

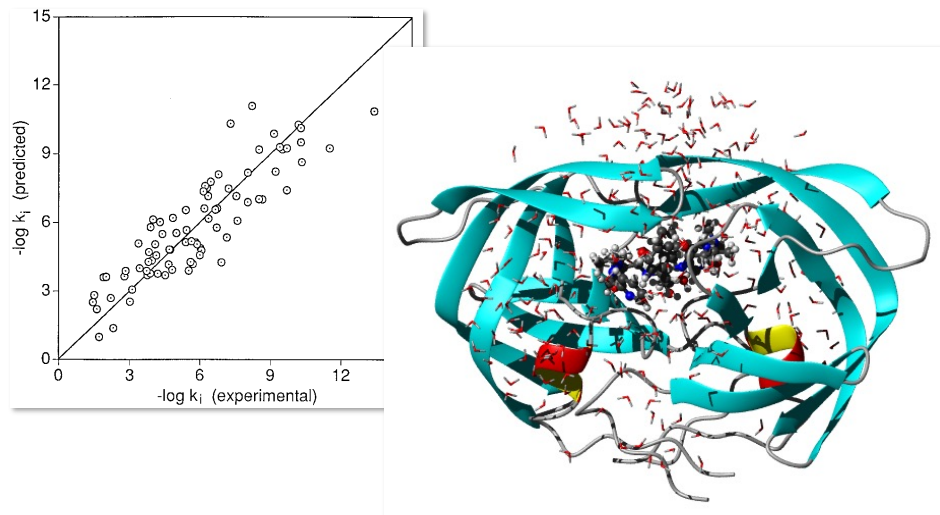
1

Lectures & Practices Agenda

Session	Lecture	Practice
1	Prologue: molecular representation	
	Introduction to (computer-aided) drug design	
	Origin of 3D structures	
	Molecular recognition	Use of UCSF chimera to analyze protein-ligand complexes
2	Binding free energy estimation	
	Introduction to molecular docking	Ligand-protein docking with AutoDock Vina
3	Introduction to molecular (virtual) screening	Ligand-based virtual screening with SwissSimilarity
4	Short introduction on target prediction of small molecules	Use of SwissTargetPrediction to perform reverse screening.
5	Introduction to ADME, pharmacokinetics, druglikeness	Estimate physicochemical, pharmacokinetic, druglike and related properties with SwissADME
6	Short introduction to bioisosterism	Use of SwissBioisostere to perform bioisosteric design

2

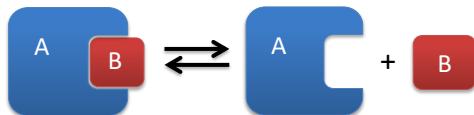
Binding free energy estimation



3

Binding free energy – Link between in silico and experimental worlds

Link between experiment and modeling



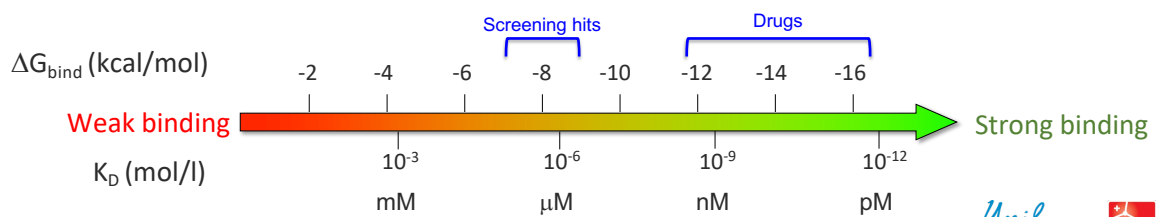
$$K_D = \frac{[A][B]}{[AB]}$$

K_D : dissociation constant

Accessible by
computer-
aided methods

Experimentally
measured

$$\Delta G_{\text{bind}} = RT \ln(K_D) = \Delta H - T\Delta S$$



4

Binding free energy – The computational methods

Ligand-based

- **Machine Learning approaches:**
- 2D QSAR. Ex: Hansch equations
- 3D QSAR. Ex: CoMFA

Structure-based

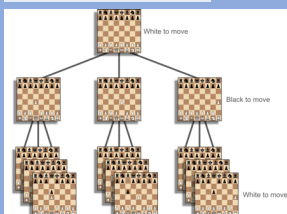
- **Force field** methods:
 - Free energy simulation (FEP, TI)
 - MM-PBSA, MM-GBSA
 - Linear interaction energy (LIE)
- **Empirical** scoring functions (regression based approaches). Ex: LUDI score
- **Knowledge**-based approaches (Potential of Mean Force). Ex: PMF score

5

Binding free energy – The computational methods – Artificial intelligence?

Artificial Intelligence

```
if(object_ahead) then:
  turn()
else:
  go_ahead()
```

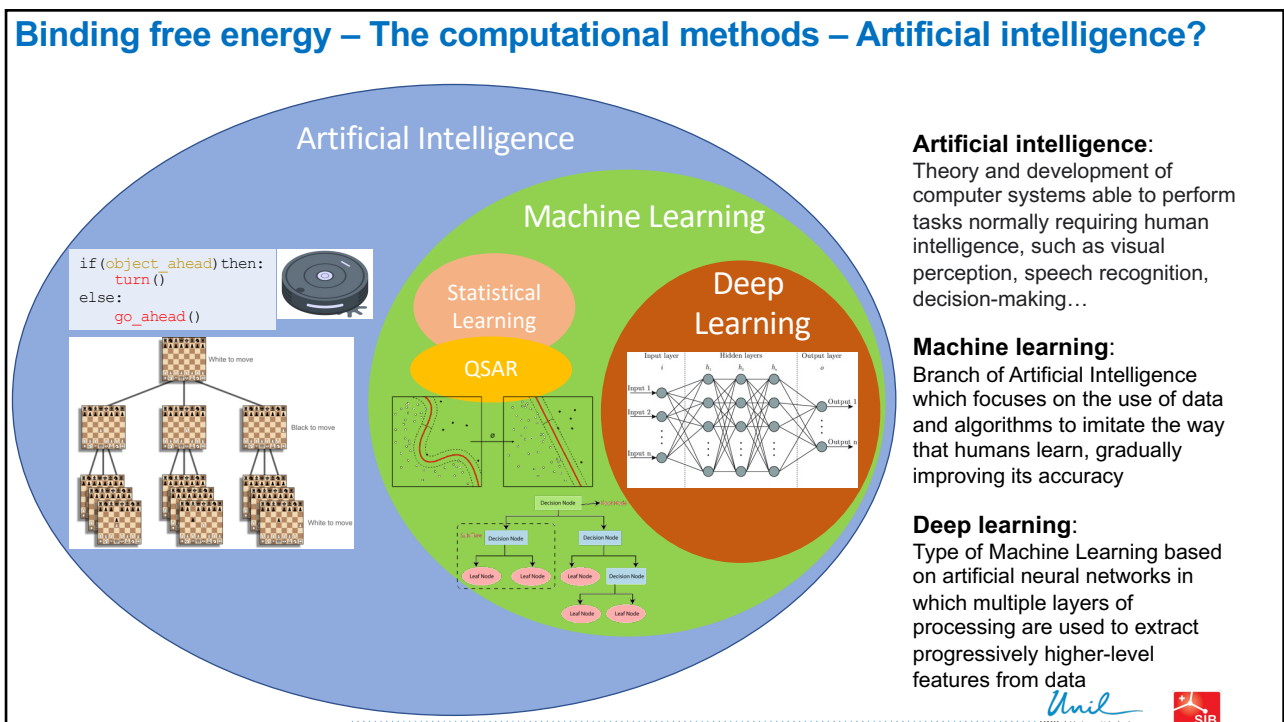


Artificial intelligence:

Theory and development of computer systems able to perform tasks normally requiring human intelligence, such as visual perception, speech recognition, decision-making...

