j}(r)}{g(r)}$$

$$g(r) = \frac{\sum_i \sum_j g_{i,j}(r)}{i*j}$$

Short-range (6 Å) contributions only – ignoring solvation

# Docking Preparation

- Receptor
  - Identification of relevant structure
  - Structure preparation (missing atoms, hydrogen assignment)
- Ligand
  - Structure preparation
  - Isomers, conformations
- Other tasks
  - Water
  - Flexibility

# Receptor Preparation

- Where
  - identification of binding site
- Good structure
  - Low R (accuracy)
  - Low B-factors (flexibility)
  - Low R-free (correctness)
- Flexibility
  - Rigid docking into several structures
    - Molecular Dynamics
    - more Xtals
  - Flexible docking

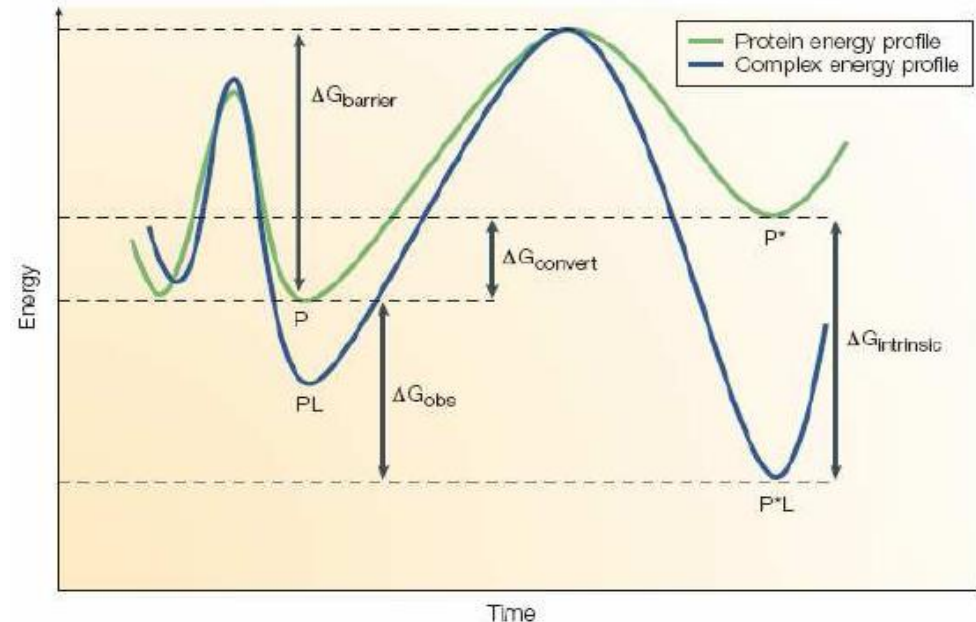
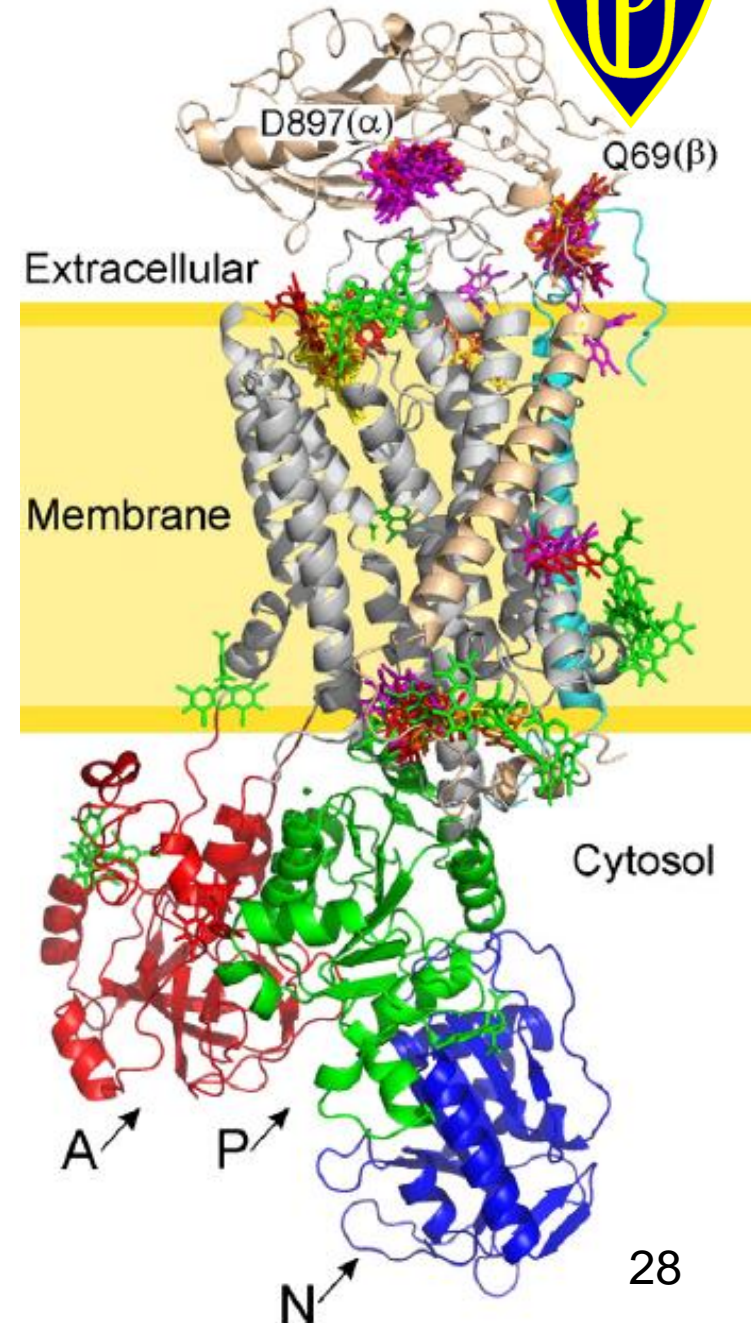
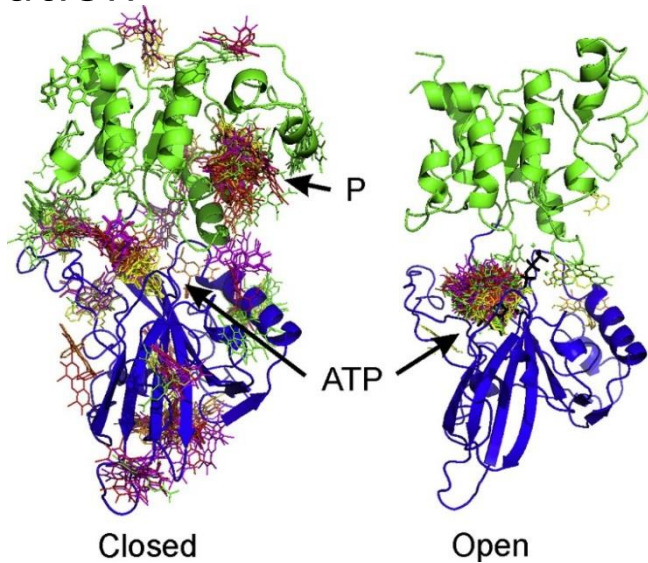


Figure 1 | **Protein mobility and ligand binding.** A protein is considered to exist in two conformations (P and P\*) with an energy difference  $\Delta G_{\text{convert}}$ . The ligand (L) can bind the protein (P) to give a complex (PL), or bind to P\* to give a complex (P\*L). Although P\* has a higher free energy, it might offer greater scope for interaction with L. For instance, P\* might represent a conformer in which the binding site has opened and exposed hydrophobic patches. This is energetically unfavourable, but offers the potential for favourable interactions with the hydrophobic moiety of a suitable incoming L, thereby giving rise to a large, favourable interaction  $\Delta G_{\text{intrinsic}}$ . The resulting complex (P\*L) has a lower energy than that of the complex PL. The observed affinity of L for the protein conformational ensemble is governed by  $\Delta G_{\text{obs}}$ . Slow binding kinetics might well be observed, as P\* is a higher-energy conformer than P and an energy barrier ( $\Delta G_{\text{barrier}}$ ) must be surmounted before optimal binding to L can take place.

