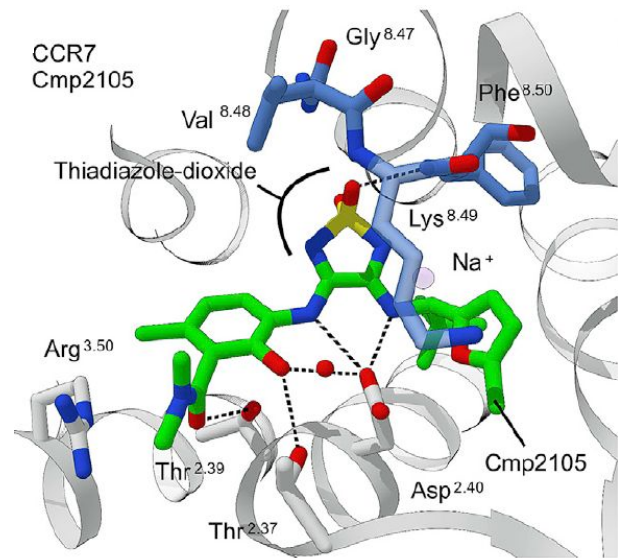
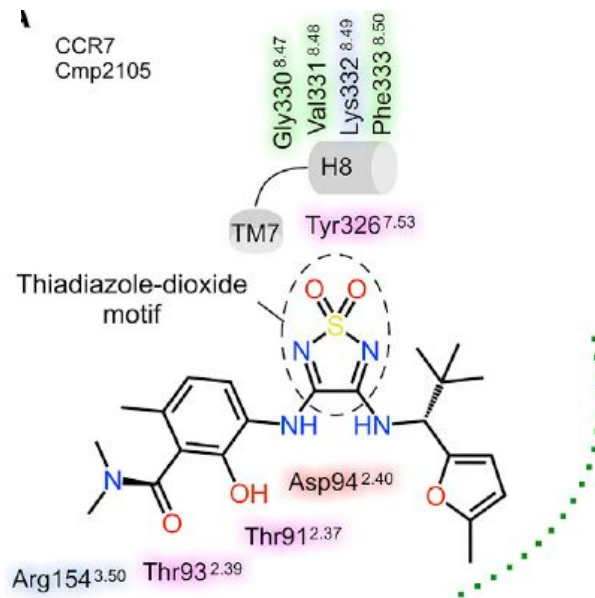
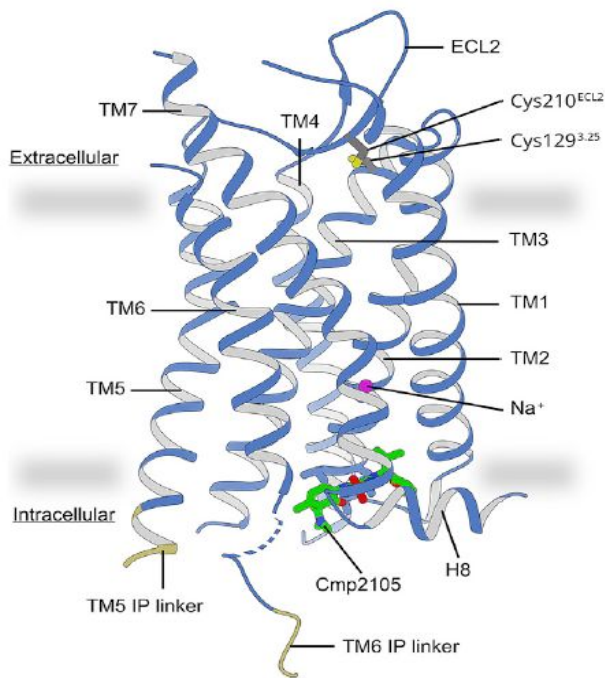


# AMIDD Lecture 6: Structure- and ligand-based drug design



**Dr. Jitao David Zhang, Computational Biologist**

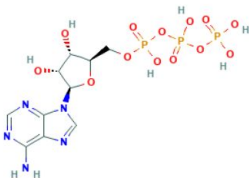
<sup>1</sup> **Pharmaceutical Sciences, Pharma Research and Early Development, Roche Innovation Center Basel, F. Hoffmann-La Roche**

<sup>2</sup> **Department of Mathematics and Informatics, University of Basel**

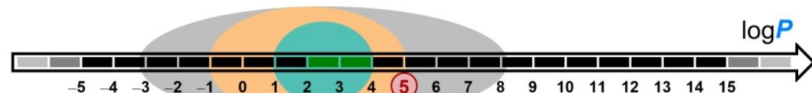
Jaeger, Kathrin, Steffen Bruenle, Tobias Weinert, Wolfgang Guba, Jonas Muehle, Takuya Miyazaki, Martin Weber, et al. "[Structural Basis for Allosteric Ligand Recognition in the Human CC Chemokine Receptor 7.](#)" Cell 178, no. 5 (August 22, 2019): 1222-1230.e10..

# Lipinski's Rule of Five of small-molecule drugs

- **HBD $\leq$ 5**: No more than **5 hydrogen-bond donors**, e.g. the total number of nitrogen–hydrogen and oxygen–hydrogen bonds.
- **HBA $\leq$ 10**: No more than **10 hydrogen-bond acceptors**, e.g. all nitrogen or oxygen atoms
- **MW $<$ 500**: A **molecular weight** less than **500 Daltons**, or 500 g/mol. Reference: ATP has a molecular mass of ~507.
- **logP $\leq$ 5**: An **octanol-water partition coefficient (log P)** that does not exceed **5**. (10-based)



ATP



- approved marketed drugs
- optimal oral drugs
- optimal CNS drugs
- ⑤ Lipinski's Rule of Five

Source: [cheminfographic.com](http://cheminfographic.com)

Table 1. New FDA Approvals (2014 to Present)<sup>a</sup> of Oral bRo5 Drugs

drug	year approved	therapeutic area	MW	cLogP	HBD	N+O
velpatasvir	2016	HCV	883.02	2.5	4	16
venetoclax	2016	oncology	868.44	10.4	3	14
elbasvir	2016	HCV	882.0	2.6	4	16
grazoprevir	2016	HCV	766.90	-2.0	3	15
cobimetinib	2015	oncology	531.31	5.2	3	5
daclatasvir	2015	HCV	738.88	1.3	4	14
edoxaban	2015	cardiovascular	548.06	-0.9	3	11
ombitasvir	2014	HCV	894.13	1.3	4	15
paritaprevir	2014	HCV	765.89	1.1	3	14
netupitant	2014	nausea from chemotherapy	578.59	6.8	0	5
ledipasvir	2014	HCV	889.00	0.9	4	14
ceritinib	2014	oncology	558.14	6.5	3	8

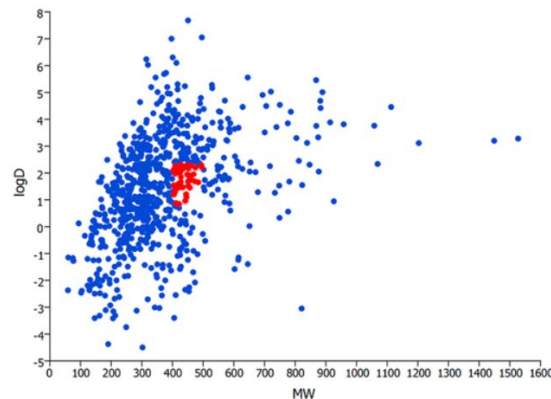
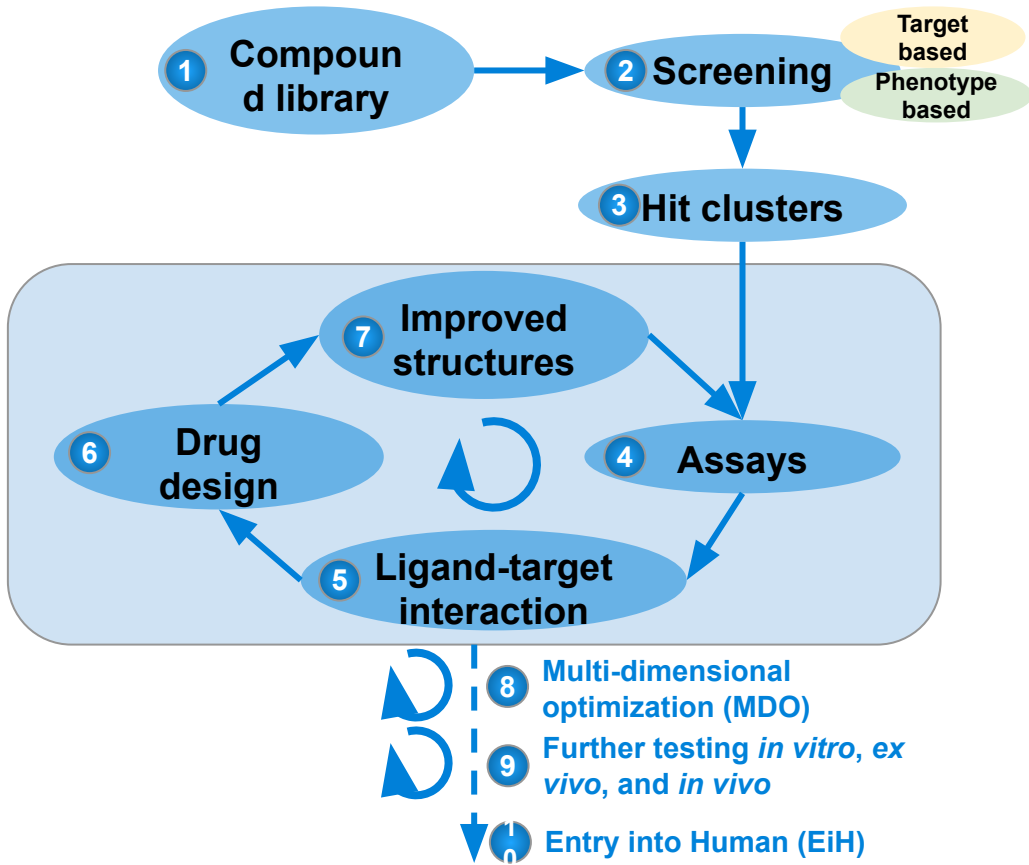


Figure 7: Plot of MW vs cLogD of FDA approved oral drugs. Red points: 'high probability area' supposed by (questionable) data analysis. Shultz, Michael D. 2019. "[Two Decades under the Influence of the Rule of Five and the Changing Properties of Approved Oral Drugs.](#)" Journal of Medicinal Chemistry 62 (4): 1701–14.

DeGoey, *et al.*. 2018. "[Beyond the Rule of 5: Lessons Learned from AbbVie's Drugs and Compound Collection.](#)" Journal of Medicinal Chemistry 61 (7): 2636–51.