6

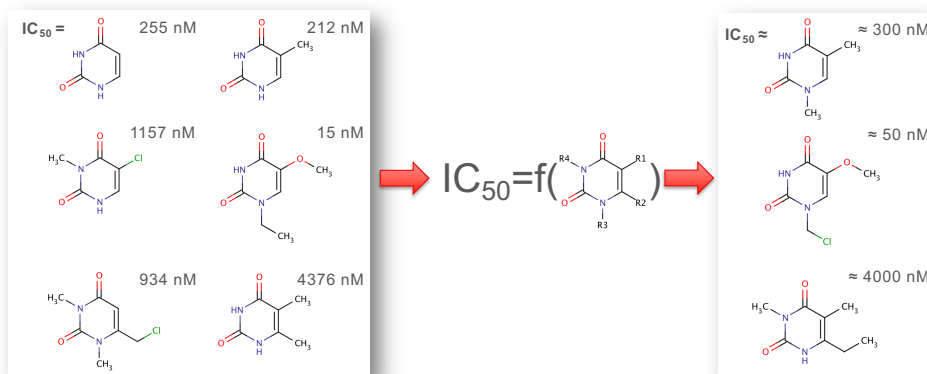
Binding free energy – The computational methods – Artificial intelligence?



Binding free energy – Ligand-based - QSAR

Quantitative Structure-Activity Relationship :

1. Built a set of **molecules** with **known experimental affinities** (activities).
2. Define a **mathematical relationship** between the **structure** (properties) of molecules and their **activities**.
3. Use this equation to **predict the binding affinities** (activity) of new molecules.



Unil



9

9

Binding free energy – Ligand-based - QSAR

Quantitative Structure-Activity Relationship :

1. Built a set of **molecules** with **known experimental affinities** (activities).
2. Define a **mathematical relationship** between the **structure** (properties) of molecules and their **activities**.
3. Use this equation to **predict the binding affinities** (activity) of new molecules.

Assumptions

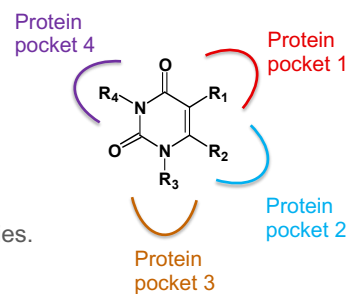
Chemical similarity of ligands → Similarity of biological response
 Affinity is a function of the differences in the ligand properties.
 The binding mode is similar.

Advantages

No need for structural information about the target
 Once established, extremely quick calculation. Suitable for very large libraries.

Limitations

The set of molecules need to be **large**, with a broad **spread of activity**.
 Limited to **structurally similar** molecules (applicability domain).



Unil

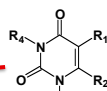


10

10

Binding free energy – Ligand-based – 2D QSAR

n structurally related molecules
characterized experimentally
(i.e. the training set)



Quantitative description

Measured activities

	X_1	X_2	...	X_m	ΔG_{bind}
1					
2					
3					
...					
n					

MLR

$$\Delta G_{\text{bind}} = k_0 + \sum k_i X_i$$

Multiple Linear Regression, for instance.

**Global or substituent “2D”
Descriptors (X_i):**

- molecular weight
- $\log P$
- polar surface area (PSA)
- simple count of atoms
- electronegativity
- charges

**Rules for statistical relevance
and model quality:**

- At least 5 molecules per descriptor
- Descriptors should **not** be intercorrelated
(should not contain the same information)

C. Hansch and T. Fujita, JACS, 1964, 86, 1616

Unil

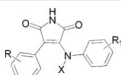


11

11

Binding free energy – Ligand-based – 2D QSAR

Ex: Probing the physicochemical and structural requirements for glycogen synthase kinase-3 α
inhibition: 2D-QSAR for 3-anilino-4-phenylmaleimides.
Sivaprakasam, P.; Xie, A.; Doerksen, R. J.; *Bioorg Med Chem* 2006, 14, 8210–8218.



Compound	R	R ₁	X	IC ₅₀ ^a	pIC ₅₀ ^b	H _A R	f _{ortho}	σ_{metaR}	I _{SCIR1}	I _{LCIR1}	I _{SOHR1}	I _{5CH3}	E _{orthoR}	π_{metaR1}
1	H	H	H	529	6.28	0	0	0	0	0	0	0	0	0
2	2-Cl	H	H	216	6.67	0	0.41	0	0	0	0	0	-0.97	0
3	2-OCH ₃	H	H	216	6.67	1	0.26	0	0	0	0	0	-0.55	0
4	3-NO ₂	H	H	141	6.85	1	0	0.71	0	0	0	0	0	0
5	4-Cl	H	H	514	6.29	0	0	0	0	0	0	0	0	0
6	4-OCH ₃	H	H	390	6.41	1	0	0	0	0	0	0	0	0
7	H	3-Cl	H	301	6.52	0	0	0	1	0	0	0	0	0.71
8	2-Cl	3-Cl	H	195	6.71	0	0.41	0	1	0	0	0	-0.97	0.71
9	2-OCH ₃	3-Cl	H	114	6.94	1	0.26	0	1	0	0	0	-0.55	0.71
10	2-NO ₂	3-Cl	H	104	6.98	1	0.67	0	1	0	0	0	-1.01	0.71

Table 4. Statistical parameters of 2D-QSAR models for 3-anilino-4-phenylmaleimides

QSAR model	n	r	r^2	F	s	q^2
1	64	0.815	0.665	23.0	0.216	0.590
2	64	0.860	0.739	26.9	0.192	0.675
3	64	0.862	0.743	27.5	0.191	0.679
4	67	0.847	0.718	31.1	0.231	0.671
5	67	0.873	0.762	32.0	0.214	0.706
6	67	0.901	0.812	36.5	0.191	0.761
7	67	0.902	0.814	36.9	0.190	0.762
8	67	0.806	0.650	39.1	0.253	0.608
9	67	0.852	0.726	41.0	0.226	0.688
10	67	0.862	0.743	35.2	0.220	0.702
11	67	0.863	0.745	35.7	0.219	0.705
12	67	0.880	0.774	41.8	0.206	0.732
13	67	0.909	0.827	47.6	0.182	0.788
13a ^a	66	0.922	0.850	55.6	0.165	0.814
14	67	0.911	0.829	48.6	0.181	0.789
14a ^a	66	0.922	0.850	55.9	0.165	0.813

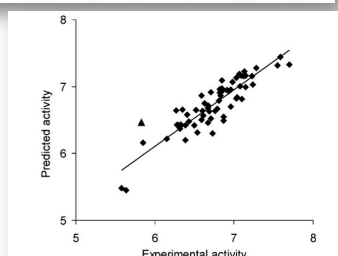


Figure 1. Correlation between the observed and predicted GSK-3 α inhibitory activity from QSAR model 14 showing compound 19 as outlier (triangle).