# Example 1: Na<sup>+</sup>/K<sup>+</sup>-ATPase



- Ion pump
- Search for binding site
  - Fluorescent probes
  - RH241 probe
- Docking is highly sensitive to protein conformation



Havlikova M, ... **Berka K**, ... et al. *BBA*, 1828(2), 568, 2013

Huličiak M, ... **Berka K**, ... et al. *submitted*, 2014

# Protein Conformations

- **Rigid Receptor Approximation**

- Most docking programs use rigid receptor for speed

- but...

- Protein can deform in order to accept several ligands  
**(ligand-induced fit)**

- Amino acids - several conformations

- **Flexible Receptor docking**

- Increase of search size – higher computational cost

1. Side chains only

(docking selected sidechains together with ligands)

2. Docking into several structures of protein

Larger movements can be taken into account

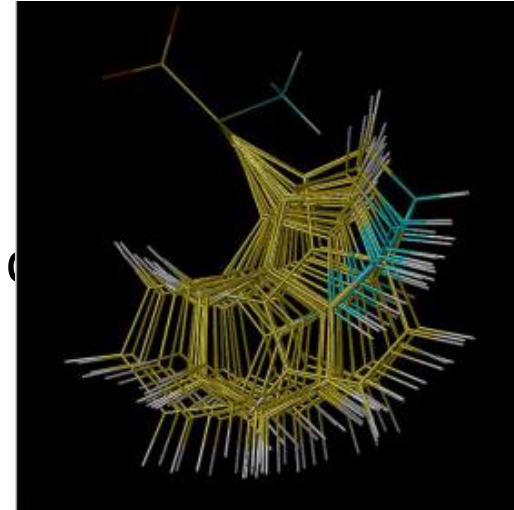


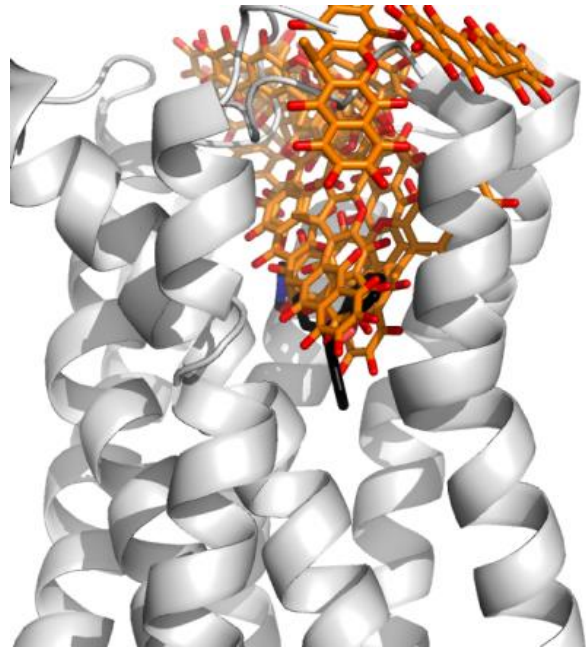
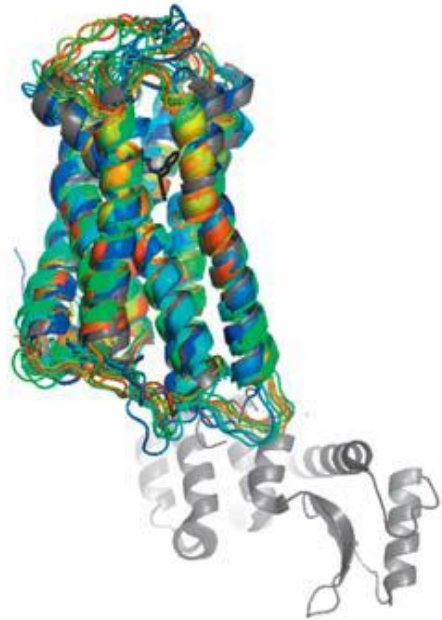
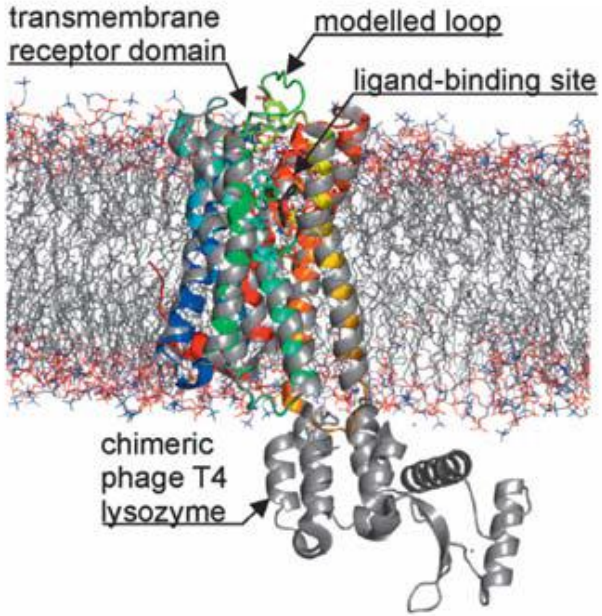
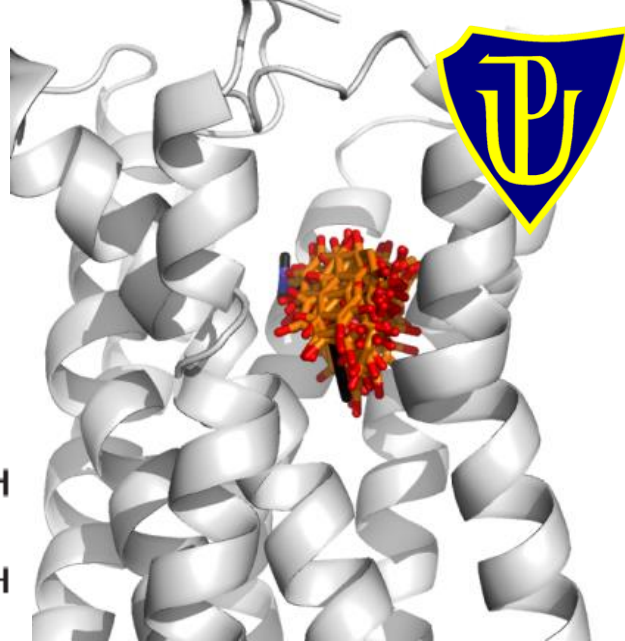
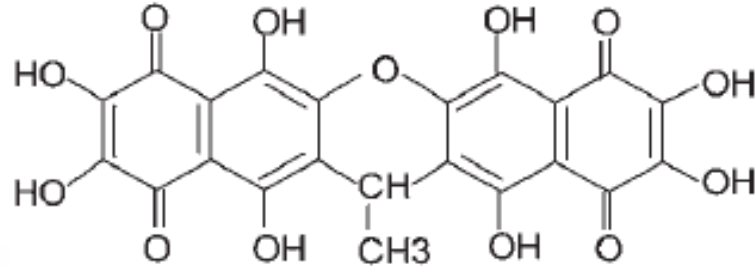
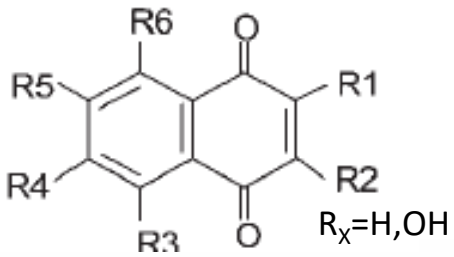
Image: Cláudio M. Soares, Protein Modelling Laboratory,  
<http://www.itqb.unl.pt/labs/protein-modelling/activities/psccip-pf>