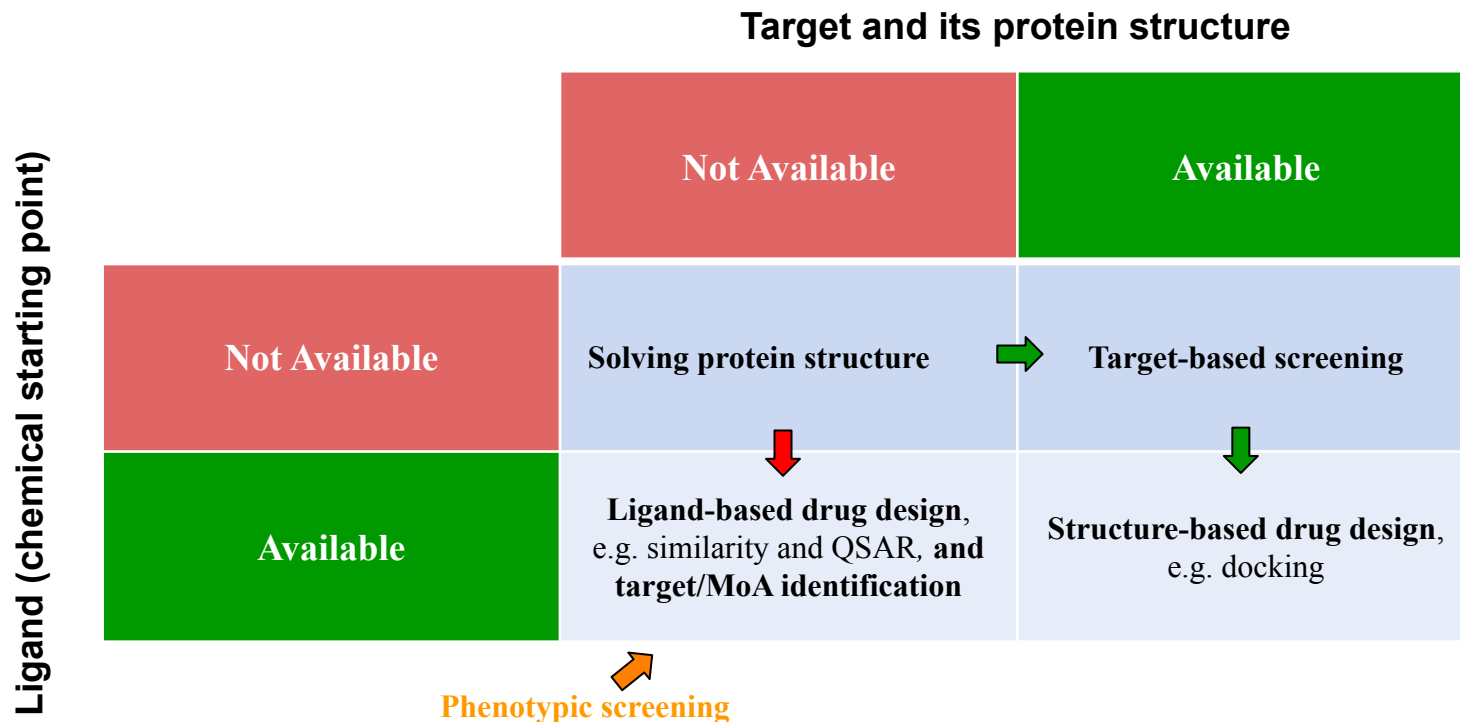
# Workflow in a typical drug-discovery program

1. Compound library construction;
2. Screening compounds with **bioassays**, or **assays**, which determine potency of a chemical by its effect on biological entities: proteins, cells, *etc*;
3. Hit identification and clustering;
4. More assays, complementary to the assays used in the screening, maybe of lower throughput but more biologically relevant;
5. Analysis of ligand-target interactions, for instance by getting the co-structure of both protein (primary target, and off-targets if necessary) and the hit;
6. *Drug design*, namely to modify the structure of the drug candidate;
7. Analog synthesis and testing (back to step 4);
8. Multidimensional Optimization (MDO), with the goal to optimize potency, selectivity, safety, bioavailability, *etc*;
9. Further *in vitro*, *ex vivo*, and *in vivo* testing, and preclinical development;
10. Entry into human (Phase 0 or phase 1 clinical trial).



A schematic presentation of structure-based drug discovery

# Structure-based and ligand-based drug design



QSAR= quantitative structure activity relationship; MoA= mechanism of action, or mode of action

# Questions about *Evaluation of the Biological Activity of Compounds: Techniques and Mechanism of Action Studies*

Q1. An important chemical and mathematical concept was not described in the book chapter: what does *the Law of Mass Action* mean? (An ODE model of reaction rate and reactant mass)

Q2: Which quantity measures binding affinity directly: dissociation constant ( $K_D$ ) or the concentration of the test compound that produces 50 percent inhibition ( $IC_{50}$ )? ( $K_D$ )

Q3: In Figure 2.3, what do x- and y-axis represent in panel (A) and panel (B), respectively? (concentrations in x-axis; y-axis: counts per minute of radioactivity (A), percentage of binding of the labelled compound)

Q4: What is a sigmoidal curve? (A S-shaped, logistic or logit curve)

Q5: Do  $IC_{50}$  values indicate a particular mechanism of action (MoA)? (No)

Q6: In a certain enzymatic assay, two compounds have the following  $pIC_{50}$  values: 7.2 (Compound A), 9.3 (Compound B). If all other conditions are held constant, what is the relationship between binding affinities of the two compounds with regard to the target? (B>A)

Q7: Why is DMSO often used in bioassays? (solvent, control)

Q8: Can you use your own language to describe what is the Hill function? (discussed in Lecture 5)

Q9: What statistical measure is used to measure the signal-noise ratio in screening? Can you use your own language explaining it? (how well can we separate positive controls from negative controls)

Q10: Why logarithm (usually base 10) transformation is often preferred to represent quantities such as  $IC_{50}$  and  $K_i$ ? (presentation, as well as statistical mechanistics)

## Questions from you:

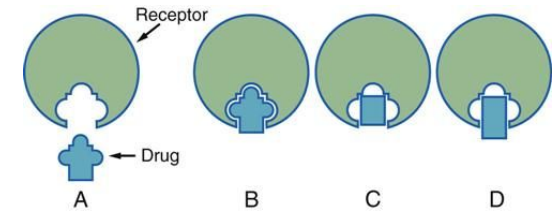
1. On page 19: what is meant with "displacement of a labelled ligand"? (I do not know what 'displacement' means in that context)
2. I didn't quite understand the application of the Z value and when it usually is used

# Outline

- **Affinity**
  - The (bio)physical view
  - The (bio)chemical view
- The **Michaelis-Menten model** and enzymatic kinetics
- Example of structure-based drug design: **molecular docking**
- Example of ligand-based drug design: **similarity and quantitative structure-activity relationship (QSAR)**

# The biophysical and biochemical views of ligand-target binding

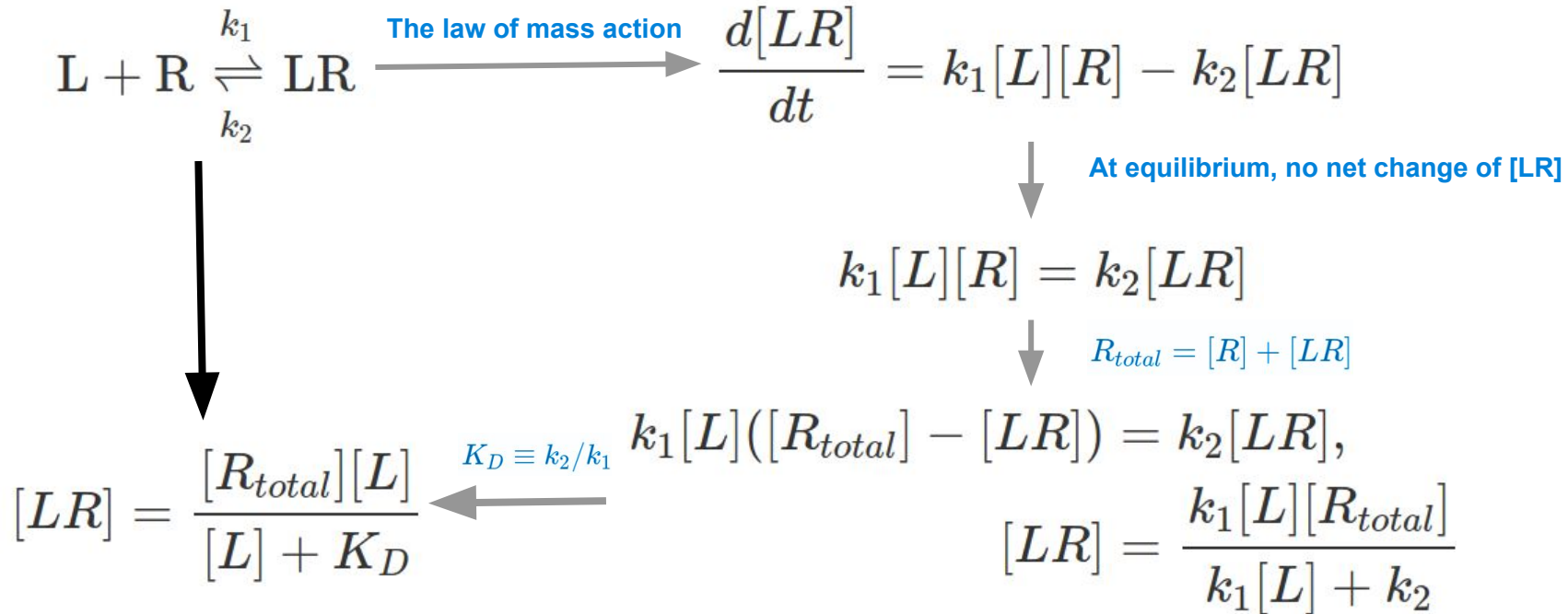
- A **ligand** is a substance that forms a complex with a biomolecule to serve a biological purpose. For instance, a drug can produce a signal by binding to a site on a target protein.
  - A ligand that binds to and alters the function of the receptor that triggers a physiological response is called a receptor **agonist**.
  - A ligand that binds to a receptor but fail to activate the physiological response is a receptor **antagonist**.
- **The biophysical view of binding:** Binding occurs in favourable steric, *i.e.* spatial, configurations (**the ‘lock-and-the-key’ model**) and is mediated by intermolecular forces, such as electrostatic interactions (ionic bonds, hydrogen bonds), Van der Waals forces (dipole interactions),  $\pi$ -effects (interactions of  $\pi$ -orbitals of a molecular system), and hydrophobic effect. Both enthalpy and entropy contribute to the binding energy.
- **The biochemical view of binding:** The *rate* of binding is called affinity, often expressed in  $K_d$  or, for inhibitors,  $K_i$ . A closely related, and often confusing, concept is  $IC_{50}$ . We will talk about them in the next lecture when we talk about the Michaelis-Menten model, the dose-response curve, and the Hill function.
- **Binding affinity data alone does not determine the overall potency of a drug.** Potency depends on binding affinity, the ligand efficacy, and many other factors.



Competitive	Uncompetitive
$  \begin{array}{c}  E + S \xrightleftharpoons{K_m} ES \xrightarrow{k_{cat}} E + P \\  \downarrow \text{EI} \\  E + I \xrightleftharpoons{K_i} EI  \end{array}  $	$  \begin{array}{c}  E + S \xrightleftharpoons{K_m} ES \xrightarrow{k_{cat}} E + P \\  \downarrow \text{ESI} \\  E + S + I \xrightleftharpoons{K_i} ESI  \end{array}  $
Non-competitive	Mixed
$  \begin{array}{c}  E + S \xrightleftharpoons{K_m} ES \xrightarrow{k_{cat}} E + P \\  \downarrow \text{EI} \\  E + I \xrightleftharpoons{K_i} EI \\  EI + S \xrightleftharpoons{K_m} ESI  \end{array}  $	$  \begin{array}{c}  E + S \xrightleftharpoons{K_m} ES \xrightarrow{k_{cat}} E + P \\  \downarrow \text{EI} \\  E + S + I \xrightleftharpoons{\alpha K_i} ESI \\  EI + S \xrightleftharpoons{\alpha K_m} ESI  \end{array}  $

Four basic types of kinetic mechanism of inhibition, source: [sciencesnail.com](https://www.sciencesnail.com)

# From the law of mass action to ligand-target interaction





# Four classical classes of mathematical models

## Compartment models

$$\frac{d[LR]}{dt} = k_1[L][R] - k_2[LR]$$

Kinetics of ligand-target interaction

$$\frac{dx}{dt} = \alpha x - \beta xy,$$

$$\frac{dy}{dt} = -\gamma y + \delta xy,$$

The Lotka-Volterra equations modelling predator-prey relationships.