Unil



12

12

Binding free energy – Ligand-based – 2D QSAR

Limited to structurally related molecules

Needs the experimental activity of a series ligands

➡ Not for *ab initio* studies

Overfitting

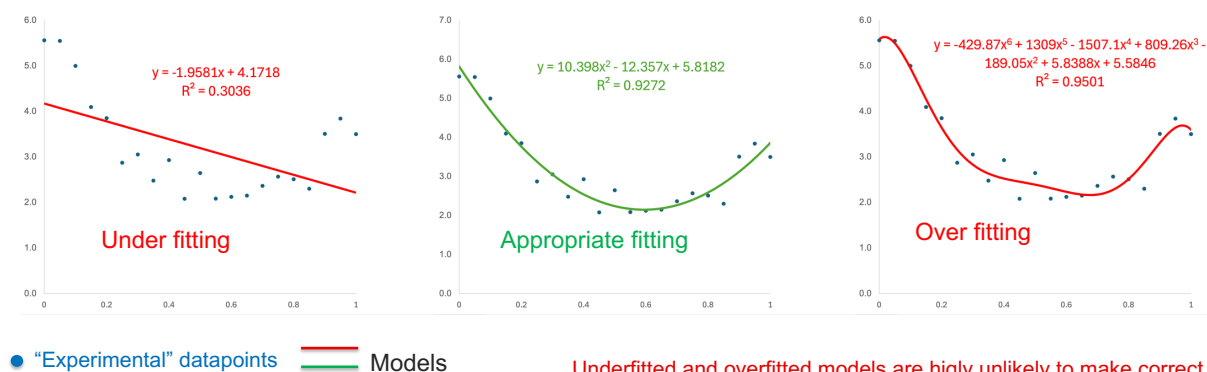
- ➡ - Method for selecting the descriptors (genetic algorithm)
- Estimation of the predictive ability (external test set, Y-randomization, Cross-validation, ...)

13

Binding free energy – Ligand-based – 2D QSAR

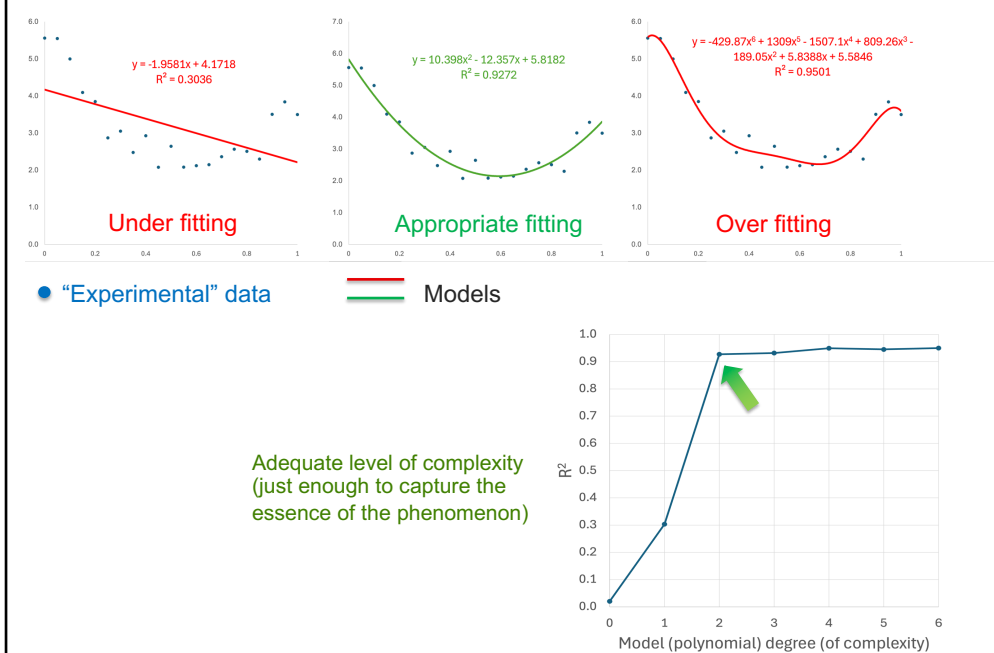
Exercise:

- Create a polynomial function $y=f(x)$, of degree 2 (i.e. a parabola), and add noise to generate 'experimental datapoints'
- Then, make models to fit the datapoints, and see how polynomial function of degrees 1, 2, 3 and more perform in terms of underfitting, appropriate fitting and overfitting



14

Binding free energy – Ligand-based – 2D QSAR



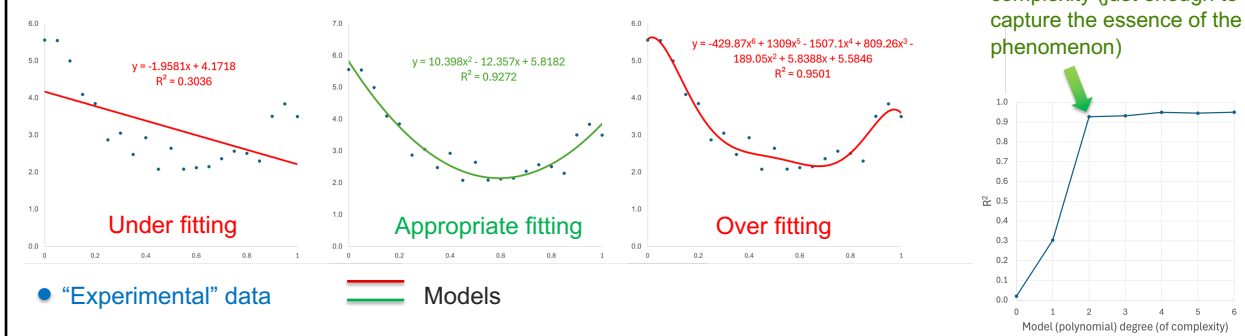
Unil

SIB

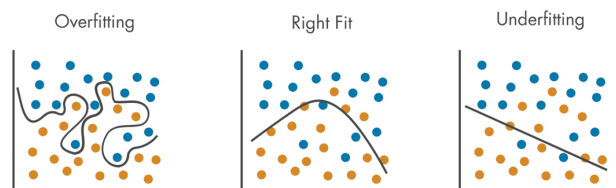
15

15

Binding free energy – Ligand-based – 2D QSAR



Also applies to classification:



Adapted from MathWorks

Unil

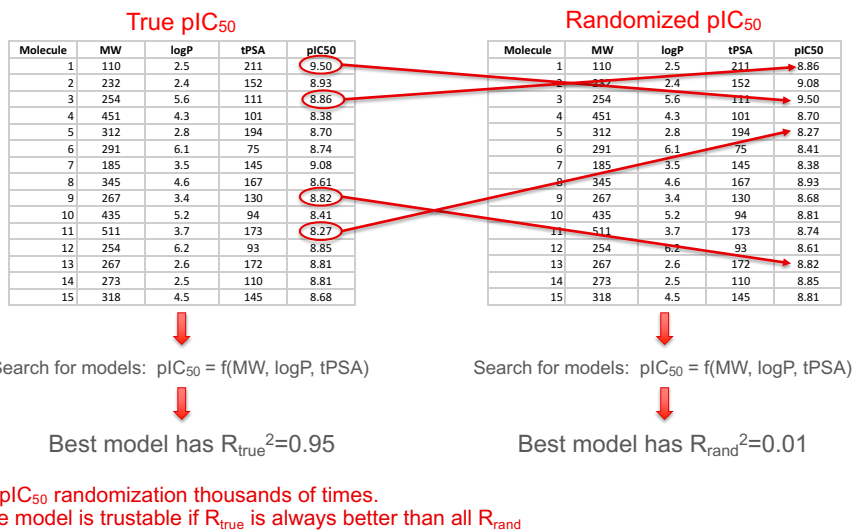
SIB

16

16

Binding free energy – Ligand-based – 2D QSAR

Principle of Y-randomization



Unil



17

17

Binding free energy – Ligand-based – 2D QSAR

Limited to structurally related molecules

Needs the experimental activity of a series ligands

➡ Not for *ab initio* studies

Overfitting

➡ - Method for selecting the descriptors (genetic algorithm)
- Estimation of the predictive ability
(external test set, Y-randomization, Cross-validation, ...)

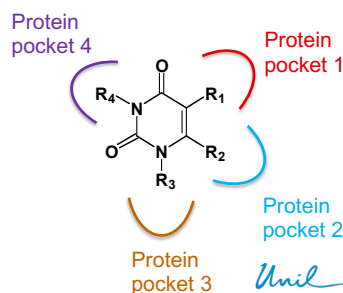
Validity domain. Interpretability strictly linked to the descriptor and training set:

If only hydrophobic groups at R₁ in the training set

➡ Influence of a hydrophilic group at R₁ ?