# Example 2: H1R receptor



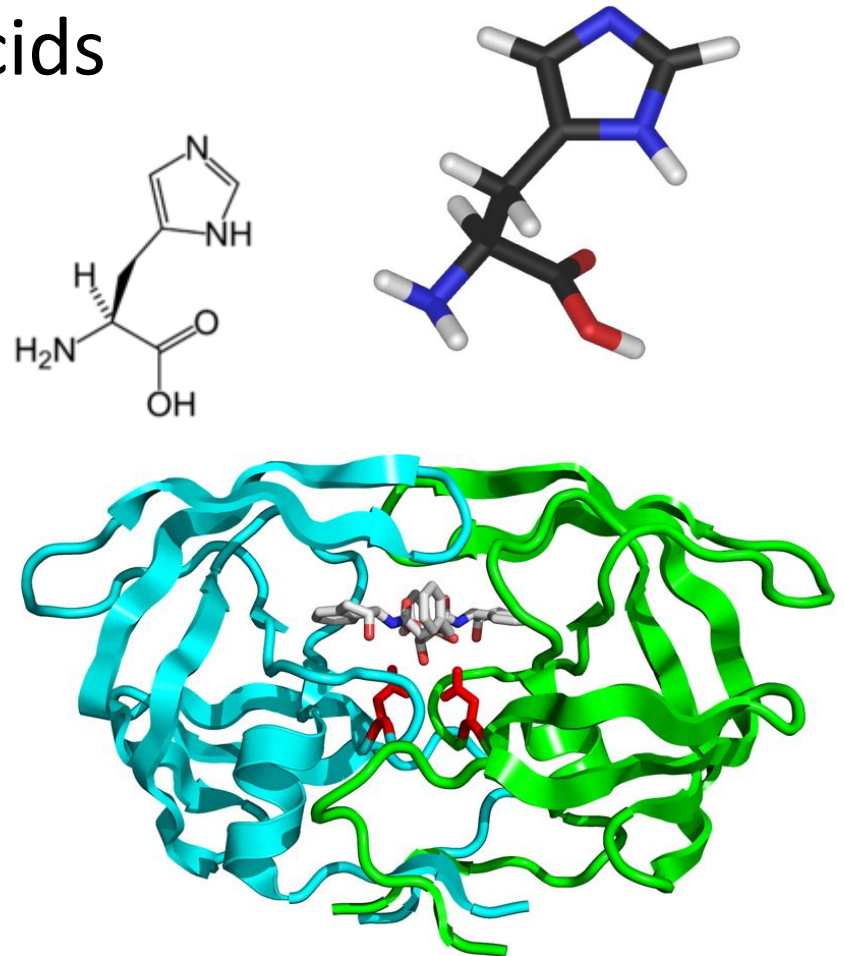
- Antiallergic compounds

dG/n(atoms) – monomers are more active than dimers



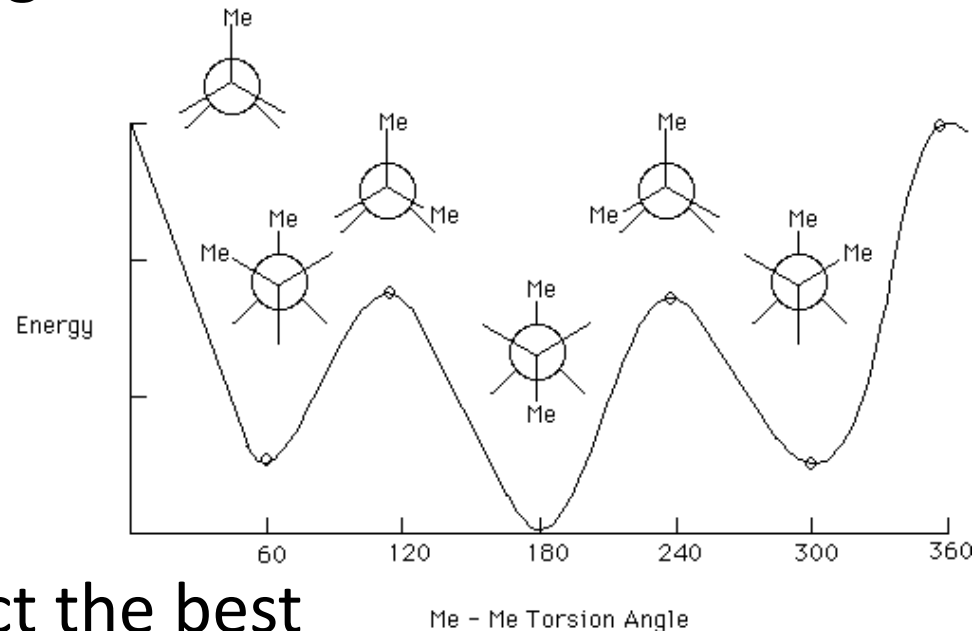
# Receptor Preparation

- protonation of aminoacids
  - His (pKa ~ 6.04)
  - Surroundings pKa shifts (Asp in HIV protease)
- tautomerization
- rotamers
- pre-selection change results significantly



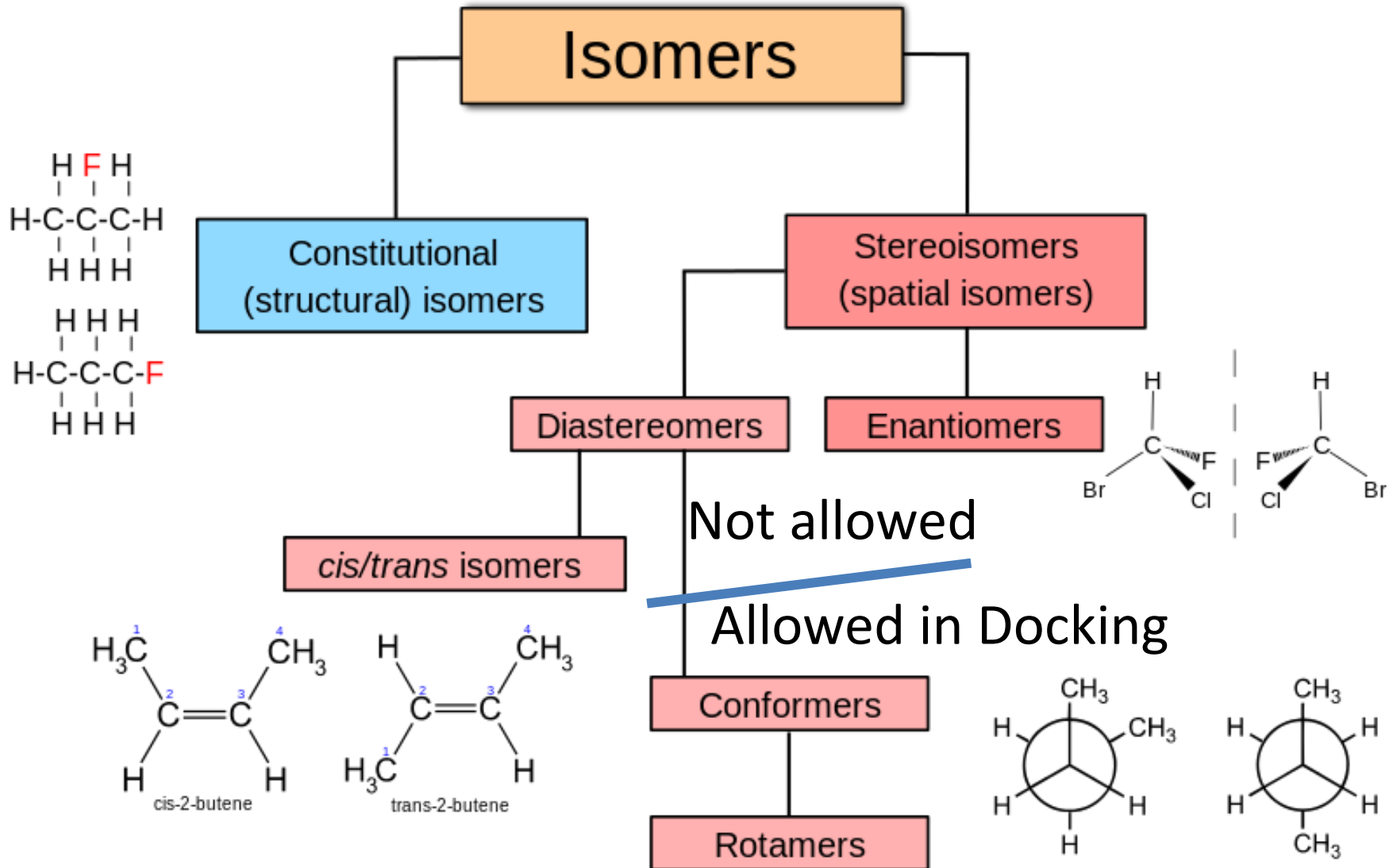
# Ligand Preparation

- Ligand Flexibility
  - Ensemble of all possible ligand conformations
  - rotation C-C bonds, but not C=C or rings
  - Angles and bonds fixed
- Isomerization
  - Charge and tautomers
  - Prepare all and then select the best
    - Relative energy
  - Ask an expert! (organic chemists)