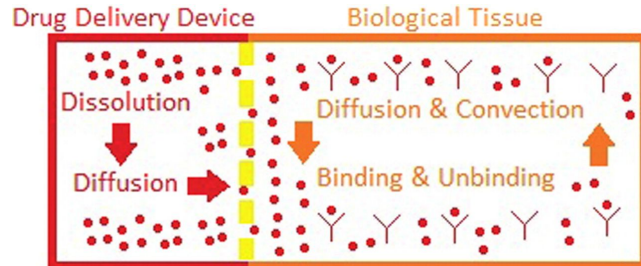
$$\frac{dS}{dt} = -\frac{\beta IS}{N},$$

$$\frac{dI}{dt} = \frac{\beta IS}{N} - \gamma I,$$

$$\frac{dR}{dt} = \gamma I$$

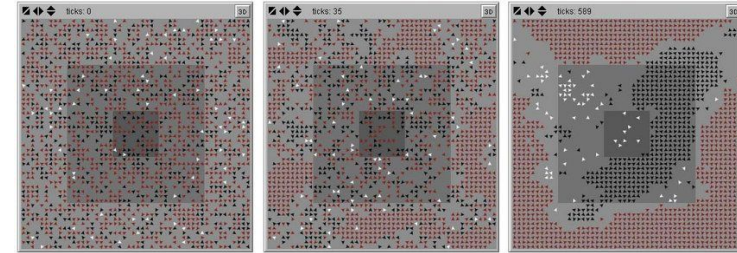
The SIR (S=susceptible, I=infectious, R=removed) model of epidemiology

## Transport models



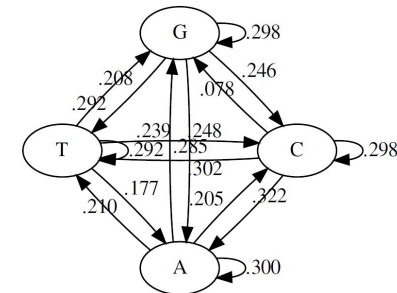
McGinty, Sean, and Giuseppe Pontrelli. 2015. "[A General Model of Coupled Drug Release and Tissue Absorption for Drug Delivery Devices](#)." *Journal of Controlled Release* 217 (November): 327–36.

## Particle models



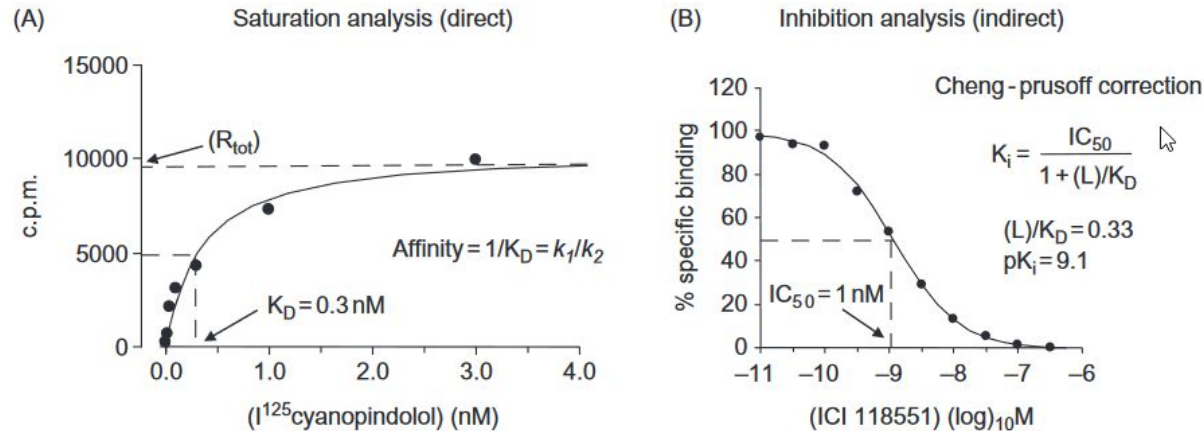
A Study on Socio-spatial Segregation Models Based on Multi-agent Systems by Quadros *et al.* (2012). 10.1109/BWSS.2012.14.

## Finite state models



A finite-state Markov chain modelling DNA sequences

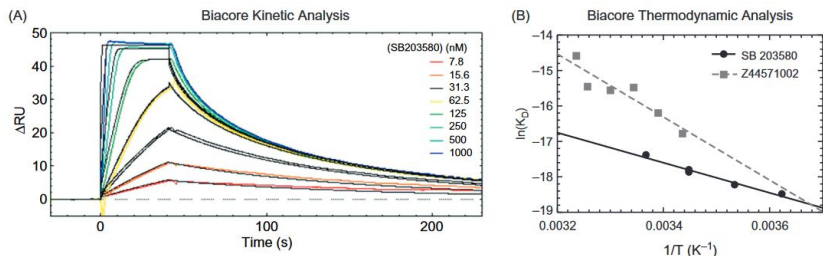
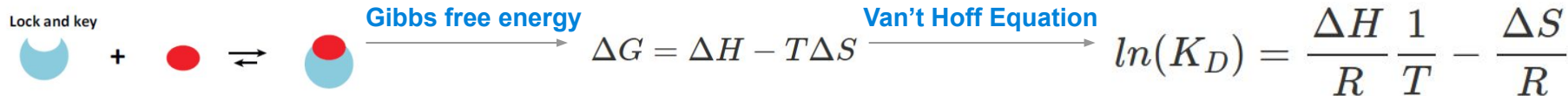
# The biochemical (kinetic) view of binding affinity: the hyperbola curve and the dissociation constant $K_D$



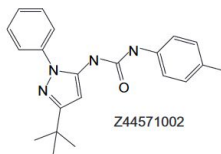
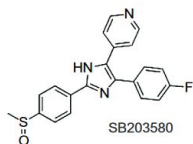
**Binding assays with direct and indirect measurements.** (A) A direct binding assay using I<sup>125</sup> labelled cyanopindolol as a  $\beta_2$ -adrenoceptor ligand. The curve describes a rectangular hyperbola which saturates at high ligand concentration. The ligand dissociation constant ( $K_D$ ) was estimated as 0.3 nM and is a measure of the ligand affinity. (B) A typical inhibition analysis using membranes expressing the human  $\beta_2$ -adrenoceptor and employing 0.1 nM I<sup>125</sup> cyanopindolol as the labeled ligand. The displacing ligand, the selective  $\beta_2$ -adrenoceptor antagonist ICI 118551, produces complete inhibition of the specific binding yielding an  $IC_{50}$  of 1 nM. From *Evaluation of the Biological Activity of Compounds: Techniques and Mechanism of Action Studies*, by Iain G. Dougall and John Unitt.

**Questions:** (1) how can we interpret the hyperbola curve? (2) if  $f(x)$  is a function with the form of  $Ax/(k+x)$ , what will be the form of function  $g(f(x))$  where  $g(x)=Bx/(k'+x)$ ? What implications does this have?

# The biophysical (thermodynamic) view of binding affinity: enthalpy and entropy



Compound Name	k1	k2	KD	ΔG	ΔH	TΔS
Z44571002	2.2e4 ± 3e2	0.001 ± 8.0e6	5.2e-8	-40	-75	-35
SB203580	1.7e6 ± 1.7e5	0.130 ± 0.014	7.8e-8	-43	-36	7.5

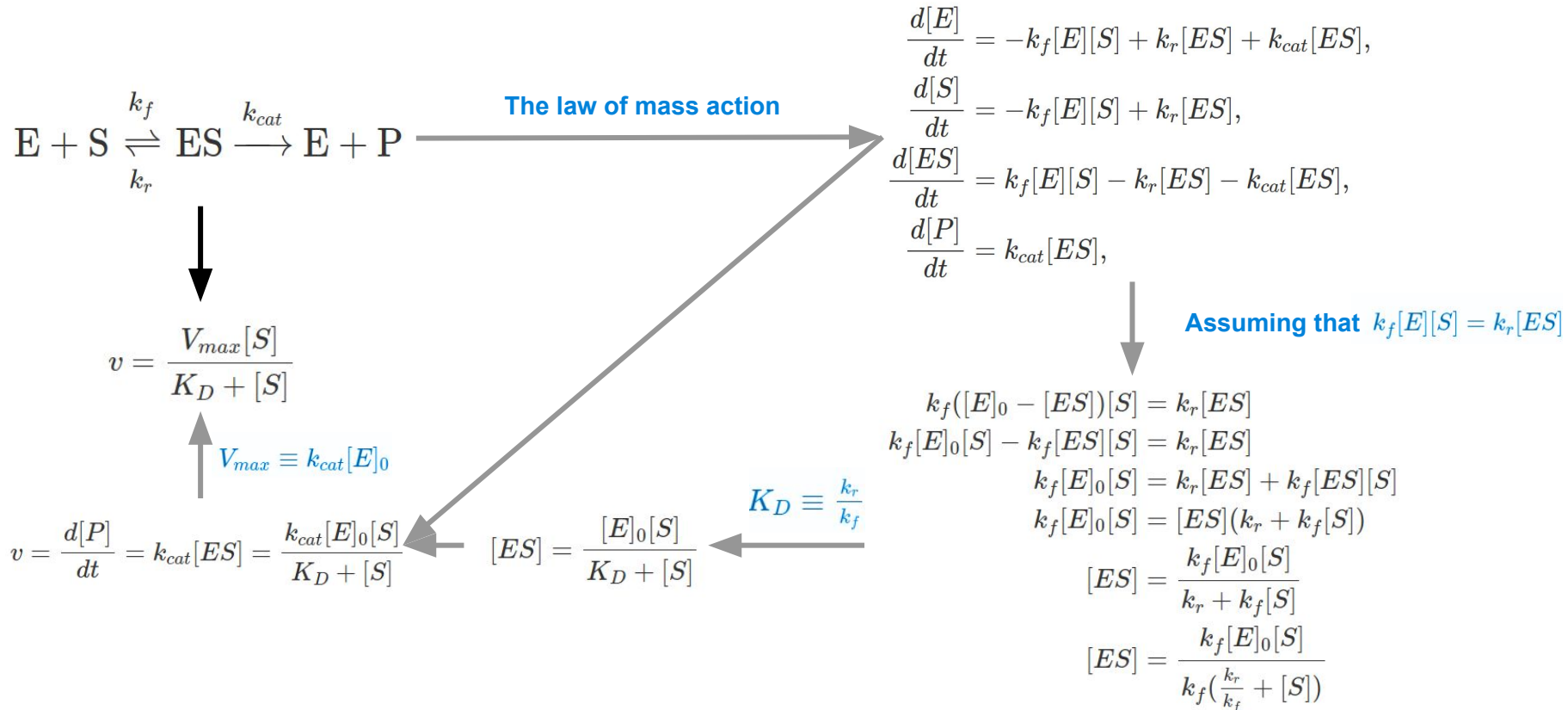


## Kinetic and thermodynamic measurements of two p38α inhibitors.

(A) The time course of SB203580 binding to immobilized mitogen activated kinase p38α. The y-axis shows the mass change resulting from compound binding to p38α. At t=0, a range of SB203580 concentrations were passed across the immobilized p38α to measure net association, and then at t=50s, the compound is replaced with buffer to initiate dissociation. The table shows the association and dissociation rate constants as well as the equilibrium dissociation constants (KD(M)) for two compounds. (B) Thermodynamic analysis. Enthalpy and entropy components of binding derived from the Van't Hoff analysis are detailed in the attached table. ΔG, ΔH and TΔS values are in kJ/mol.