If only methyl, ethyl, propyl, butyl at R₁ in the training set

➡ Contribution of pentyl, hexyl, etc... ?



Unil



18

18

Binding free energy – Ligand-based – 3D QSAR

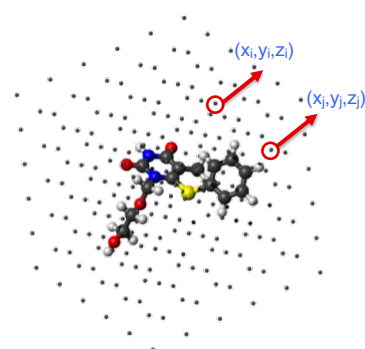
Example : Comparative molecular field analysis (CoMFA)

R.D. Cramer *et al.*, JACS, 1988, 110, 5959

Molecules superposed in a 3D grid

“3D” Descriptor = Molecular Fields

	(x_1, y_1, z_1)	(x_2, y_2, z_2)	(x_1, y_1, z_1)	(x_2, y_2, z_2)		
	S_1	S_2	...	E_1	E_2	...
1						
2						
3						
...						
n						
	Steric field (Lennard-Jones)			Electrostatic field (Coulomb)		
	ΔG_{bind}					



MLR $\rightarrow \Delta G_{\text{bind}} = k_0 + \sum \alpha_i S_i + \sum \beta_i E_i$ **Multiple Linear Regression**, for instance.

Unil

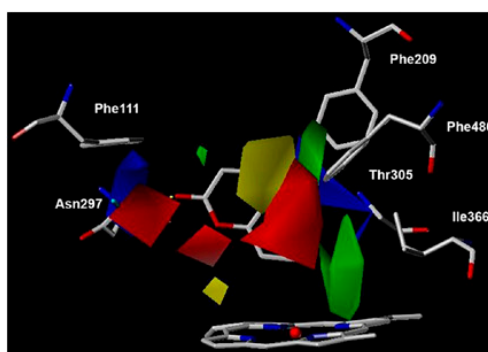


19

19

Binding free energy – Ligand-based – 3D QSAR

“inverse” image of the receptor binding site



CoMFA **electrostatic** fields:

- **blue**, negative-charge disfavored area
- **red**, negative-charge detrimental area

CoMFA **steric** field:

- **green**, bulk favorable area
- **yellow**, bulk detrimental area

Identification of inhibitors of the nicotine metabolising CYP2A6 enzyme: an *in silico* approach
M. Rahnasto, C. Wittekindt, R. O. Juvonen, M. Turpeinen, A. Petsalo, O. Pelkonen, A. Poso, G. Stahl, H-D Hotje and H. Raunio
Nature, 2008, 8(5), 328-338

Unil



20

20

Binding free energy – Ligand-based – 3D QSAR

Requires the **experimental activity** of a series of ligands

Risks of overfitting

Needs to respect the **domain of validity** when used for prediction

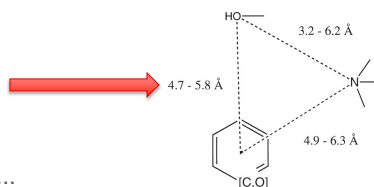
Not limited to structurally related molecules

Main limitation: **Alignment** of the molecules in their (guessed) **bioactive conformation**.

Possible help of:

- Structure of a protein-ligand complex available
→ alignment over cocrystallized ligand or by docking.
- Set including rigid molecules
→ alignment over rigid molecules
- Functional groups in agreement with a pharmacophore hypothesis
→ alignment over pharmacophoric points.

Others : CoMSIA, HASL, Compass, APEX-3D, YAK, ...



Unil



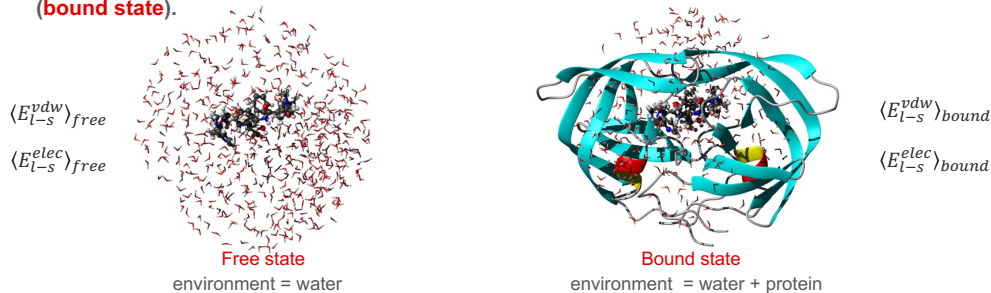
21

21

Binding free energy – Structure-based – Force field (Physics-based)

Example: Linear Interaction Energy (LIE)

- compute the **van der Waals (Lennard-Jones)** and **electrostatic interaction (force field) energies** of the ligand with water (**free state**) and of the ligand with protein and water (**bound state**).



$$\Delta G_{\text{bind}} = \alpha \left(\langle E_{l-s}^{vdw} \rangle_{\text{bound}} - \langle E_{l-s}^{vdw} \rangle_{\text{free}} \right) + \beta \left(\langle E_{l-s}^{elec} \rangle_{\text{bound}} - \langle E_{l-s}^{elec} \rangle_{\text{free}} \right)$$

J. Åqvist, *J. Phys. Chem.*, **1994**, 98, 8253
 $\alpha=0.165$ and $\beta=0.5$

T. Hansson *et al.*, *J. Comp.-Aided Molec. Design*, **1998**, 12, 27
 $\alpha=0.181$ and $\beta=0.5, 0.43, 0.37, 0.33$

W. Wang, *Proteins*, **1999**, 34, 395
 α function of binding site hydrophobicity

Unil



22

22

Binding free energy – Structure-based – Force field (Physics-based)

Example: Linear Interaction Energy (LIE)

Modifications :

- Additional term proportional to buried surface upon complexation
D.K. Jones-Hertzog and W.L. Jorgensen, *J. Med. Chem.*, **1997**, *40*, 1539
- Use of continuum solvent model instead of explicit solvent
R. Zhou and W.L. Jorgensen *et al.*, *J. Phys. Chem.*, **2001**, *105*, 10388
- Replace molecular dynamics simulations by simple minimization
Huang, D.; Caffisch, A. *J. Med. Chem.* **2004**, *47*, 5791

Advantages :

- Can treat more structurally different ligands than QSAR. But still generally restricted to rather similar ligands.

Shortcomings:

- Slower than scores based on a single conformation (LUDI, PMF, ...)
- Not really universal (α and β are system-dependent)
- Need experimental binding affinities of known complexes

Unil

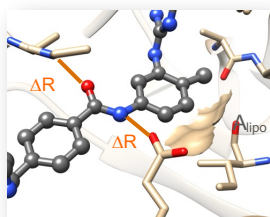


23

23

Binding free energy – Structure-based – Empirical methods

Example: LUDI score



Evaluation of ΔG_{bind} : a simple **count** of various type of interactions between ligand and protein.

$$\Delta G_{\text{bind}} = \Delta G_0 + \Delta G_{\text{polar}} + \Delta G_{\text{apolar}} + \Delta G_{\text{solv}} + \Delta G_{\text{flexi}}$$

Developed using a 82 protein-ligand complexes dataset with known experimental ΔG_{bind}

Polar interactions

$$\begin{aligned} \Delta G_{\text{polar}} = & \Delta G_{\text{hb}} \sum_{\text{hb}} f(\Delta R, \Delta \alpha) \times f(N_{\text{neighb}}) \times \text{fpcs} \\ & + \Delta G_{\text{ion}} \sum_{\text{ion}} f(\Delta R, \Delta \alpha) \times f(N_{\text{neighb}}) \times \text{fpcs} \\ & + \Delta G_{\text{esrep}} N_{\text{esrep}} \end{aligned}$$

$$\Delta G_{\text{hb}} = -0.81, \Delta G_{\text{ion}} = -1.41 \text{ and } \Delta G_{\text{esrep}} = +0.10 \text{ kcal/mol}$$