# Isomers



# Ligand conformation

- Conformation – rotation around torsion angles
  - N rotational bonds – rotate by  $\theta$  degrees ( $5^\circ$ )
  - Conformations:  $(360^\circ / \theta)^N$
- Question
  - If the torsion angles are incremented in steps of  $30^\circ$ , how many conformations does a molecule with 5 rotatable bonds have, compared to one with 4 rotatable bonds?
- Having too many rotatable bonds results in “combinatorial explosion”
- Also ring conformations

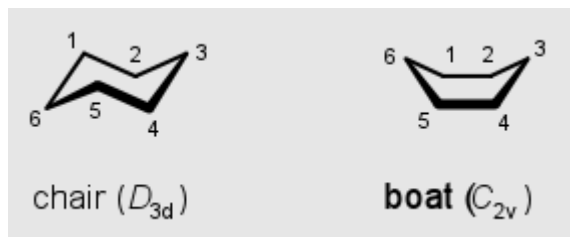
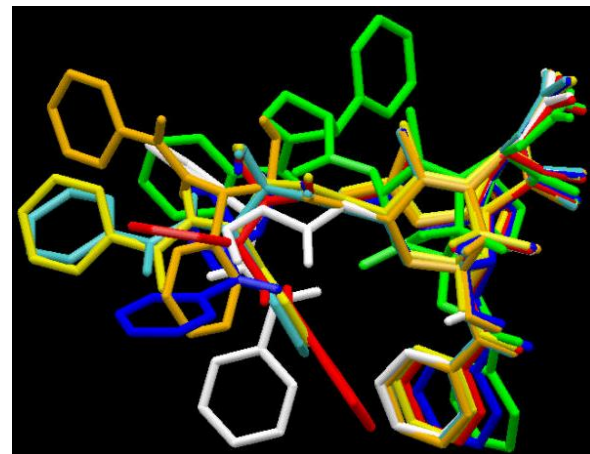
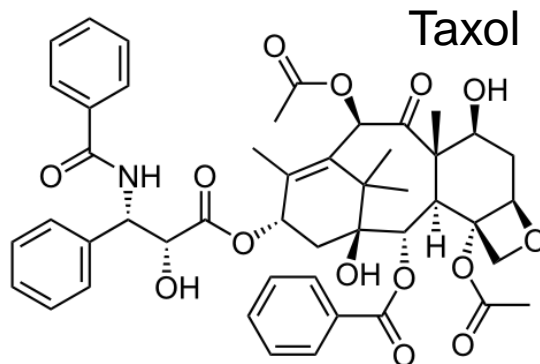


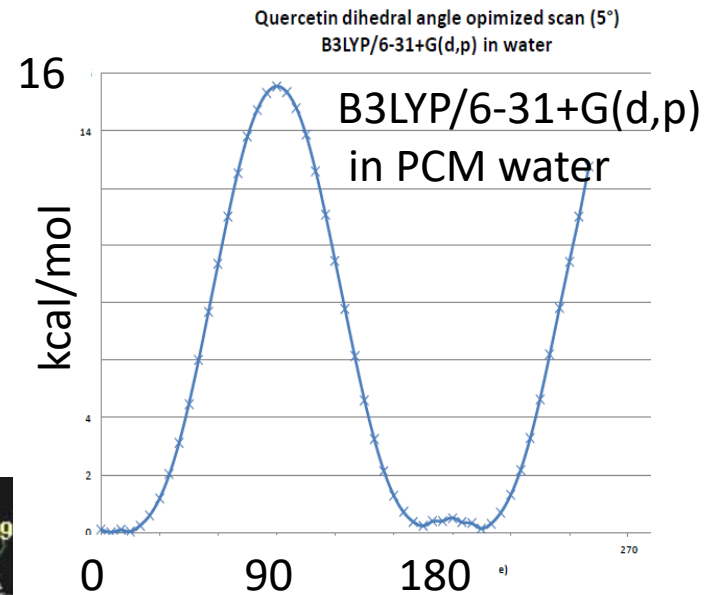
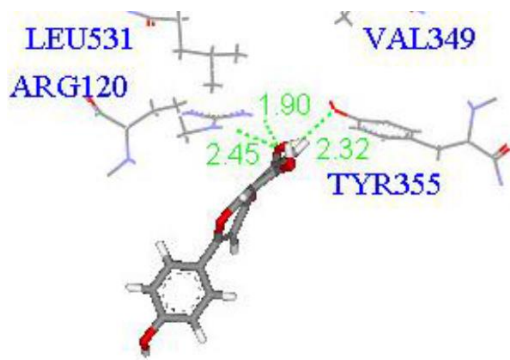
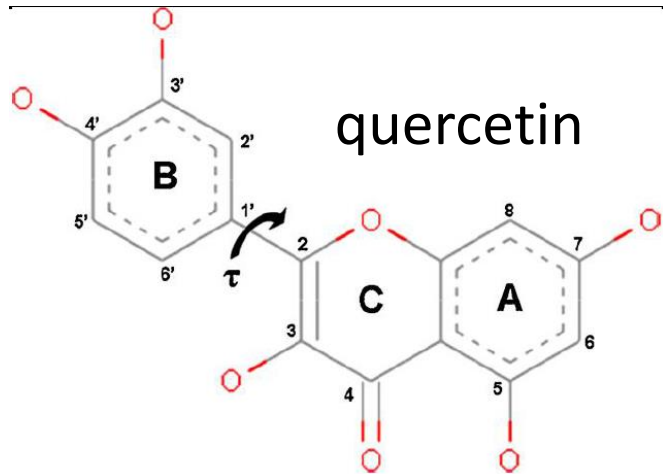
Image: IUPAC Gold Book



Lakdawala *et al.* *BMC Chemical Biology*  
2001 1:2

# Ligand Structure Generation

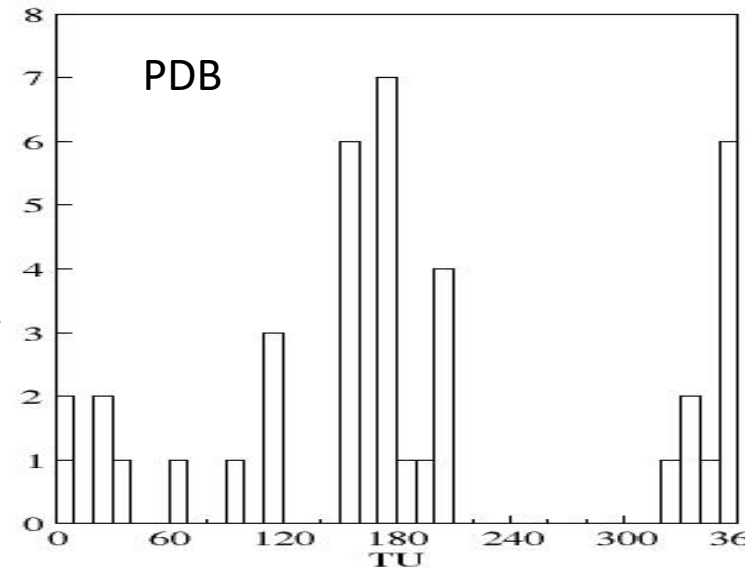
## Torsion angles



**Accelrys DSmodelling 1.2** Wu,  
Chien-Ming et al *Int.J. Mol. Sci.*  
8 (2007): 830–841.

**GLIDE**  
D'mello, P et al *Int.J.Pharm. Sci.*  
3 (2011): 33–40.

LigandFit, FlexX, DOCK 6.0, Autodock 4.0, MOE,  
Discover in Insight II, FlexiDock, Gold 3, ...



