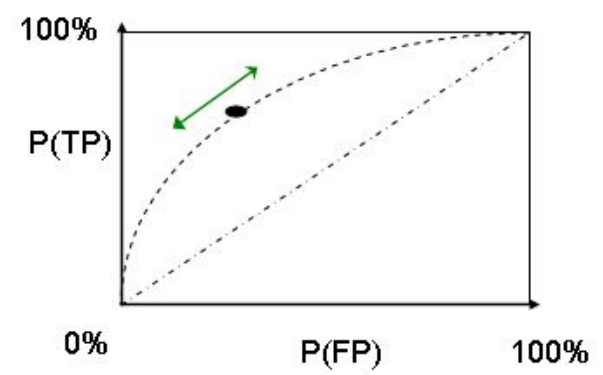
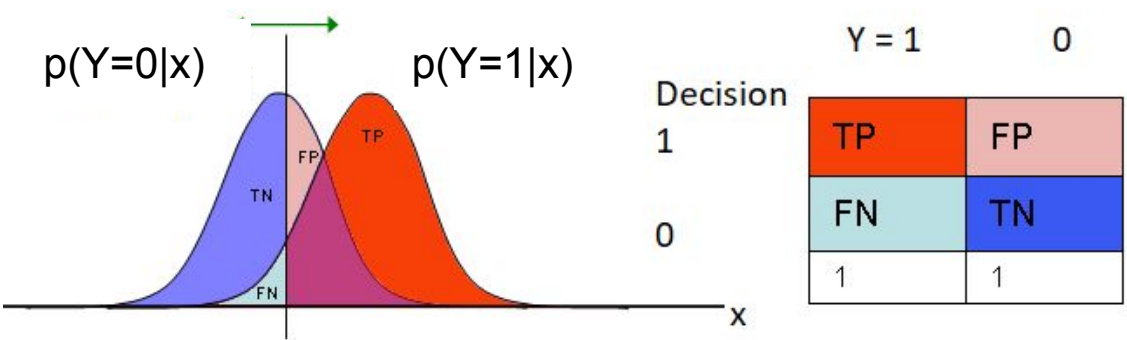