Apolar interactions

$$\Delta G_{\text{apolar}} = \Delta G_{\text{lipo}} A_{\text{lipo}} + \Delta G_{\text{aro}} \sum_{\text{aro}} f(R)$$

$$\Delta G_{\text{lipo}} = -0.81 \text{ and } \Delta G_{\text{aro}} = -0.62 \text{ kcal/mol}$$

Desolvation effect

Active site filled with water molecules $\xrightarrow{\text{MD}}$ Unbound water molecules

$$\Delta G_{\text{solv}} = \Delta G_{\text{lipo wat}} \sum \text{unbound water}$$

$$\Delta G_{\text{lipo water}} = -0.33 \text{ kcal/mol}$$

Ligand flexibility

$$\Delta G_{\text{flex}} = \Delta G_{\text{rot}} N_{\text{rot}}$$

$$\Delta G_{\text{rot}} = +0.26 \text{ kcal/mol} \quad N_{\text{rot}} : \text{number of rotatable bonds}$$

H.J. Bohm, *J. Comput.-Aided Mol. Des.*, **1994**, *8*, 623
H.J. Bohm, *J. Comput.-Aided Mol. Des.*, **1998**, *12*, 309

Unil

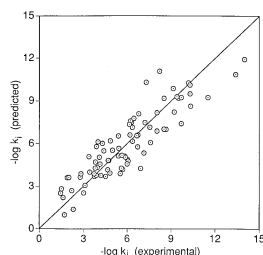


24

24

Binding free energy – Structure-based – Empirical methods

Example: LUDI score



82 complexes of the training set
SD ~2 kcal/mol

Advantages :

- Allows identification of high affinity ligands
- **Rapid** estimation of the affinities
- Structurally **diverse ligands**
- Different proteins
→ can be used routinely for **docking/virtual screening**

Others : ChemScore, VALIDATE

Shortcomings:

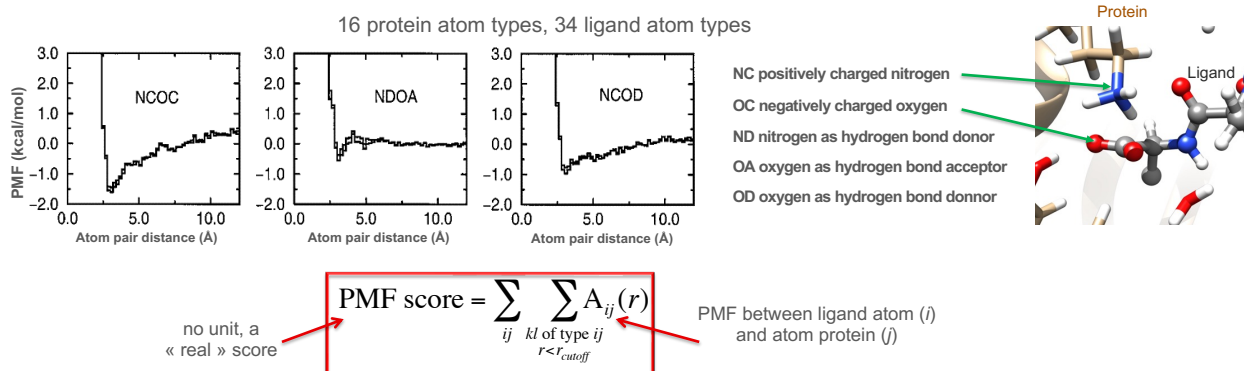
- Somewhat large errors
- Method **biased by training**:
 - certain type of proteins
 - only good complementarity protein/ligand
- Some interactions ignored:
 - cation – π
 - I...O, halogens
 - ...

25

Binding free energy – Structure-based – Knowledge-based

Example: PMF score

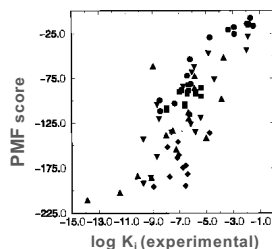
- Statistical **observations** of intermolecular **close contacts** in large 3D databases (e.g. the Protein Data Bank, **PDB**).
- The **more frequent** a protein-ligand contact between given atom types the **more favorable** to binding affinity.
- Trained on 697 PDB complexes.
- → derivation of **potential energy** ("potential of mean force", **PMF**), no need to train on ΔG_{bind}



26

Binding free energy – Structure-based – Knowledge-based

Example: PMF score



77 complexes, 5 different proteins
SD ~2 kcal/mol

Advantages :

- Allows identification of high affinity ligands
- Rapid estimation of the affinities
- Structurally different ligands
- "Universal" (different ligand and protein types)
- No fitting parameters to measured ΔG_{bind}

Drawbacks :

- Somewhat large errors
- No measure of directionality of H-bonds
- Does not estimate directly binding in kcal/mol
- Still a bias if a rare interaction is not observed.

Others : DrugScore, GoldScore (trained on CSD)

Unil



27

27

Binding free energy – Conclusion

- Large variety of methods to estimate binding free energies
- None really satisfying in terms of predictive ability versus speed
 - ➡ Consensus scoring: evaluation of ΔG_{bind} with different scoring function, select virtual ligands that are predicted to be of high affinity by several scores.
- But, can be efficient to rank similar putative ligands
- Still a "limiting problem" of molecular docking methods and more generally for computational drug design

S.S. So and M. Karplus, *J. Comput. Aided. Mol. Des.*, **2001**, 15, 613

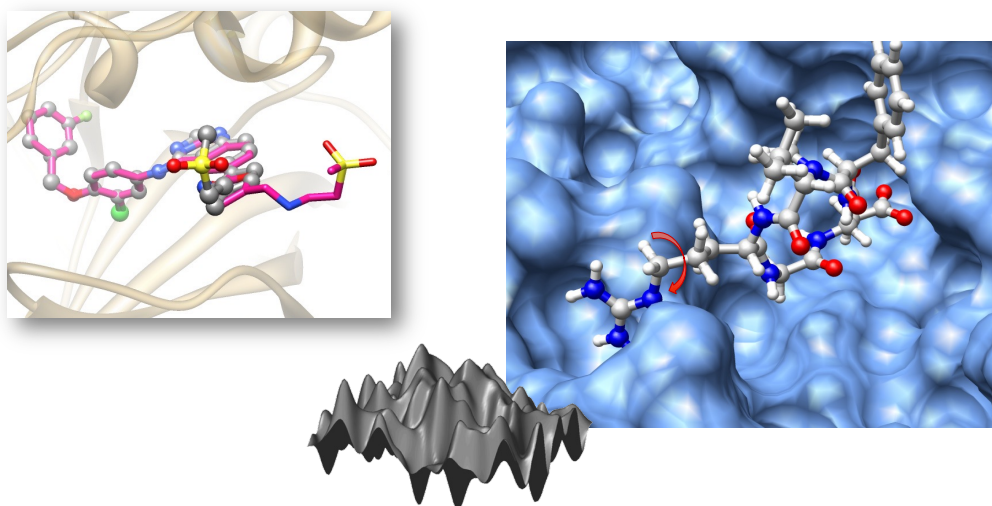
Unil



28

28

Small molecule docking



Unil



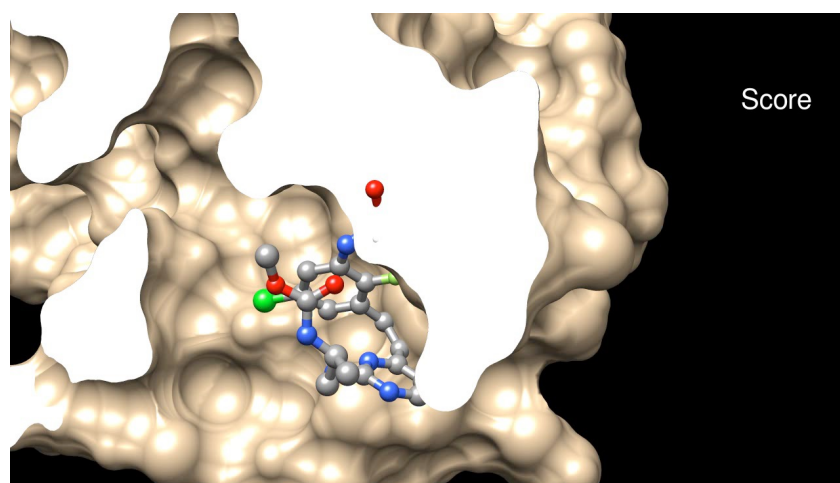
29

29

Docking – Objective

Docking small molecules into protein cavities:

Predict the **binding mode** (location, orientation and the geometry) of the **small molecule in the protein**
 = "How the small molecule is **recognized by its macromolecular target**".



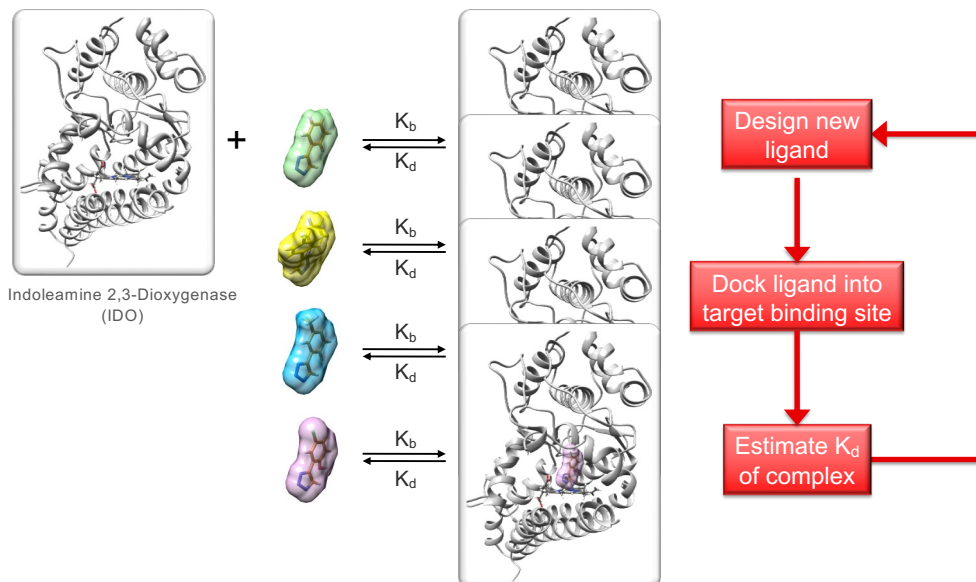
Unil



30

30

Docking – The cornerstone of structure-based ligand design



Courtesy of Dr. Ute Röhrig

Unil

SIB

31

31

Docking – Definitions

Pose: location, orientation and conformation of a small molecule on a macromolecule surface (cavity, pocket of groove) ~ tentative binding modes.

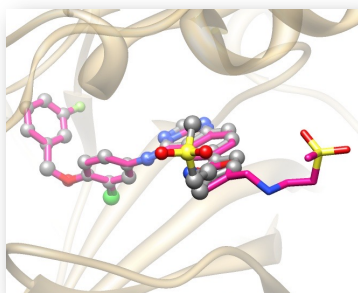
Native binding mode: experimentally defined binding mode (X-ray, NMR). Expected to be the best binding mode in term of binding free energy.

Docking: predicting the (native) binding mode using molecular modeling approaches.

Re-docking : docking on the X-ray structure of the receptor obtained in complex with the studied ligand (i.e. perfect induced fit). Used for exercise or benchmark.

Cross-docking : docking on a X-ray structure of the receptor obtained without the studied ligand (*apo* protein → no induced fit, or complex with another ligand → different induced fit)

Success: ability to predict a binding mode close to the native binding mode (when known, i.e. exercise or benchmark of the approach). Generally, RMSD < 2 Å.



Unil

SIB

32

32

Docking – The root mean square deviation (RMSD)

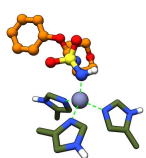
Root Mean Square Deviation (RMSD, in Å [10^{-10} m]) :

- is the average distance between the pair of atoms (normally heavy atoms).
- is the measurement of superimposition of two poses of the same molecule.
- the greater the less superimposed. RMSD = 0 means perfect overlay.

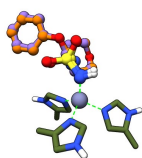
$$RMSD = \sqrt{\frac{1}{N} \sum_{i=1}^N d_{ii}^2}$$

N : the number of atoms.

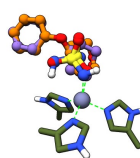
d_{ii} : the distance between atom i in the first pose and the same atom i in the second pose [Å]



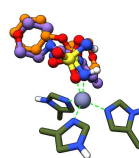
X-ray Structure



RMSD = 0.2 Å



RMSD = 1.1 Å



RMSD = 4.5 Å

Unil



33

33

Docking – The root mean square deviation (RMSD)

Root Mean Square Deviation (RMSD, in Å [10^{-10} m]) :

- is the average distance between the pair of atoms (normally heavy atoms).
- is the measurement of superimposition of two poses of the same molecule.
- the greater the less superimposed. RMSD = 0 means perfect overlay.

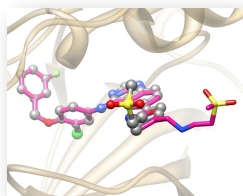
$$RMSD = \sqrt{\frac{1}{N} \sum_{ij=1}^N d_{ij}^2}$$

N : the number of pairs of atoms.

i : atom of pose 1

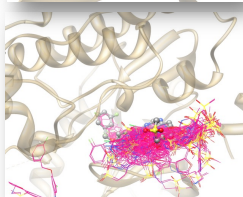
j : atom of pose 2

d_{ij} : distance between atom i and atom j [Å]



Measure of success (redocking):

- Measure all distances between heavy atom pairs of the crystallographic pose (C in grey ball & stick) and the docking pose (C in pink stick).
- Calculate RMSD.
- Rule-of-thumb: if RMSD < 2Å: SUCCESS!
- Here: ~2Å



Clustering (help at selecting the numerous poses)

- Each pose is compared to all poses by computing RMSDs
- Poses are classified into clusters (e.g. 2Å).
- All members of a given cluster have not more than 2Å RMSD to any member of the same cluster.
- A cluster represents all poses with similar binding mode.