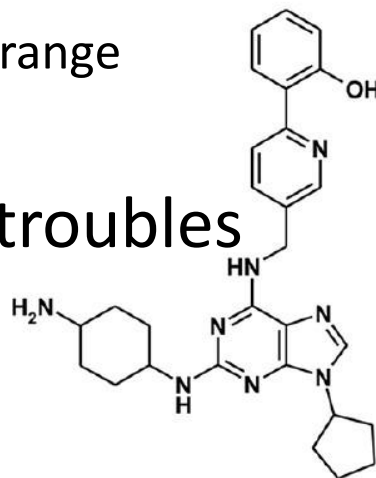
# Example: CDK2 kinase



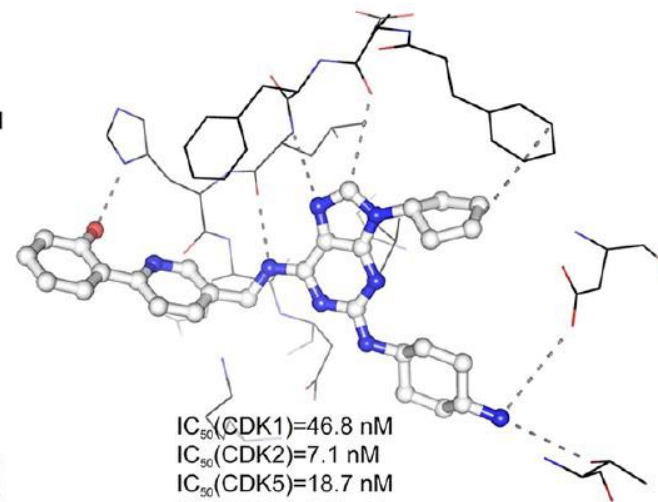
- Cell cycle regulation
  - Result: Inhibitors of CDK2 in nM range
  - Autodock Vina – speed

## – Ligand conformational troubles

(planar NH close to aromatic ring, torsional angle of biphenyl moiety)

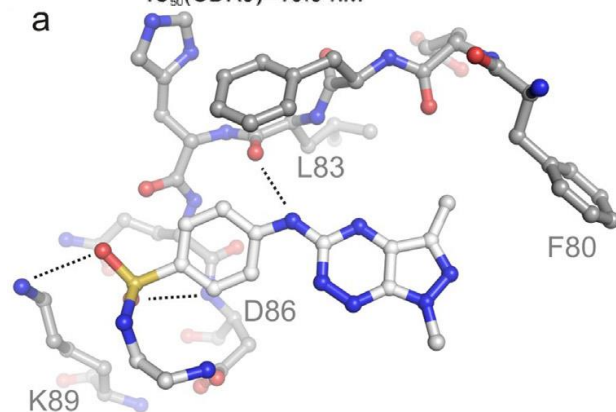
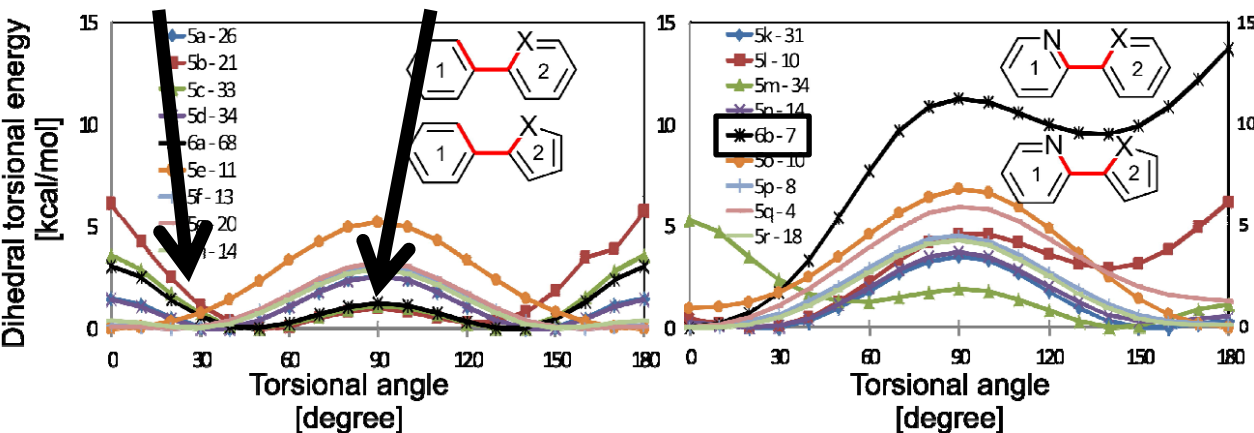


Compound 6b



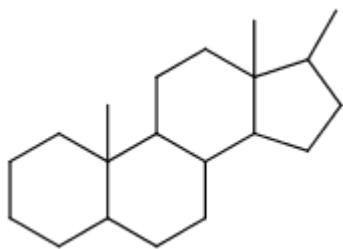
IC<sub>50</sub>(CDK1)=46.8 nM  
 IC<sub>50</sub>(CDK2)=7.1 nM  
 IC<sub>50</sub>(CDK5)=18.7 nM  
 IC<sub>50</sub>(CDK7)=306 nM  
 IC<sub>50</sub>(CDK9)=10.9 nM

## DFT scan vs docking programs

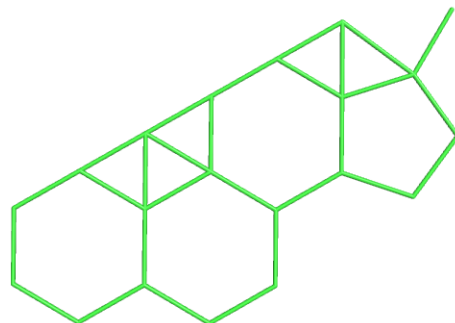


# Ligand Structure Generation Stereochemistry

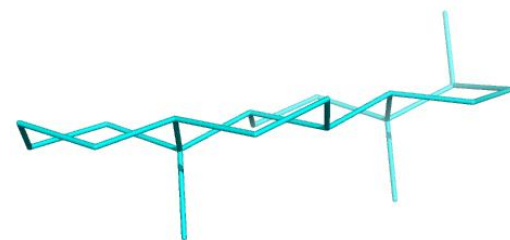
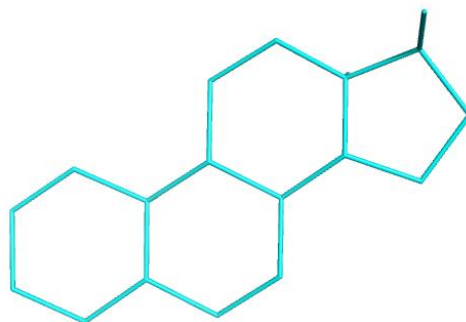
Experimentalist drawing



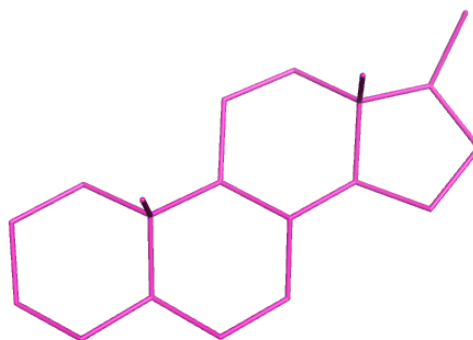
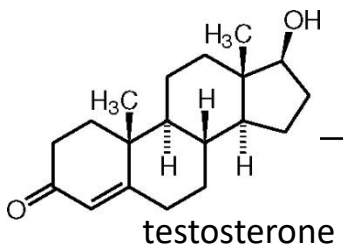
SDF