For a thorough discussion about enthalpic and entropic contributions to molecular interactions, see [A Medicinal Chemist's Guide to Molecular Interactions](#) (Journal of Medicinal Chemistry 53 (14): 5061–84) by Bissantz et al.

# Modelling enzyme kinetics with the Michaelis-Menten model



# The dose-response curve and IC50: The Hill function and *in vitro* pharmacology

- The Hill function is one of the mostly useful non-linear functions to model biological systems.
- In its general form,  $H_{max}$  indicates the maximal value to which the function is asymptotic,  $n$  is the shape parameter (known as the Hill's coefficient), and  $k$  is the reflection point, often abbreviated as  $XC_{50}$  (X=I, E, C, ...), the half-saturation constant.
- The Michaelis-Menten model is a special case of the Hill function with  $n=1$ .

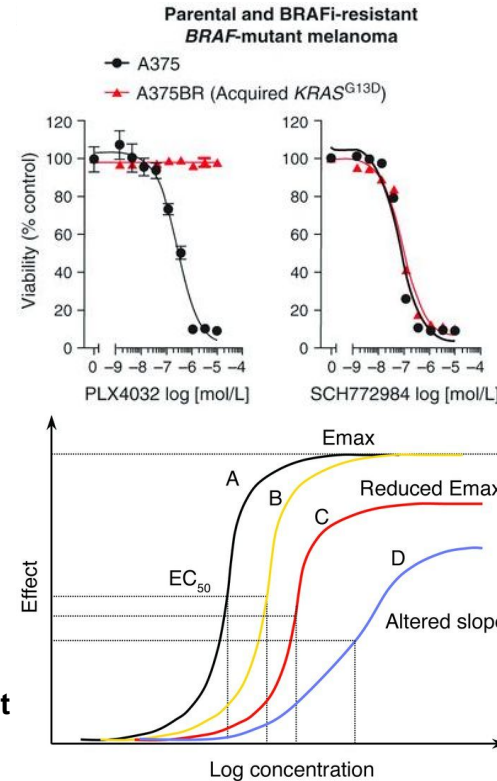
$$H = H_{max} \frac{x^n}{k^n + x^n}$$

**General form of the Hill function**

$$E = E_{max} \frac{[L]^n}{EC_{50}^n + [L]^n}$$

$$= E_{max} \frac{1}{1 + \left(\frac{EC_{50}}{[L]}\right)^n}$$

**Modelling dose-dependent effect**



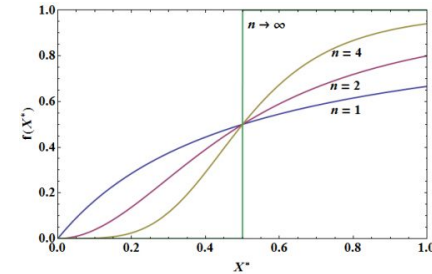
[Morris et al. Cancer Discov. 3\(7\): 742–50. ©2013 AACR.](#)

White. *J Clin Invest.* 2004;113(8):1084-1092. <https://doi.org/10.1172/JC121682>.

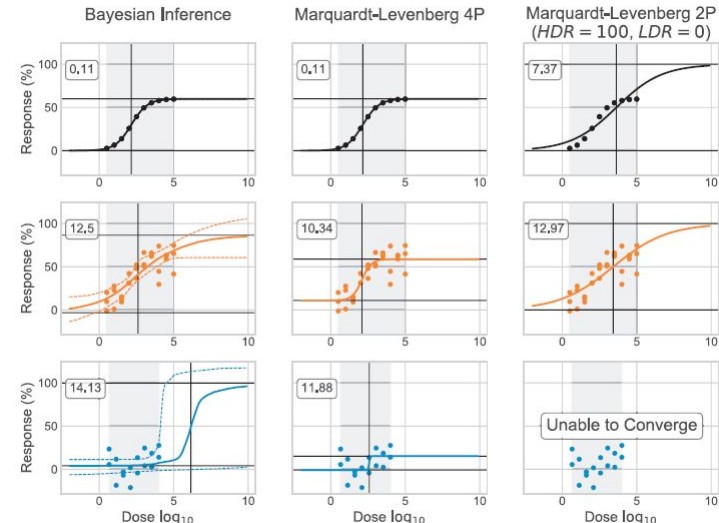
Suppose it is an antiviral drug, compared with curve B, what does curve A, C, and D suggest?

# More about the Hill function and dose-response curves

- The Hill function is often used to model either *target occupancy* or *tissue response*. In pharmacology, it is often used to model the tissue response.
- The Hill function can be approximated by a step function when  $n$  goes towards infinity (top panel). This can be seen as one of the theoretical foundations of Boolean network modelling.
- The Hill function can be deduced from statistical mechanics of binding, a particle modelling approach. See for instance [an article on Biophysics Wiki by Andreas Piehler](#) for details.
- Data needs to be fit to the model, and in reality data can look quite different from the ideal curve (bottom panel). By setting priors, it is possible to perform inference even with ill-looking data.



From [the biophysics wiki article](#) by Andreas Piehler



The Bayesian inference approach versus the Marquardt-Levenberg algorithm for non-linear regression fitting (an alternative to gradient descent and Gauss-Newton methods). 4P: four parameter model; 2P: two parameter model (IC50 and  $n$ ). Numbers in boxes are root mean square errors of fitting. Figure 2 from Labelle, Caroline, Anne Marinier, and Sébastien Lemieux. 2019. [“Enhancing the Drug Discovery Process: Bayesian Inference for the Analysis and Comparison of Dose–Response Experiments.”](#) *Bioinformatics* 35 (14): i464–73.



# The principle of molecular docking, a case study of structure-based drug design

- **Docking is like a discotheque: it is all about posing and scoring – Roger Sayle (NextMove Software Limited)**
- Three basic methods to represent target and ligand structures *in silico*
  - **Atomic:** used in conjunction with a potential energy function, computational complexity high
  - **Surface:** often used in protein-protein docking
  - **Grid representation:**
    - Basic idea: to store information about the receptor's energetic contributions on grid points so that it only needs to be read during ligand scoring.
    - In the most basic form, grid points store two types of potentials: **electrostatic** and **van der Waals forces**.

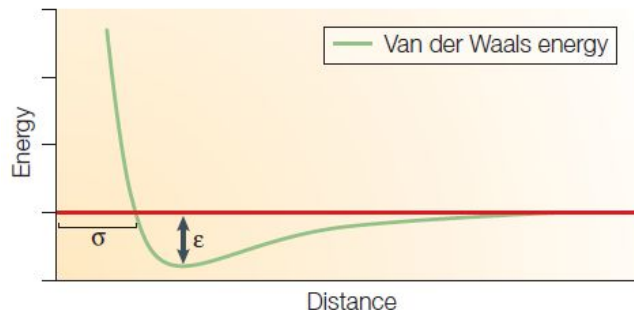
$$E_{coul}(r) = \sum_{i=1}^{N_A} \sum_{j=1}^{N_B} \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}}$$

**Coulombic interactions**

$$E_{vdW}(r) = \sum_{j=1}^N \sum_{i=1}^N 4\epsilon \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right]$$

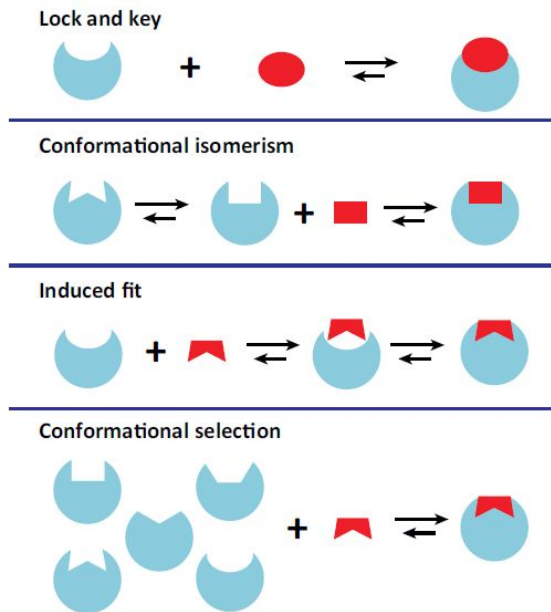
**Lennard–Jones 12–6 function**

- $\epsilon$  is the **well depth** of the potential
- $\sigma$  is the **collision diameter** of the respective atoms  $i$  and  $j$ .



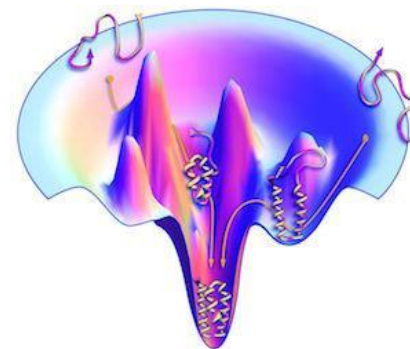
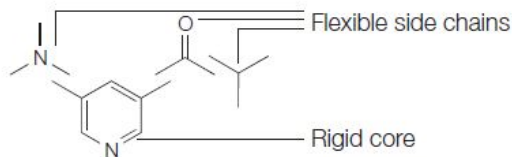
Kitchen, Douglas B., Hélène Decornez, John R. Furr, und Jürgen Bajorath. „Docking and Scoring in Virtual Screening for Drug Discovery: Methods and Applications“. *Nature Reviews Drug Discovery* 3, Nr. 11 (November 2004): 935–49. <https://doi.org/10.1038/nrd1549>.

# Posing: dealing with flexibility of ligand and of protein



*TRENDS in Pharmacological Sciences*

Chen, Yu-Chian. „Beware of docking!“ *Trends in Pharmacological Sciences* 36, Nr. 2 (1. Februar 2015): 78–95.  
<https://doi.org/10.1016/j.tips.2014.12.001>.