Unil



34

34

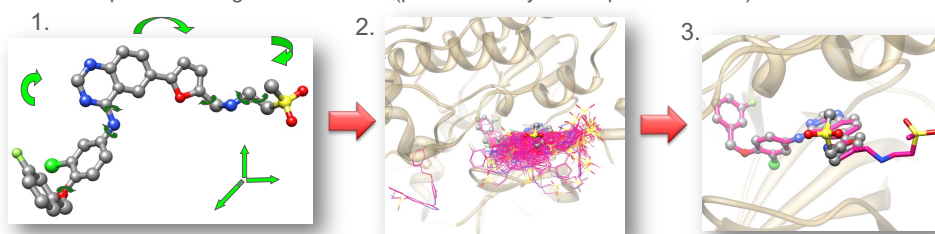
Docking – General approach

The “standard” methodology:

1. **Generate a large number of poses.** Sampling the posing space (orientational/translational/conformational) of the ligand into the protein binding site
2. **Assess the binding strength.** Scoring each possible ligand pose (~ fast evaluation of the ligand affinity)
3. **Selecting the pose with the most favorable binding (best score)**
→ **predicted binding mode**

Different levels of approximation:

- protein and ligand are rigid (the past!)
- the protein is rigid, the ligand is flexible (today, except HTVS)
- protein and ligand are flexible (possible today at computational cost)



Unil



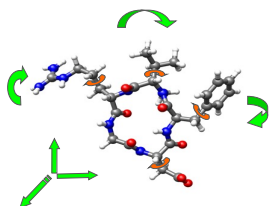
35

35

Docking – Handling ligand and protein flexibility

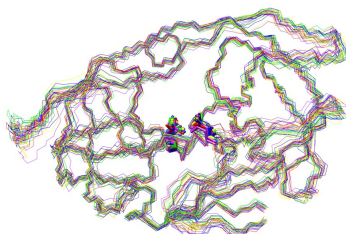
Ligand flexibility

- Pre-generated conformational libraries.
Generation of several conformers, and rigid docking of each conformer
- Conformational search “**on the fly**” by the **posing algorithm**.



Protein flexibility

- Selection of **several** different protein **conformations** (from experimental structures, or MD simulations) and parallel dockings (**ensemble docking**).
- Use of an “**averaged protein structure**”: Incorporating multiple receptor structures into a **grid** (energy-weighted grid, evt with reduced atom size).
- Conformational search “**on the fly**”. Selected side chains can take preferential known conformations (**rotamer libraries**).



Unil



36

36

Docking – Existing programs

Many programs exist. They differ in:

- the **posing** algorithm
- the handling of ligand and protein **flexibility**
- the **scoring** function

Program	Posing algorithm	Scoring function	Protein flexibility
Autodock	EA	Force field	Flexible side chains
UCSF Dock	Incremental build	Force field / contact score	Protein side chain and backbone flexibility
Autodock Vina (swissdock.ch)	MC + local search	Empirical + knowledge-based	Flexible side chains
FlexX	Incremental build	Empirical score	Ensemble of protein structures
Gold	EA	Empirical / Knowledge-based	Selected side chain / ensemble docking
Glide	Exhaustive search	Empirical score	-
EADock 2	EA	Force field	Protein side chain and backbone flexibility
EADock DSS (old.swissdock.ch)	Incremental build	Force field	Protein side chain and backbone flexibility
Attracting Cavities (swissdock.ch)	Energy minimizations	Force field	Protein side chain and backbone flexibility

Unil



37

37

Docking – Posing algorithms

Three types of sampling algorithms:



Unil



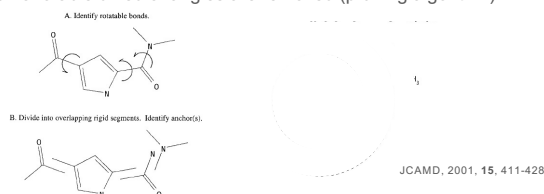
38

38

Docking – Posing algorithms – Systematic search

- Anchor and grow (Ex: EADock DSS [Swissdock.ch], FlexX, DOCK)

- The ligand is divided into rigid (core fragments) and flexible parts (side chains)
- An anchor is selected among the rigid fragments and docked into the target
- The ligand is rebuilt incrementally, starting from the anchor, through systematic dihedral angle exploration. Unfavorable dihedral angles are removed (pruning algorithm)



Methods differ in the docking of the anchor and in the pruning algorithm

- Fragment growing (Ex: Hammerhead)

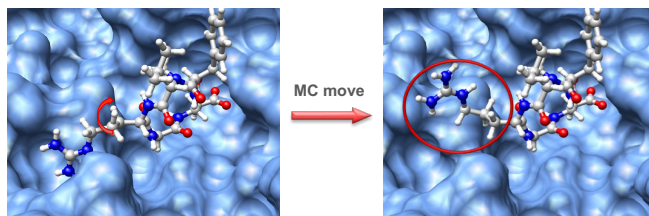
- The ligand is divided into rigid (core fragments) and flexible parts (side chains)
- All rigid fragments are docked
- The ligand is rebuilt from fragments that have acceptable initial scores

➡ More “reconstruction algorithms” than “systematic search”

Docking – Posing algorithms – Stochastic search

The monte Carlo algorithm:

- Generate an **initial pose** (ligand random conformation, translation and rotation) and **score it**
- Generate a **new pose from the previous one** (through random conformational change, translation, rotation) and **score it**
- Use **metropolis criterion**(*) to determine whether the new pose is retained
- Repeat steps 2-3 until the number of desired poses is obtained (typically >100,000)



(*) Metropolis criterion:

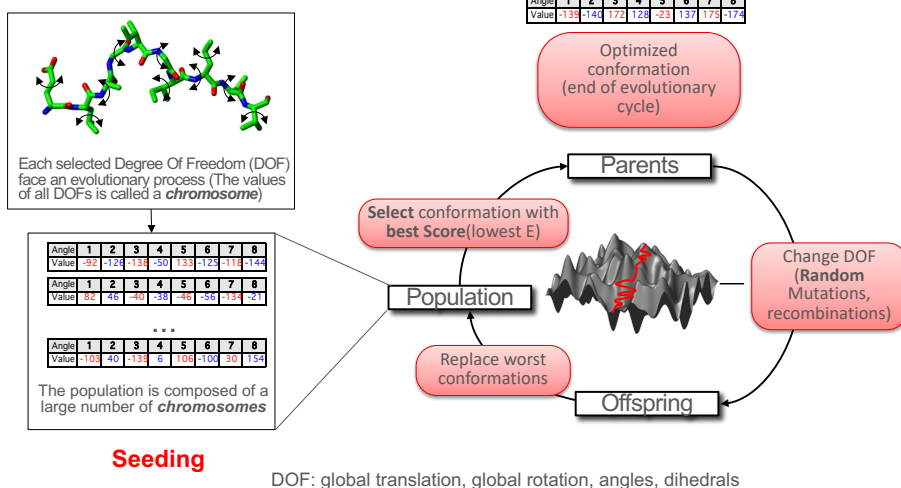
If the difference in energy between the new and previous pose (ΔE) is negative (i.e. the new pose has better interactions with the protein), then the new pose is accepted.

If ΔE is positive, a random number between 0 and 1, $0 < X < 1$, is generated and the new pose is accepted only if $\exp(-\Delta E/RT) > X$.

Docking – Posing algorithms – Stochastic search

Genetic or evolutionary algorithms:

Ex: Gold, Autodock, EADock 2



Unil



41

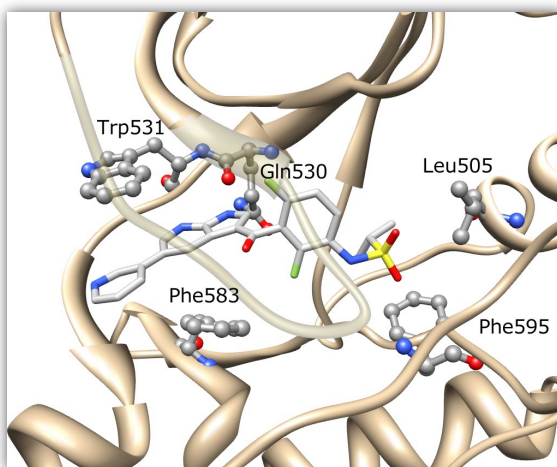
41

Docking – Posing algorithms – Deterministic search

Minimizations and Molecular Dynamics (MD) simulations cannot cross easily energy barriers.



Generally restricted to **local search around the starting pose**
May be useful to **after** systematic or stochastic process to refine poses
More useful in **post-processing**, computationally demanding!



MD simulation of a BRAF/inhibitor complex.