3D optimization



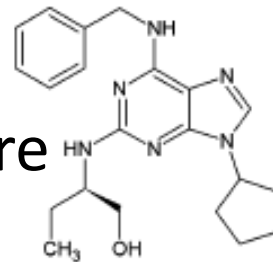
right stereoisomer



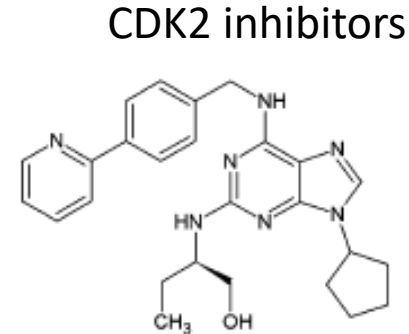
# Ligand Structure Generation

## Ligand size

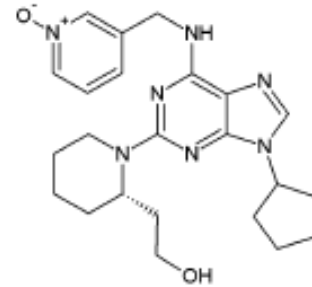
- Ligands in series
  - Typically similar size –  
easy comparison by score  
(e.g. Vina  $dG_{\text{bind}}$ )



R-roscovitine



R-CR8



dinaciclib

- Ligands in library

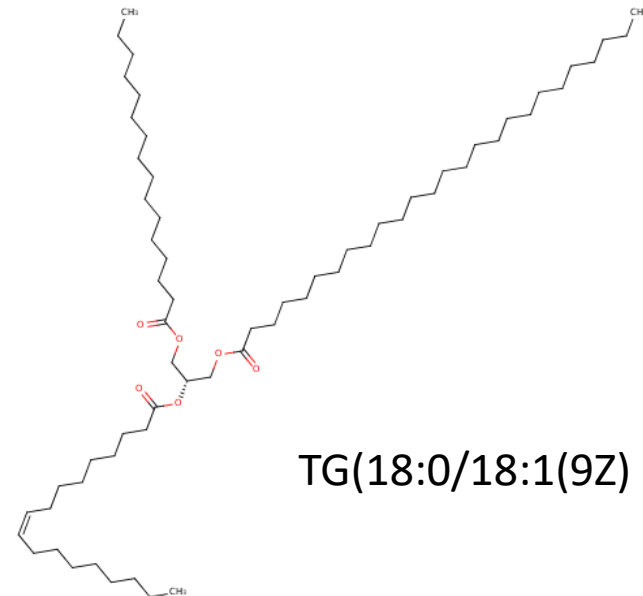
- Diverse sizes

ZINCdb – propane (MW = 44 Da)

- TG(18:0/18:1(9Z)) (MW = 1000 Da)

Larger ligand = more interactions =  
stronger binding (albeit artificially -  
entropy)

- Size incorporating measures



TG(18:0/18:1(9Z))

# Programs For Docking

- **DOCK** (I. D. Kuntz, UCSF)
- **AUTODOCK** (Arthur Olson, The Scripps Research Institute)
- **Vina** (Arthur Olson, The Scripps Research Institute)
- RosettaDOCK (Baker, Washington Univ., Gray, Johns Hopkins Univ.)
- ArgusLabs
- **GOLD**
- **FlexX**
- Hex
- Glide (Schrodinger)

# Virtual Screening

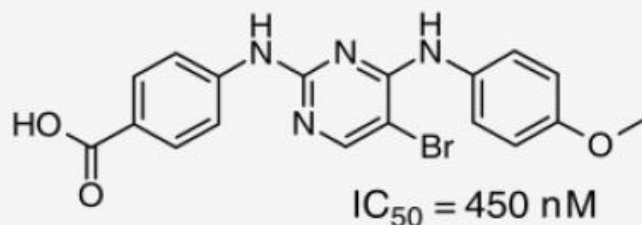


- Equivalent of biological screening (HTS – high throughput screening)
- Testing of thousands of compounds in silico
  - For further testing
  - For lead optimization
  - For leading organic synthesis



# Virtual Screening Examples

## Checkpoint kinase 1



560 000 slouč.

Drug-like compound selection

199 000 slouč.

FlexX-Pharm docking

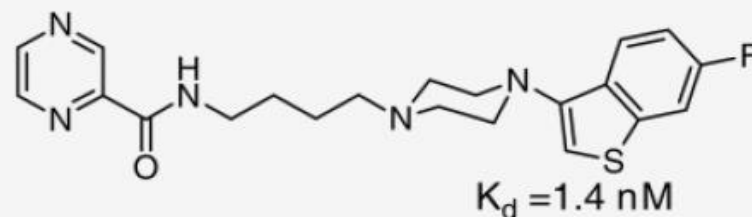
250 sloučenin

Hand picked

103 tested, 36 successful

J. Med. Chem. 47, 1962 (2004).

## $\alpha_{1A}$ adrenergic receptor



mnoho sloučenin

Drug-like compound selection

22 950 slouč.

GOLD docking

300 sloučenin

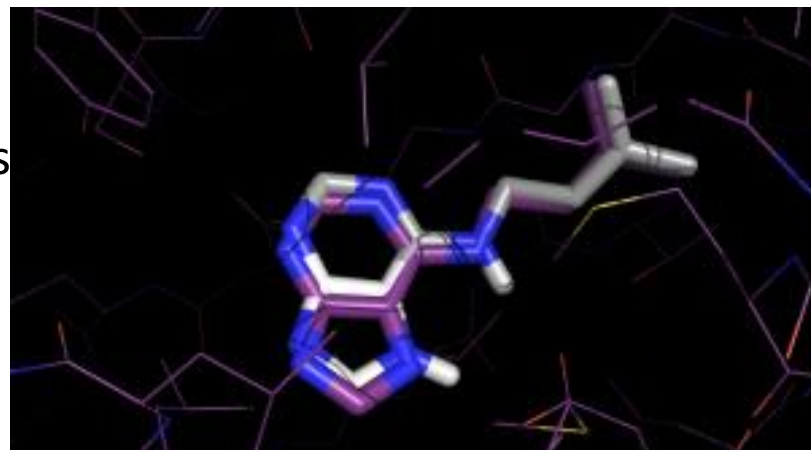
Statistical analysis

80 tested, 37 successful

J. Med. Chem. 48, 1088 (2005).

# Quality control

- Redocking (back to Xtal)
  - RMSD < 2Å
  - flexible ligand docking ~70%
  - test for scoring functions/docking programs
- Correlation plot ( $r^2 > 0.5$ )
  - $\Delta G_{\text{eff}}$
- test sets – validation
  - GOLD test set, Astex set
  - decoys – ZINC, DUD (similar phys-chem., different structures)
- Virtual Screening
  - Enrichment factor
  - (BED)ROC curves