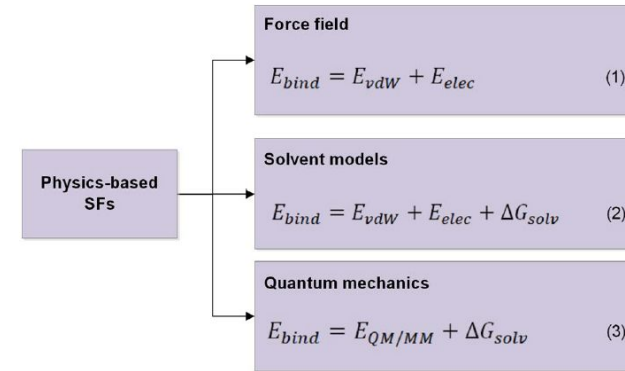
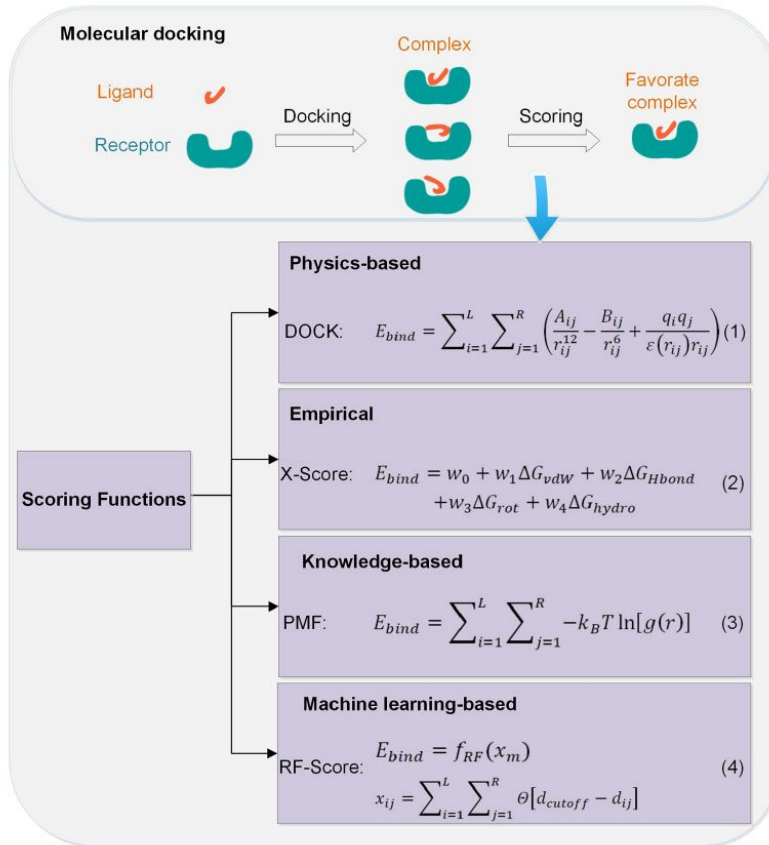


## Methods to deal with ligand and protein flexibility

- Systematic search
- Random search, such as Monte-Carlo and genetic algorithms
- Simulation methods, such as molecular dynamics



# Types of scoring functions



- Empirical scoring functions estimate the binding affinity of a complex by **summing up the important energetic factors for protein–ligand binding**, such as hydrogen bonds, hydrophobic effects, steric clashes, etc. It relies on training set and regression analysis.
- Knowledge-based scoring functions derive the desired pairwise potentials from three-dimensional structures of a large set of protein–ligand complexes based **on the inverse Boltzmann distribution**. It is assumed that the frequency of different atom pairs in different distances is related to the interaction of two atoms and converts the frequency into the distance-dependent potential of mean force.
- Machine learning-based scoring functions are usually used for rescoring to improve the initial docking.

Li, Jin, Ailing Fu, and Le Zhang. „An Overview of Scoring Functions Used for Protein–Ligand Interactions in Molecular Docking“. *Interdisciplinary Sciences: Computational Life Sciences* 11, Nr. 2 (1. Juni 2019): 320–28. <https://doi.org/10.1007/s12539-019-00327-w>.

# Interested in learning more about molecular modelling?

PROTOCOL

## Computational protein–ligand docking and virtual drug screening with the AutoDock suite

Stefano Forli, Ruth Huey, Michael E Pique, Michel F Sanner, David S Goodsell & Arthur J Olson

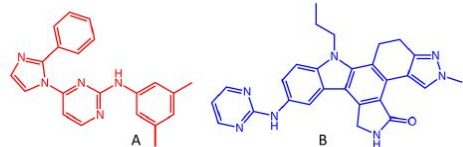
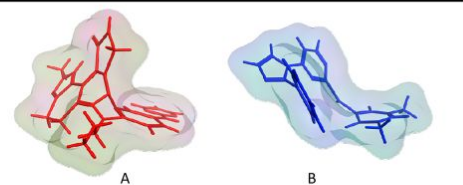
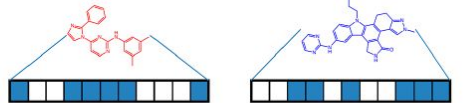
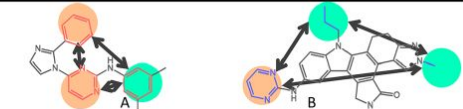
Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, California, USA. Correspondence should be addressed to A.J.O. (olson@scripps.edu).

Published online 14 April 2016; doi:10.1038/nprot.2016.051

Computational docking can be used to predict bound conformations and free energies of binding for small-molecule ligands to macromolecular targets. Docking is widely used for the study of biomolecular interactions and mechanisms, and it is applied to structure-based drug design. The methods are fast enough to allow virtual screening of ligand libraries containing tens of thousands of compounds. This protocol covers the docking and virtual screening methods provided by the AutoDock suite of programs, including a basic docking of a drug molecule with an anticancer target, a virtual screen of this target with a small ligand library, docking with selective receptor flexibility, active site prediction and docking with explicit hydration. The entire protocol will require ~5 h.

- Try docking yourself by following this protocol: Forli, Stefano, Ruth Huey, Michael E. Pique, Michel F. Sanner, David S. Goodsell, und Arthur J. Olson. „Computational Protein–Ligand Docking and Virtual Drug Screening with the AutoDock Suite“. *Nature Protocols* 11, Nr. 5 (Mai 2016): 905–19. <https://doi.org/10.1038/nprot.2016.051>.
- In-depth reading: Sliwoski, Gregory, Sandeepkumar Kothiwale, Jens Meiler, und Edward W. Lowe. „Computational Methods in Drug Discovery“. *Pharmacological Reviews* 66, Nr. 1 (1. Januar 2014): 334–95. <https://doi.org/10.1124/pr.112.007336>.
- A more advanced talk by Arthur Olson can be found [here](#), Workshop on the Mathematics of Drug Design/Discovery, June 4 - 8, 2018, The Fields Institute.
- Courses available at the University of Basel and beyond.

# Molecular similarity and similarity measures

Chemical similarity	<table><tr><td></td><td>Mol. weight</td><td>LogP</td><td>Rotatable bonds</td><td>Aromatic rings</td><td>Heavy atoms</td></tr><tr><td>A</td><td>341.4</td><td>5.23</td><td>4</td><td>4</td><td>26</td></tr><tr><td>B</td><td>463.5</td><td>4.43</td><td>4</td><td>5</td><td>35</td></tr></table>		Mol. weight	LogP	Rotatable bonds	Aromatic rings	Heavy atoms	A	341.4	5.23	4	4	26	B	463.5	4.43	4	5	35
	Mol. weight	LogP	Rotatable bonds	Aromatic rings	Heavy atoms														
A	341.4	5.23	4	4	26														
B	463.5	4.43	4	5	35														
Molecular similarity																			
2D similarity																			
3D similarity																			
Biological similarity	<table><tr><td></td><td>Vascular endothelial growth factor receptor 2</td><td>Tyrosine-protein kinase TIE-2</td></tr><tr><td>A</td><td>active</td><td>inactive</td></tr><tr><td>B</td><td>active</td><td>active</td></tr></table>		Vascular endothelial growth factor receptor 2	Tyrosine-protein kinase TIE-2	A	active	inactive	B	active	active									
	Vascular endothelial growth factor receptor 2	Tyrosine-protein kinase TIE-2																	
A	active	inactive																	
B	active	active																	
Global similarity																			
Local similarity																			

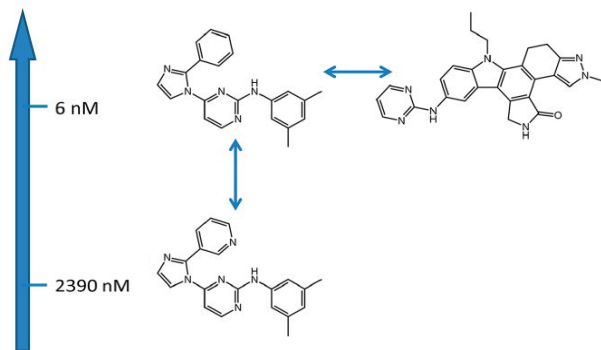
**Table 2 Formulas for the various similarity and distance metrics**

Distance metric	Formula for continuous variables <sup>a</sup>	Formula for dichotomous variables <sup>a</sup>
Manhattan distance	$D_{A,B} = \sum_{j=1}^n  x_{jA} - x_{jB} $	$D_{A,B} = a + b - 2c$
Euclidean distance	$D_{A,B} = \left[ \sum_{j=1}^n (x_{jA} - x_{jB})^2 \right]^{1/2}$	$D_{A,B} = [a + b - 2c]^{1/2}$
Cosine coefficient	$S_{A,B} = \left[ \sum_{j=1}^n x_{jA} x_{jB} \right] / \left[ \sum_{j=1}^n (x_{jA})^2 \sum_{j=1}^n (x_{jB})^2 \right]^{1/2}$	$S_{A,B} = \frac{c}{[ab]^{1/2}}$
Dice coefficient	$S_{A,B} = \left[ 2 \sum_{j=1}^n x_{jA} x_{jB} \right] / \left[ \sum_{j=1}^n (x_{jA})^2 + \sum_{j=1}^n (x_{jB})^2 \right]$	$S_{A,B} = 2c/[a + b]$
Tanimoto coefficient	$S_{A,B} = \frac{\left[ \sum_{j=1}^n x_{jA} x_{jB} \right]}{\left[ \sum_{j=1}^n (x_{jA})^2 + \sum_{j=1}^n (x_{jB})^2 - \sum_{j=1}^n x_{jA} x_{jB} \right]}$	$S_{A,B} = c/[a + b - c]$
Soergel distance <sup>b</sup>	$D_{A,B} = \left[ \sum_{j=1}^n  x_{jA} - x_{jB}  \right] / \left[ \sum_{j=1}^n \max(x_{jA}, x_{jB}) \right]$	$D_{A,B} = 1 - \frac{c}{[a+b-c]}$

$S$  denotes similarities, while  $D$  denotes distances. The two can be converted to each other by *similarity* =  $1/(1 + \text{distance})$ .  $x_{jA}$  means the  $j$ -th feature of molecule A.  $a$  is the number of *on* bits in molecule A,  $b$  is number of *on* bits in molecule B, while  $c$  is the number of bits that are *on* in both molecules.

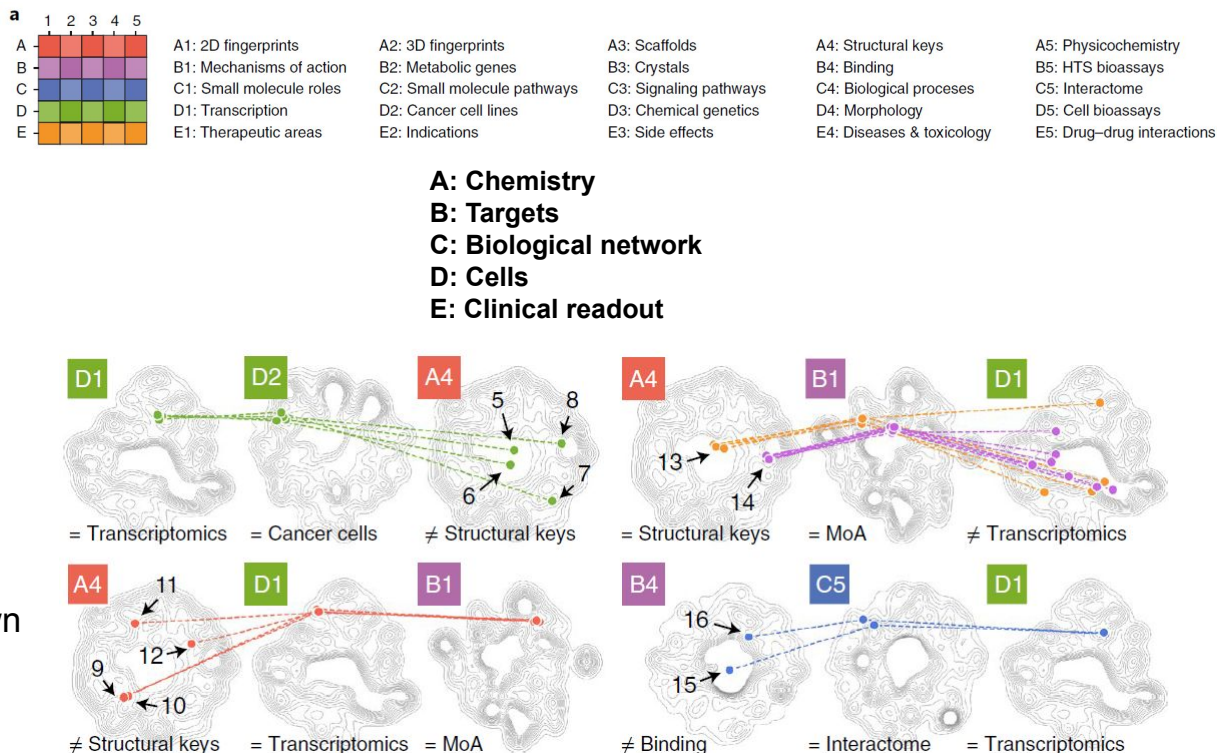
(Left) Maggiora, Gerald, Martin Vogt, Dagmar Stumpfe, und Jürgen Bajorath. „[Molecular Similarity in Medicinal Chemistry](#)“. *Journal of Medicinal Chemistry* 57, Nr. 8 (24. April 2014): 3186–3204. (Right) Bajusz, Dávid, Anita Rácz, and Károly Héberger. 2015. “[Why Is Tanimoto Index an Appropriate Choice for Fingerprint-Based Similarity Calculations?](#)” *Journal of Cheminformatics* 7 (1): 20.