Unil



42

42

Docking – Scoring

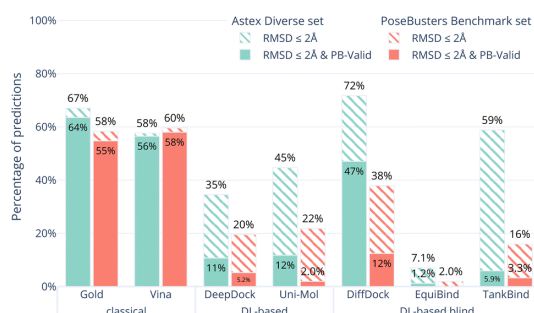
The two roles of scoring functions:

- Rank poses for one ligand in a given target (docking → predict binding mode)
 - Rank the binding modes of different ligands for a given target (compounds **selection in lead optimization** or **virtual screening** in hit-finding in large databases).
→ must be a quick estimate of the binding.
- Force-field** (physics-based) functions:
 - the most descriptive and comprehensive, potentially the most accurate,
 - slower,
 - Simplified Force field**: minimal description of interactions
→ sum of vdW and electrostatic interaction energies for each atoms (e.g. AutoDock function).
 - Exception: EADock DSS/Attracting Cavities**, affinity evaluated by computation of ΔG_{bind}
→ CHARMM force field, MMFF force field and FACTS solvation model.
 - Empirical** functions, most common for docking
 - excellent balance fastness/predictive power
 - only hard-coded interactions are accounted for
e.g. ChemScore, LUDI
 - Knowledge-based** scoring functions
 - good speed/accuracy balance
 - a “real” score: without unit nor true physical meaning.
e.g. DrugScore, GoldScore, PMF-Score

43

Docking – Beyond Physics-Based Docking: Deep Learning

Method	Authors	Date	Search space
DeepDock ⁸	Méndez-Lucio <i>et al.</i>	Dec 2021	Pocket
DiffDock ⁷	Corso <i>et al.</i>	Feb 2023	Blind
EquiBind ⁸	Stärk <i>et al.</i>	Feb 2022	Blind
TankBind ⁹	Lu <i>et al.</i>	Oct 2022	Blind
Uni-Mol ¹⁰	Zhou <i>et al.</i>	Feb 2023	Pocket



- PoseBusters: plausibility checks for generated molecule poses (i.e. are predicted binding modes follow the physics behind molecular interactions?)
- Co-folding: predict protein and ligand structure simultaneously (Umol, AlphaFold 3)

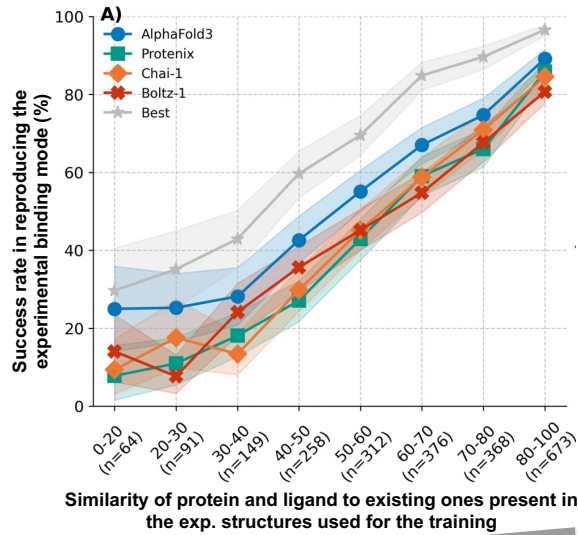
M. Buttenschon, G. M. Morris, C. Deane, **PoseBusters: AI-based docking methods fail to generate physically valid poses or generalise to novel sequences**, *Chem. Sci.* **15**, 3130 (2024)

Abramson, J., Adler, J., Dunger, J. et al. **Accurate structure prediction of biomolecular interactions with AlphaFold 3**. *Nature* (2024)

Bryant, P., Kelkar, A., Guljas, A. et al. **Structure prediction of protein-ligand complexes from sequence information with Umol**. *Nat Commun* **15**, 4536 (2024)

44

Docking – Beyond Physics-Based Docking: Deep Learning



DL cofolding approaches are overfitted: they predict/reproduce well known complexes, but fail to generalize to new, unseen, complexes

Have protein-ligand co-folding methods moved beyond memorisation?
Peter Škrinjar, Jérôme Eberhardt, Janani Durairaj, Torsten Schwede, BioRxiv, Feb 2025

45

Docking – Beyond Physics-Based Docking: Deep Learning

Current Shortcomings of Deep Learning approaches

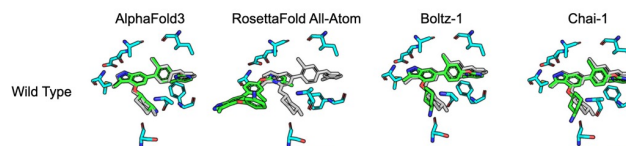


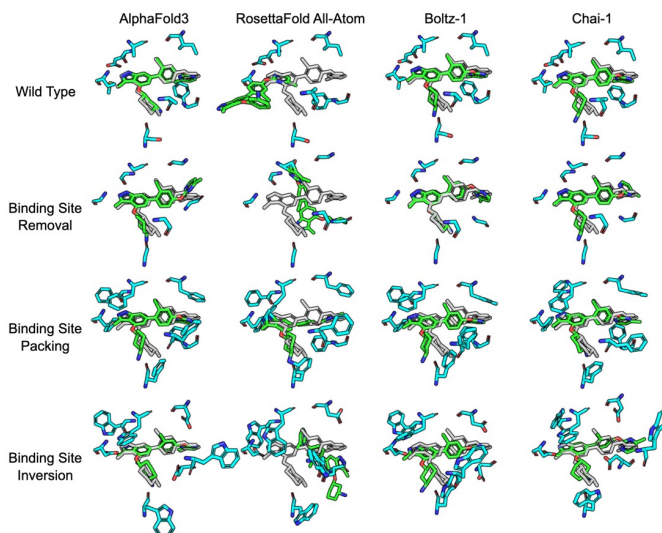
Fig. 2 | Binding site mutagenesis challenges against co-folding models using the MEK1 system (PDB: 7XLP). Predicted binding site residues are shown as cyan sticks, predicted ligand poses are shown as green sticks, and the original co-crystallized ligand pose is shown as gray sticks. Results are presented in the same manner as Fig. 1.

M.R. Masters, A.H. Mahmoud, M.A. Lill, Investigating whether deep learning models for co-folding learn the physics of protein-ligand interactions, Nat Comm. 2025

46

Docking – Beyond Physics-Based Docking: Deep Learning

Current Shortcomings of Deep Learning approaches



DL cofolding approaches do not consider correctly molecular interactions

Fig. 2 | Binding site mutagenesis challenges against co-folding models using the MEK1 system (PDB: 7XLP). Predicted binding site residues are shown as cyan sticks, predicted ligand poses are shown as green sticks, and the original co-crystallized ligand pose is shown as gray sticks. Results are presented in the same manner as Fig. 1.

M.R. Masters, A.H. Mahmoud, M.A. Lill, Investigating whether deep learning models for co-folding learn the physics of protein-ligand interactions, Nat Comm. 2025



47

47

Docking – Success

Success: ability to predict a binding mode close to the native binding mode (redocking, i.e. exercise or test of the approach). Generally, RMSD < 2 Å.

Success rate : ~ 50 to 90 % in benchmarks (re-docking)

in “real” application (cross docking) the success rate decreases by at least 30%

→ room for improvement.

→ need for high precision docking programs, handling protein flexibility (but also water, ion, ...).

Unil



48

48

Binding free energy estimation and molecular docking

Contacts: vincent.zoete@unil.ch , antoine.daina@sib.swiss