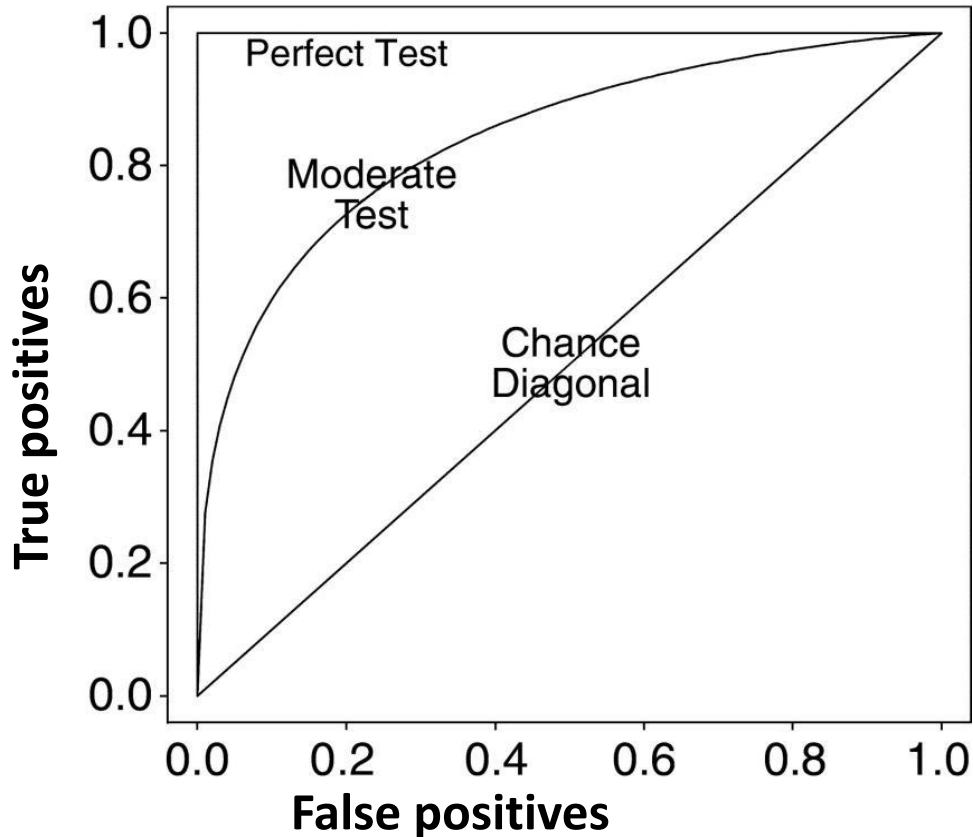
$$EF = \frac{a/n}{A/N}$$

- top (e.g. top10)  
a – active  
n - total  
- overall

$$\Delta G_{\text{eff}} = \Delta G_{\text{eff}} / N_{\text{nonHatoms}}$$

# ROC curve

- Receiver operator characteristic curve
- – signal to noise ratio



# Sorting Quality

- ROC
  - receiver operating characteristic curve
- AUAC
  - area under the accumulation curve
- average rank of actives
- EF
  - enrichment factor
- RIE
  - robust initial enhancement
- BEDROC
  - Boltzmann-enhanced discrimination of receiver operating characteristic

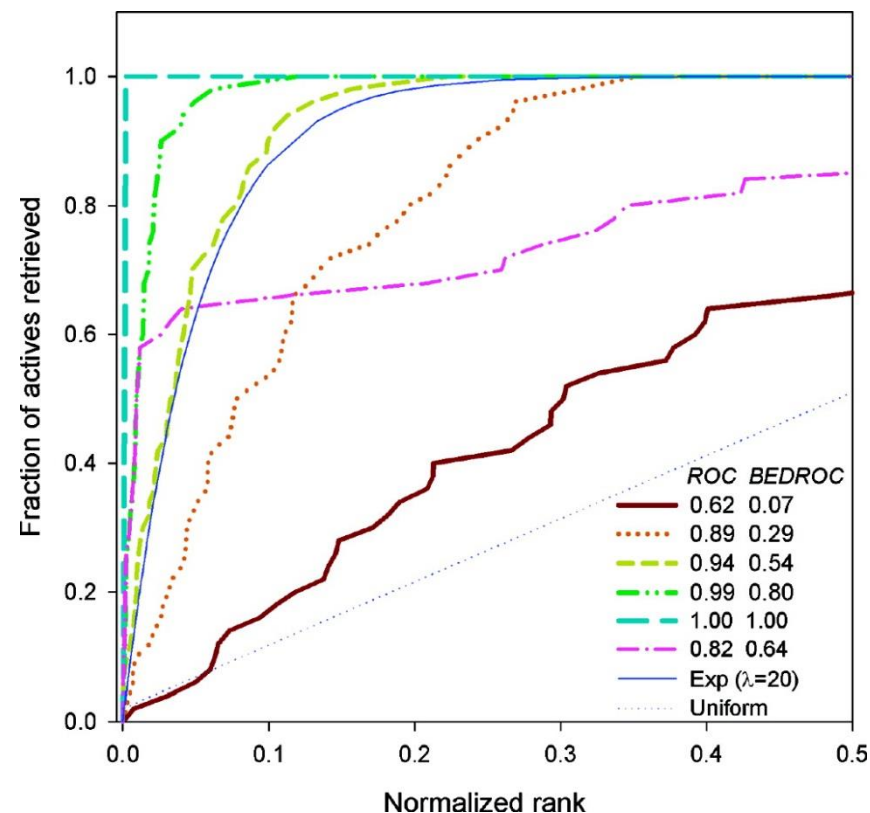


Figure 9 Different accumulation curves from sampling ( $n = 50$ ,  $N = 25000$ ) shown together with the corresponding ROC and BEDROC values where  $\alpha = 20.0$ . An exact CDF with  $\lambda = 20$  is also shown to highlight the fact that the BEDROC metric returns a value of 1/2 for a curve close to this CDF.

# Docking Summary



- Usable in SAR (structure-activity relationship)
  - explore the interactions between ligands and receptor
  - can lead drug development
- Troubles
  - Ligand preparation – 3D generation, torsion angles
  - Receptor preparation – protein flexibility
  - Scoring function – identification of right binding pose, size of ligand issue

# Tutorial preparation

Programs installation:

- Python 2.7

<http://www.python.org/download/releases/>

- Pymol

<http://www.lfd.uci.edu/~gohlke/pythonlibs/#pymol>

- Autodock Tools

<http://mgltools.scripps.edu/>

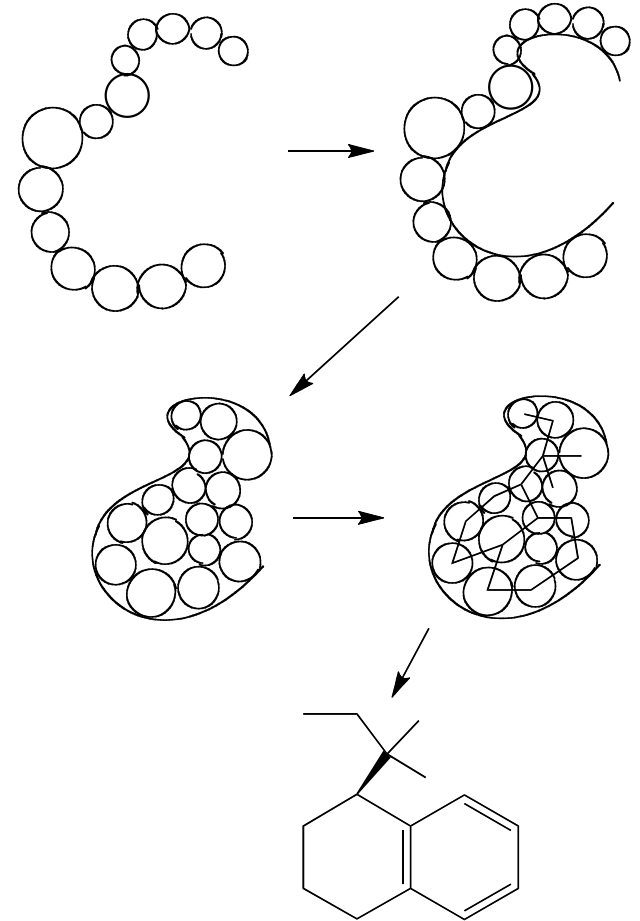
- Autodock Vina

<http://vina.scripps.edu/>



# de novo docking

- 1) Surface of receptor (SASA - Connelly)
- 2) “negative” of surface fill with spheres
- 3) Distances between spheres
- 4) Sphere distance to Bond distance conversion
- 5) Search in small molecule database
- 6) Selection of ligands with largest overlap with spheres
- 7) Scoring





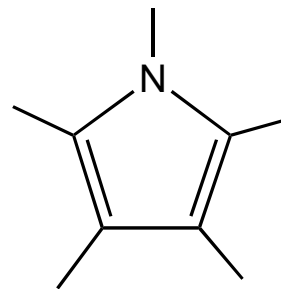
# Groupbuild

- Building of new compounds by filling of active site by random fragments

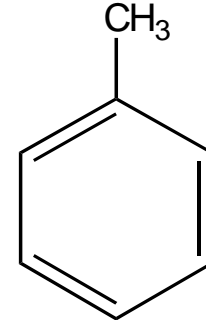
Algorithm:

- 1) Grid for receptor binding site
- 2) Structure generation
  - 1) Docking of "core" fragment
  - 2) Build up (random fragments additions to core)
  - 3) Selection of best structures
  - 4) Iterate steps 2 and 3 until final criteria fulfilled (number of steps, minimal energy, etc.)
- 3) Selection of final structures for synthesis