# Molecular similarity does not equal biological similarity



## Watch out biological activity cliffs!

Similarity does not imply activity. Three vascular endothelial growth factor receptor 2 (VEGFR2) ligands are shown that represent different similarity–activity relationships.

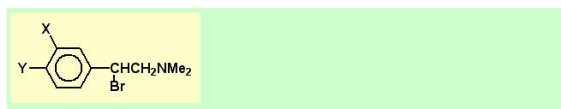


Duran-Frigola, Miquel, Eduardo Pauls, Oriol Guitart-Pla, Martino Bertoni, Víctor Alcalde, David Amat, Teresa Juan-Blanco, and Patrick Aloy. 2020. “[Extending the Small-Molecule Similarity Principle to All Levels of Biology with the Chemical Checker](#).” Nature Biotechnology, May, 1–10.

# Quantitative Structure-Activity Relationships (QSARs)

QSAR is a statistical modelling of correlation between biological activity and physicochemical properties, or  $\Delta\phi=f(\Delta S)$ , where  $\phi$  indicates a biological activity and S indicates a chemical structure (1868-1869).

An example: **The Free-Wilson analysis**. The assumption: the biological activity for a set of analogues could be described by the contributions that substituents or structural elements make to the activity of a parent structure.



**Molecular Descriptors (MD)**

	Target property	MD <sub>1</sub>	MD <sub>2</sub>	...	MD <sub>M</sub>
C <sub>1</sub>	y <sub>1</sub>	x <sub>1,1</sub>	x <sub>1,2</sub>	...	x <sub>1,M</sub>
C <sub>2</sub>	y <sub>2</sub>	x <sub>2,1</sub>	...	...	...
C <sub>3</sub>	y <sub>3</sub>	...	...	...	...
C <sub>4</sub>	y <sub>4</sub>	...	...	...	...
...	...	...	...	...	...
...	...	...	...	...	...
C <sub>N</sub>	y <sub>N</sub>	x <sub>N,1</sub>	x <sub>N,2</sub>	...	x <sub>N,M</sub>

The basic form of a QSAR model: find a function  $f$  that predicts  $y$  from  $x$ ,  $y \sim f(x)$

meta	para	meta-					para-					log 1/C	log 1/C
(X)	(Y)	F	Cl	Br	I	Me	F	Cl	Br	I	Me	obsd.	calc.a
H	H											7.46	7.82
H	F						1					8.16	8.16
H	Cl							1				8.68	8.59
H	Br								1			8.89	8.84
H	I									1		9.25	9.25
H	Me										1	9.30	9.08
F	H	1										7.52	7.52
Cl	H		1									8.16	8.03
Br	H			1								8.30	8.26
I	H				1							8.40	8.40
Me	H					1						8.46	8.28
Cl	F		1				1					8.19	8.37
Br	F			1				1				8.57	8.60
Me	F					1	1						
Cl	Cl		1						1				
Br	Cl			1						1			
Me	Cl					1					1		
Cl	Br		1										
Br	Br			1									
Me	Br					1							
Me	Me						1						
Br	Me			1									

log (1/ED<sub>50</sub>) = -

+

+

n

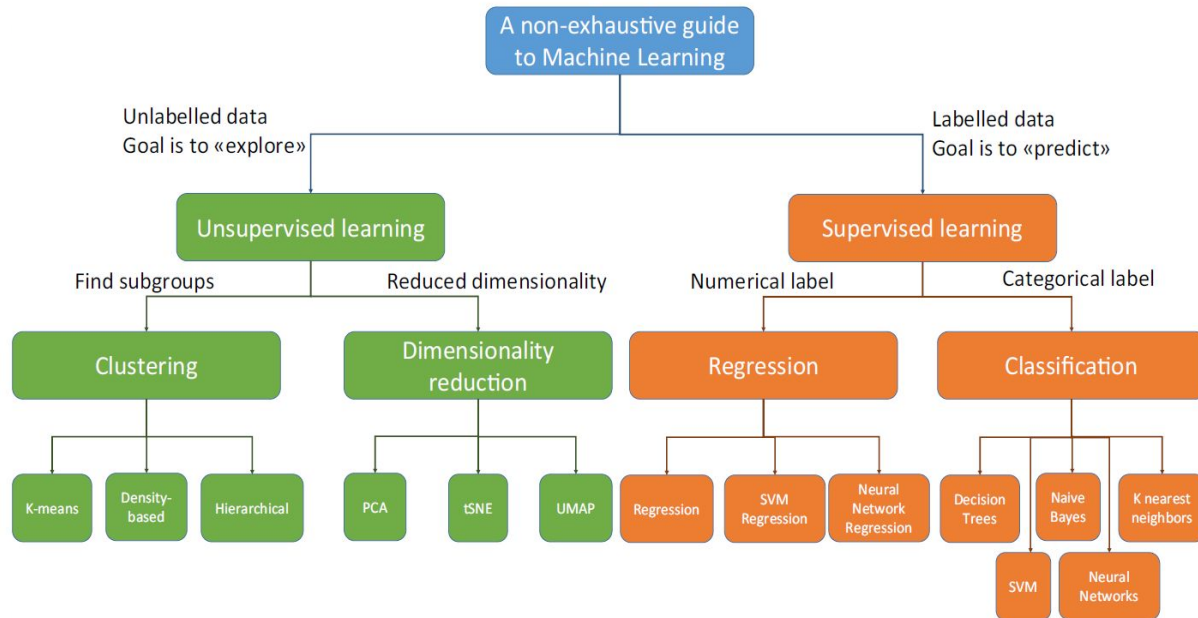
Multivariate regression analysis

$$\log (1 / \text{ED}_{50}) = -0.301[m-F] + 0.27[m-Cl] + 0.434[m-Br] + 0.579[m-I] \\ + 0.454[m-Me] + 0.340[p-F] + 0.768[p-Cl] + 1.020[p-Br] \\ + 1.429[p-I] + 1.256[p-Me] + 7.821 \\ n = 22, r^2 = 0.94, s = 0.194, F = 17.0$$

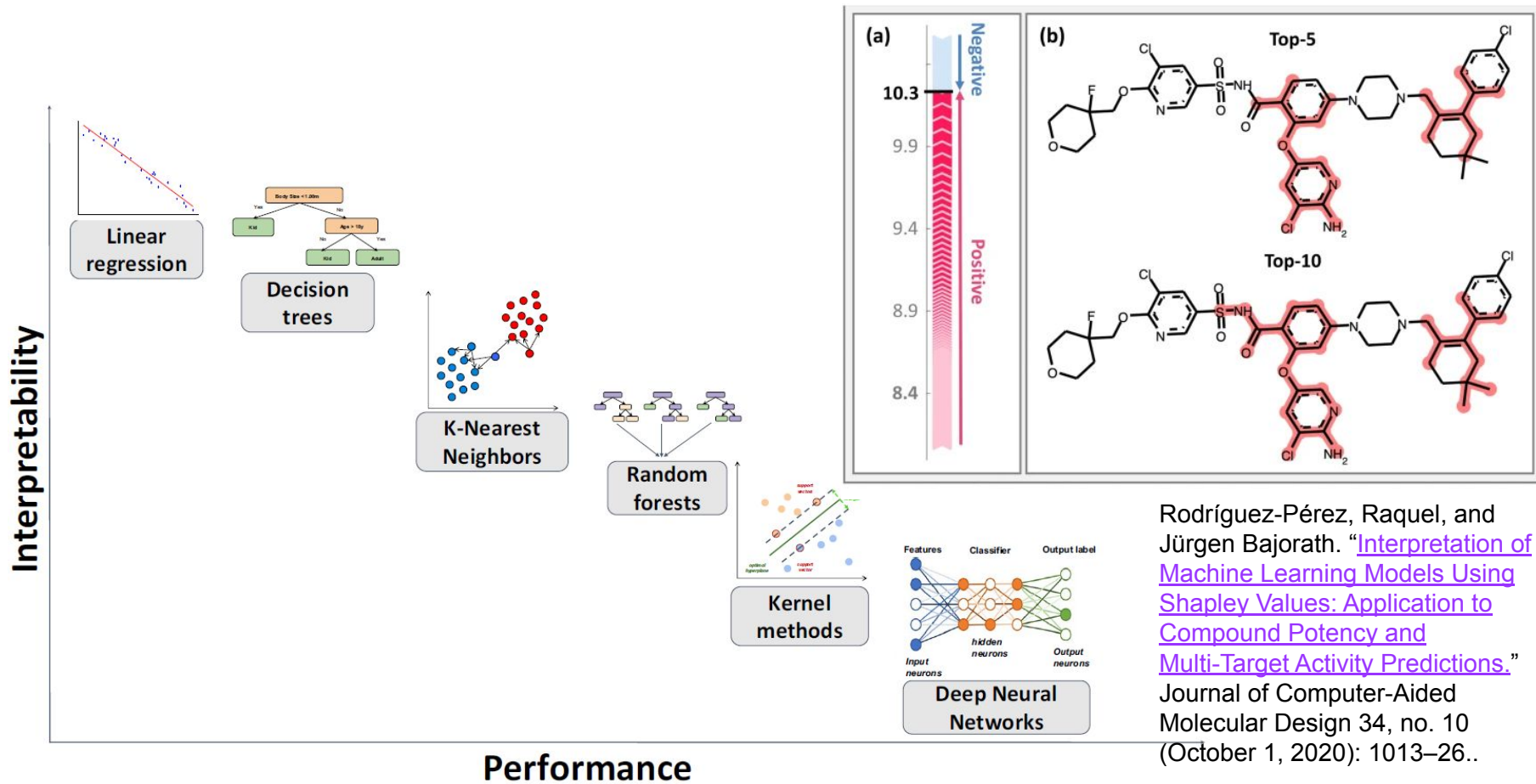


# QSAR models mark the early adoption of machine learning in drug discovery

- QSAR is among the earliest subjects that used machine learning and pattern recognition in drug discovery.
- **Advantages:** technically easy, fast, and many models are useful as filters.
- **Disadvantages:** statistical models cannot capture mechanistic aspects of biochemical interactions, limited ability to debug when a model fails to work, and findings may not be generalizable.

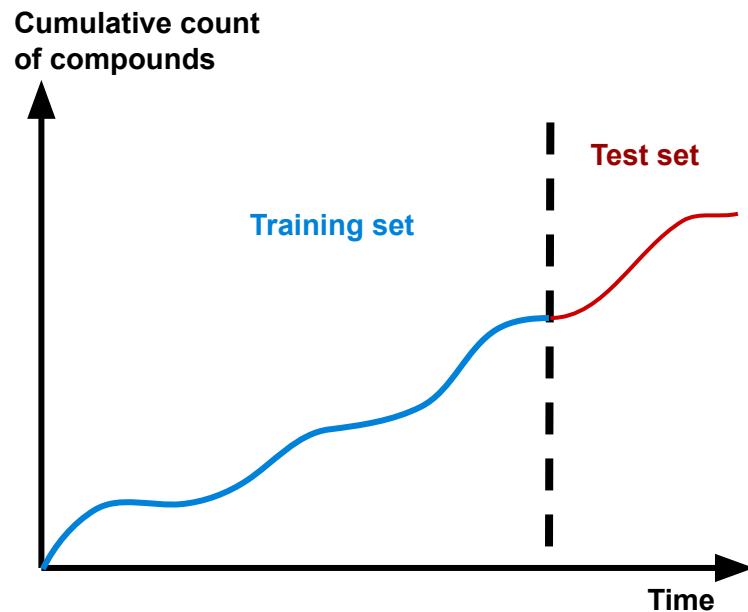
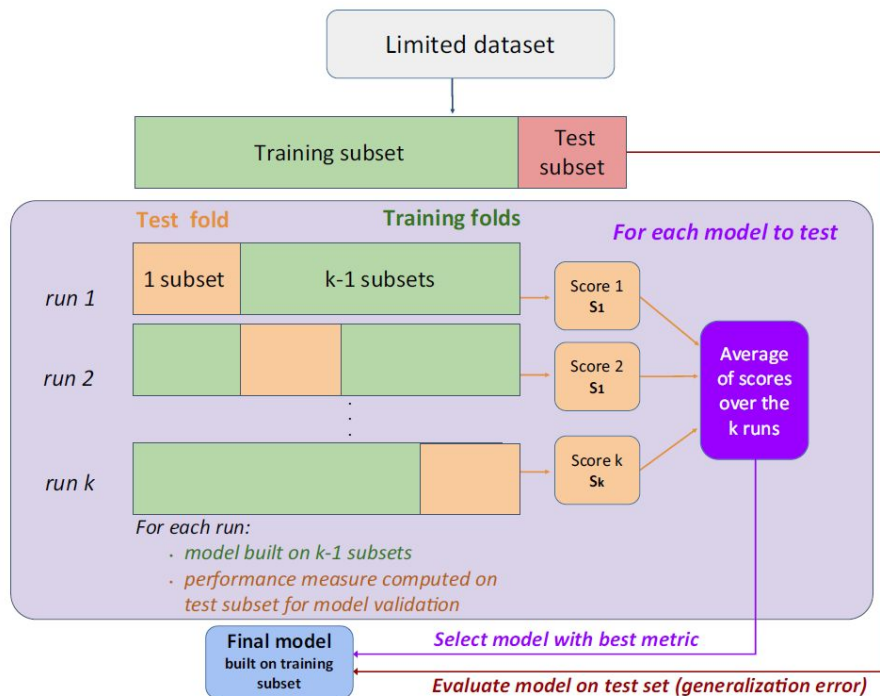


# Interpretable and Causal Models will become more important



Rodríguez-Pérez, Raquel, and Jürgen Bajorath. "[Interpretation of Machine Learning Models Using Shapley Values: Application to Compound Potency and Multi-Target Activity Predictions.](#)"  
Journal of Computer-Aided Molecular Design 34, no. 10 (October 1, 2020): 1013–26..

# The general practice of training a supervised learning model

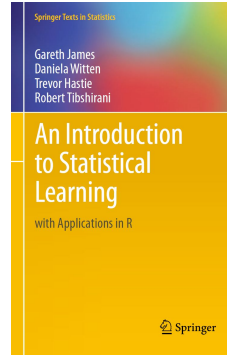
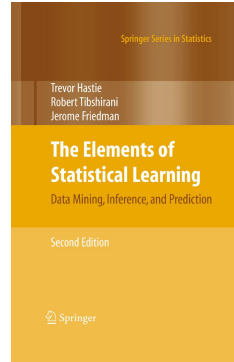


(Left) To assess the generalization ability of a supervised learning algorithm, data are separated into a training subset used for building the model and a test subset used to assess the generalization error (from Badillo *et al.*, 2020) (Right) Temporal validation is especially important for drug discovery, because chemical structures used in the training set may differ substantially from those that will be tested.

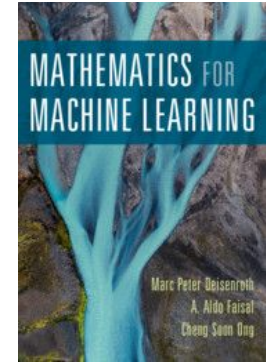


# Resources for learning about machine learning

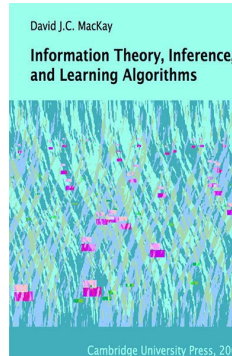
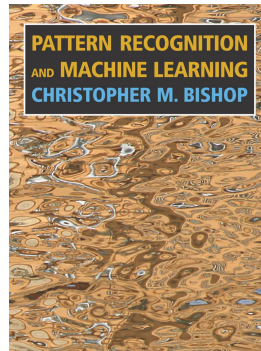
ESL and ISL: From a frequentist view (almost)



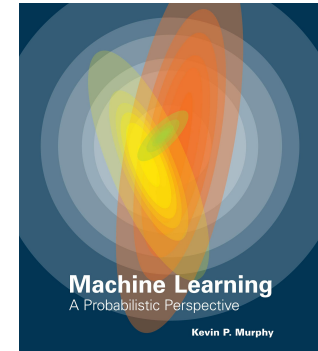
Mathematical foundations



PRML and ITILA: From a Bayesian view



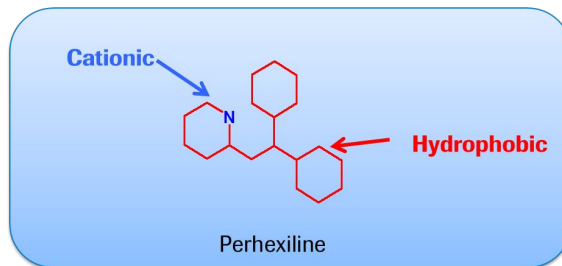
MLaPP: Application oriented, more accessible, and balanced views



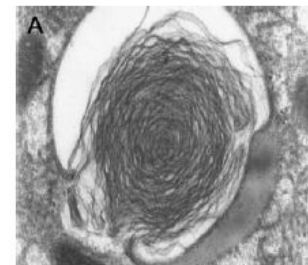
# Drug-induced phospholipidosis is correlated with amphiphilicity

- Phospholipidosis is a lysosomal storage disorder characterized by the excess accumulation of phospholipids in tissues.
- Drug-induced phospholipidosis is caused by cationic amphiphilic drugs and some cationic hydrophilic drugs.
- Clinical pharmacokinetic characteristics of drug-induced phospholipidosis include (1) very long terminal half lives, (2) high volume of distribution, (3) tissue accumulation upon frequent dosing, and (4) deficit in drug metabolism.

Fischer *et al.* (Chimia 2000) discovered that it is possible to predict the amphiphilicity property of druglike molecules by calculating the amphiphilic moment using a simple equation.



Lüllmann *et al.*, Drug Induced Phospholipidosis, *Crit. Rev. Toxicol.* 4, 185, 1975



Anderson and Borlak, Drug-Induced Phospholipidosis, *FEBS Letters* 580, Nr. 23 (2006): 5533–40.

$$\vec{A} = \sum_i d \cdot \vec{\alpha}_i$$

$\vec{A}$ : Calculated amphiphilic moment

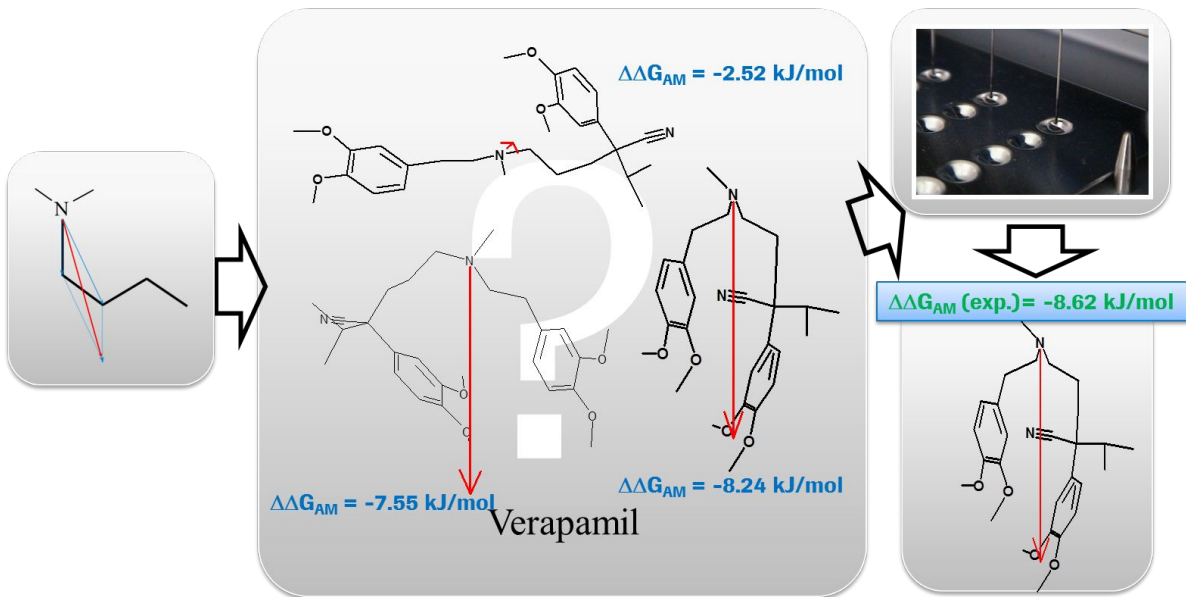
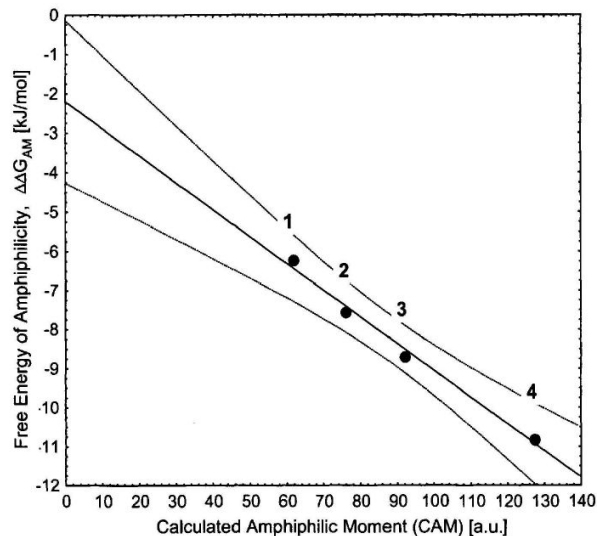
$d$ : distance between the center of gravity of the charged part of a molecule and the hydrophobic/hydrophilic remnant of the molecule

$\vec{\alpha}_i$ : the hydrophobic/hydrophilic contribution of atom/fragment  $i$