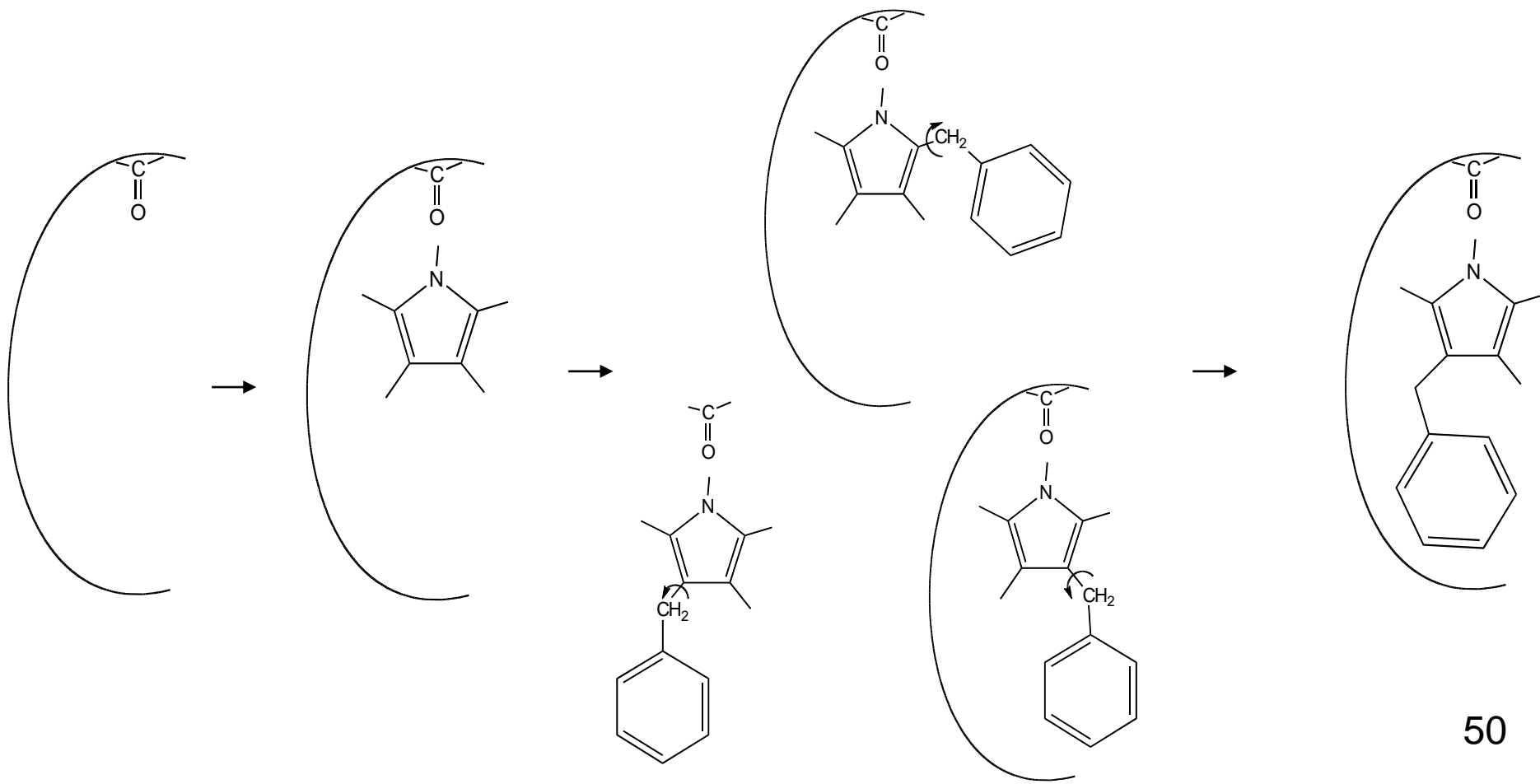
# Groupbuild



initial core fragment

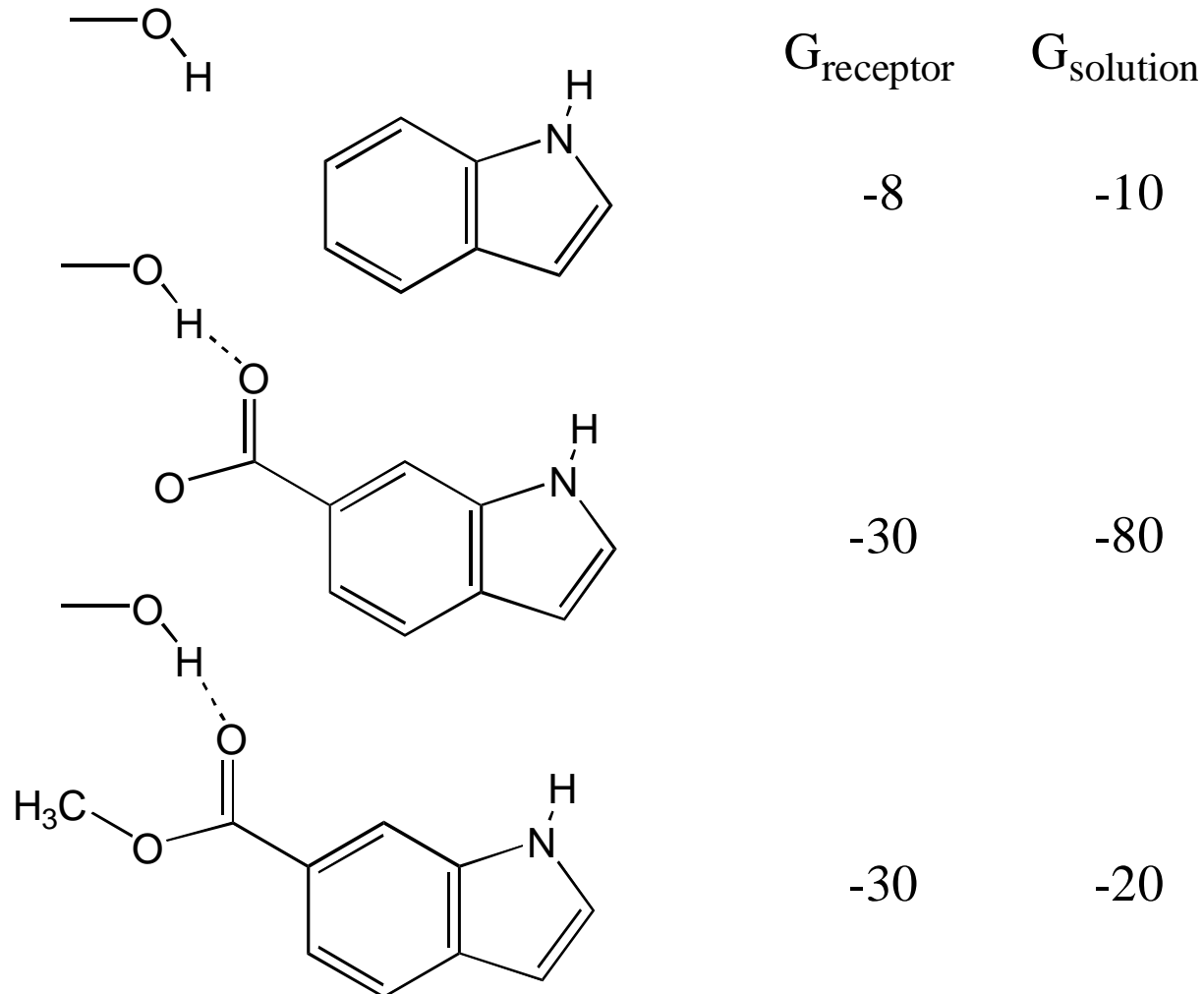


fragment to be added



# Groupbuild

## Example for Hypothetical Receptor



$$\Delta\Delta G = \Delta G_{\text{eq},2} - \Delta G_{\text{eq},1} = \Delta G_{\text{receptor}} - \Delta G_{\text{solution}}$$