***In silico* calculation of amphiphilicity property may be used to predict phospholipidosis induction potential**

# In silico prediction of amphiphilicity

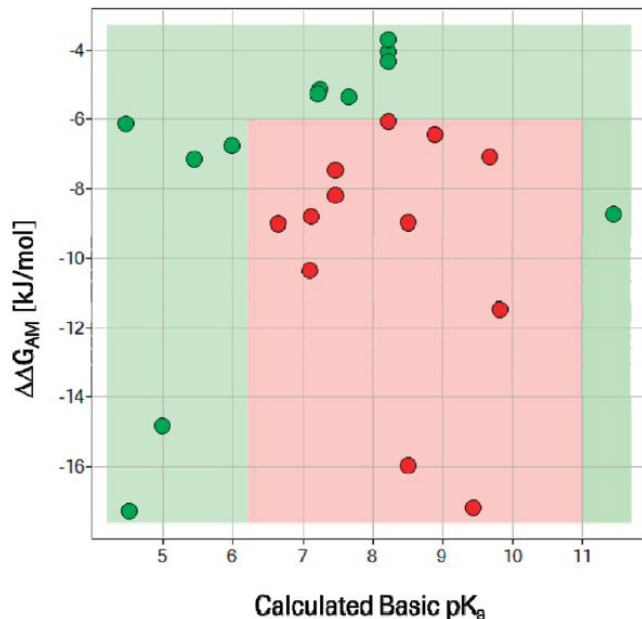
Development of CAFCA (CAlculated Free energy of amphiphilicity of small Charged Amphiphiles)



Iterative model building, experimentation, and model refining led to the predictive tool CAFCA

# Validation of in silico phospholipidosis prediction

Model Validation from 1999-2004



Plot of amphiphilicity ( $\Delta\Delta G_{AM}$ ) versus calculated basic  $pK_a$  for the training set of 24 compounds. The red area defines the region where a positive PLD response is expected, and the green area defines where a negative response is expected according to the tool.

in vitro/ in vivo	in silico/ in vivo	Exp. PC/ in vivo	In silico/ in vitro	n=36
94%	81%	89%	89%	

in vitro/in silico			n=422
Accuracy [(TP+TN)/ (P+N)]	Sensitivity [True Positive Rate]	Specificity [True Negative Rate]	Precision [TP/(TP+FP)]
86%	80%	90%	84%

Fischer et al., *J. Med. Chem.*, 55 (1),  
2012

**We gained mechanistic insights of phospholipidosis induction by cationic amphiphilic drugs with the model**

# Phospholipidosis: lessons learned (and lessons not yet learned)

- Cationic amphiphilic properties of a molecule is an early marker for safety in drug discovery and early development.
  - Phospholipidosis in dose range finding studies
  - Cardiac ion channel interactions (hERG, sodium channel, ...)
  - Receptor binding promiscuity
  - P-gp inhibition
  - Mitochondrial toxicity in case of safety relevant findings, e.g. in dose range finding studies
- Extreme basic amphiphilic properties should be avoided because of a higher risk of PLD, QT-prolongation, mitochondrial toxicity. However, basic compounds with moderate amphiphilic properties are still a preferred scaffold for many therapeutic areas (especially CNS).
- **Generally, some safety liabilities, despite complex underlying biological and chemical mechanisms, can be predicted by molecular modelling well, sometimes with surprisingly elegant models!**

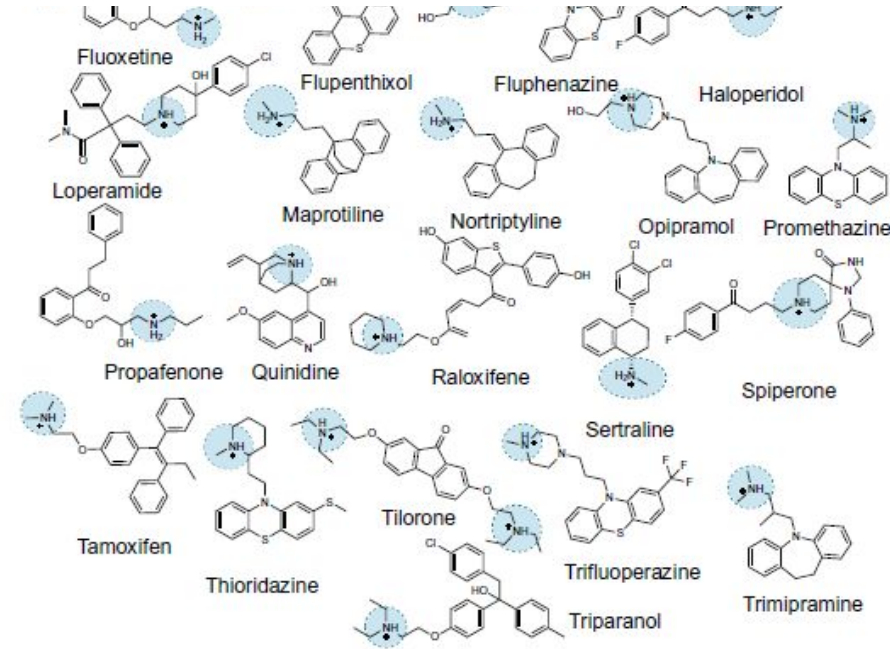
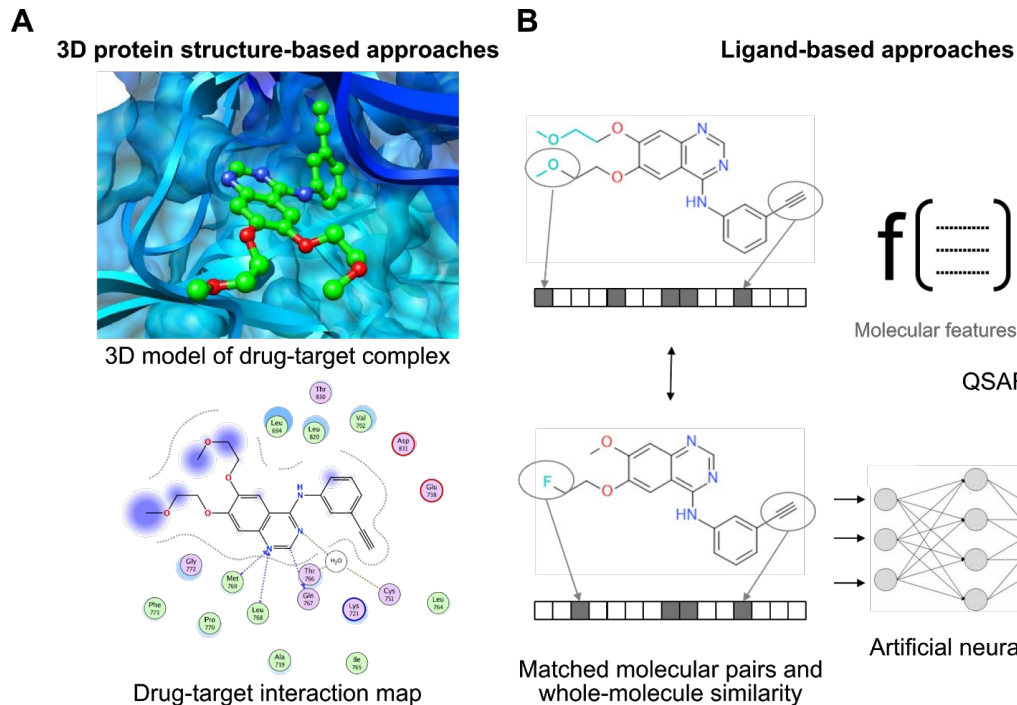


Fig. 1. Representative examples of CADs that are identified in SARS-CoV-2 drug repurposing screens.

Tummino, Tia A., Veronica V. Rezelj, Benoit Fischer, Audrey Fischer, Matthew J. O'Meara, Blandine Monel, Thomas Vallet, et al. "Drug-Induced Phospholipidosis Confounds Drug Repurposing for SARS-CoV-2." *Science* 373, no. 6554 (July 30, 2021): 541–47. <https://doi.org/10.1126/science.abi4708>.

# Summary and Q&A



Overview of  
**non-sequence-based,  
molecular-level modelling  
techniques:** (A) 3D protein  
structure-based approaches (B)  
Ligand-based approaches.

Zhang, Jitao David, Lisa  
Sach-Peltason, Christian Kramer,  
Ken Wang, and Martin Ebeling.  
2020. "[Multiscale Modelling of  
Drug Mechanism and Safety](#)."  
*Drug Discovery Today* 25 (3):  
519–34.

# Offline activities

- Read selected pages of *Computational Methods in Drug Discovery* by Sliwoski *et al.* Please submit your results to the Google Form, the link of which will be sent via a separate email.
- Optional and recommended:
  - Fill the anonymous survey #6 (link will be sent via a separate email).
  - Recommended readings:
    - Badillo *et al.* 2020. “[An Introduction to Machine Learning](#).” Clinical Pharmacology & Therapeutics.
    - Jiménez-Luna, José, Francesca Grisoni, and Gisbert Schneider. 2020. “[Drug Discovery with Explainable Artificial Intelligence](#).” Nature Machine Intelligence 2 (10): 573–84..

## More about the the Free-Wilson analysis

- [\*A Mathematical Contribution to Structure-Activity Studies\*](#) by Spencer M. Free and James W. Wilson, Journal of Medicinal Chemistry, 1964, and reviewed by [Kubinyi](#), 1988.
- A Python implementation on [GitHub](#), and a [blog post](#) going through examples, is shared by Pat Walters.
- Free-Wilson nonadditivity is a research topic, for instance see [Cramer et al., 2015](#)
- Source of the example shown in the lecture: QSAR of the [ACCVIP](#) project (The Australian Computational Chemistry via the Internet Project)



# Resources about the mathematics underlying molecular structure determination

## • Mathematical and physical foundations

- Recommended reading: [Mathematical techniques used in biophysics](#)
- [Background on imaging physics](#) at xrayphysics.com
- [Physics for life-science students](#) at U Maryland

## • X-ray diffraction by electrons

- An [AMS Feature Column](#) by Tony Phillips
- Stanford open course [Fourier transform and its applications](#)

## • Nuclear Magnetic Resonance (NMR)

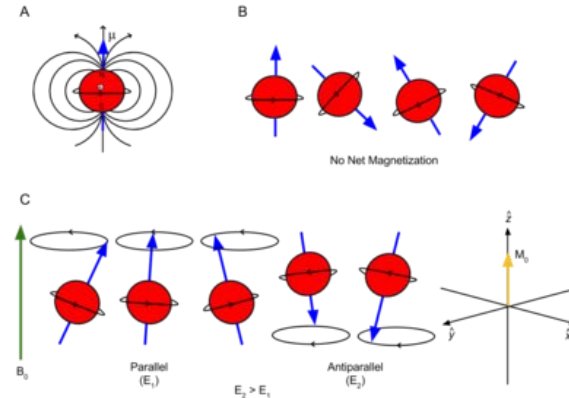
- [A beautiful video tutorial](#) about the principles of magnetic resonance imaging (MRI), which is a variant of NMR

## • Cryo-electron microscopy (CryoEM)

- [A three-minute introduction to CryoEM](#)
- [Nobel Prize Talk by Joachim Frank](#)
- [Talk on Mathematics of CryoEM](#), by Prof Amit Singer, with a manuscript available at arXiv: <https://arxiv.org/abs/1803.06714>



Swiss Light Source, the synchrotron at the Paul Scherrer Institute (PSI), copyright of PSI



Adapted from Bushberg JT, [The Essential Physics of Medical Imaging](#): Lippincott Williams & Wilkins; 2002