# Offline activities of Lecture 6

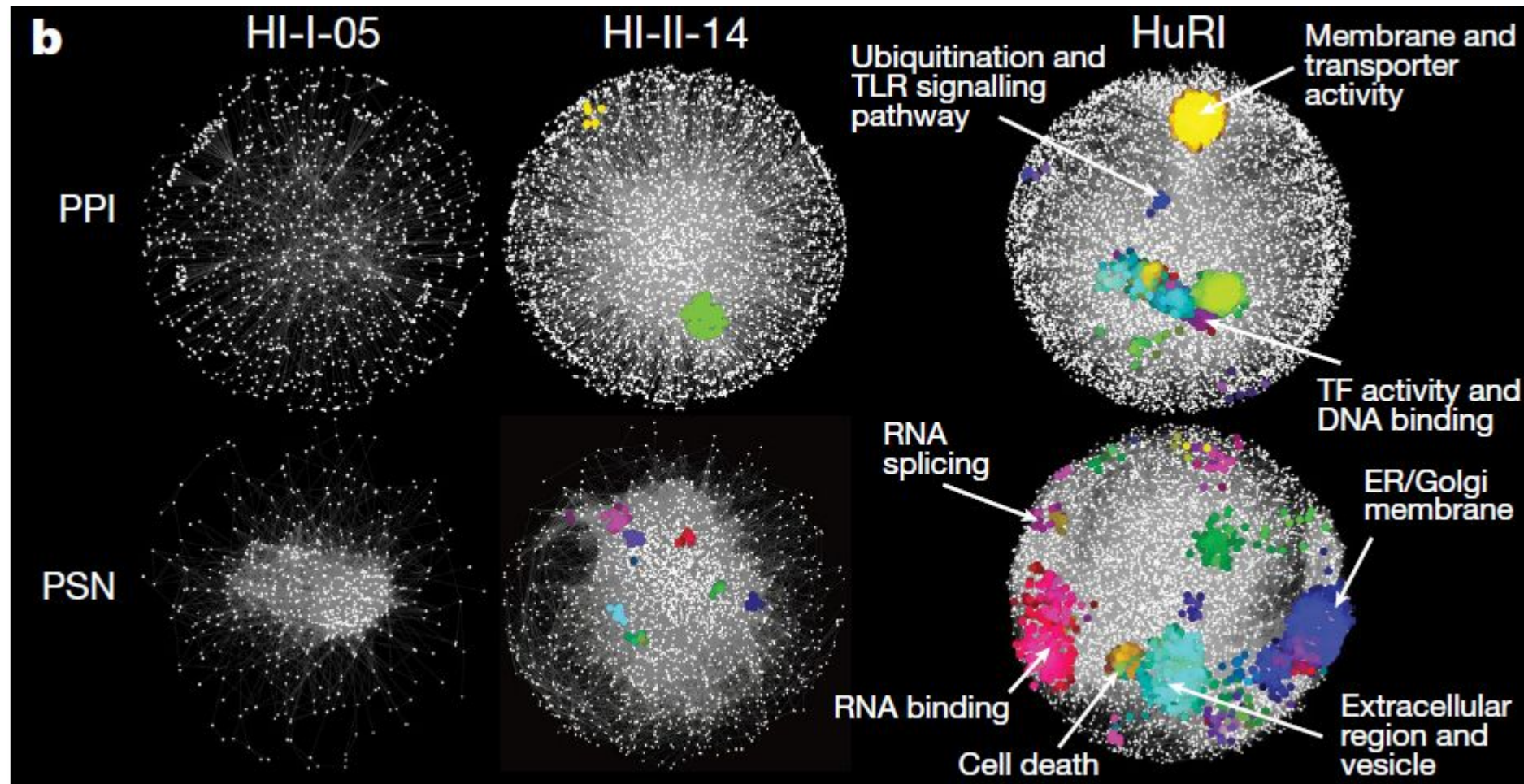
- Q1: What is the method commonly used to benchmark performance of different techniques of computer-aided drug design (CADD)? ([Receiver Operating Characteristic curves](#))
- Q2: What do we mean by molecular dynamics? ([A computer simulation method to analyze the movements of atoms and molecules using Newtonian mechanistics](#))
- Q3: What are the three basic methods to represent target and ligand structures *in silico*? ([atomic, surface, and grid representations](#))
- Q4: What sampling algorithms are there for protein-ligand docking? Can you explain one of them using your words? ([systematic algorithms, molecular dynamics simulations, Monte Carlo search with Metropolis Criterion and genetic algorithms](#))
- Q5: What are the key steps in structure-based virtual high-throughput screenings (SB-vHTS)? ([preparing structures, posing, scoring](#))
- Q6: What is the usual starting point of structure-based CADD campaign? ([Experimentally determined protein structures, preferably in complex with ligands](#))
- Q7: What do we mean by 'pharmacophore'? ([model of the target binding site which summarizes steric and electronic features needed for optimal interaction of a ligand with a target, a "subgraph" of a molecule with interesting properties for drug design/protein binding](#))
- Q8: In QSAR analysis, why it is important to select optimal descriptors/features? ([to reduce noise, to increase generalized performance, and for hypothesis generation](#))
- Q9: What do we mean by the acronyms *DMPK* and *ADMET*? ([DMPK=drug metabolism and pharmacokinetics; ADMET= absorption, distribution, metabolism, excretion, and the potential for toxicity](#))
- Q10: Why common CADD methods have difficulties handling protein-protein interaction and protein-DNA interactions? ([large interaction size, lack of user-friendly tools, and comparably little training data](#))

# Question about the ROC curve



|                                 |   | Predicted labels                |                                           |                                                                                      |
|---------------------------------|---|---------------------------------|-------------------------------------------|--------------------------------------------------------------------------------------|
|                                 |   | 1                               | 0                                         |                                                                                      |
| Actual labels<br>(observations) | 1 | True Positive (TP)              | False Negative (FN)                       | Recall=TPR<br>(True Positive Rate)<br>$TPR = \frac{TP}{TP+FN}$                       |
|                                 | 0 | False Positive (FP)             | True Negative (TN)                        | Specificity = $\frac{TN}{TN+FP}$<br>False Positive Rate:<br>$FPR = \frac{FP}{FP+TN}$ |
|                                 |   | Precision<br>$\frac{TP}{TP+FP}$ | False Negative Rate<br>$\frac{FN}{TN+FN}$ | Accuracy<br>$\frac{TP+TN}{TP+TN+FP+FN}$                                              |

# AMIDD Lecture 7: From individual interactions to networks



Luck *et al.* “[A Reference Map of the Human Binary Protein Interactome.](#)”  
Nature, 2020

**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

# Topics

- **From molecular models to cellular models**
- **Gene expression profiling: a case study of omics and cellular modelling**
- **Applications for drug mechanism of action: molecular phenotyping**



# Four classical classes of mathematical models

## Compartment models

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

Kinetics of ligand-target interaction

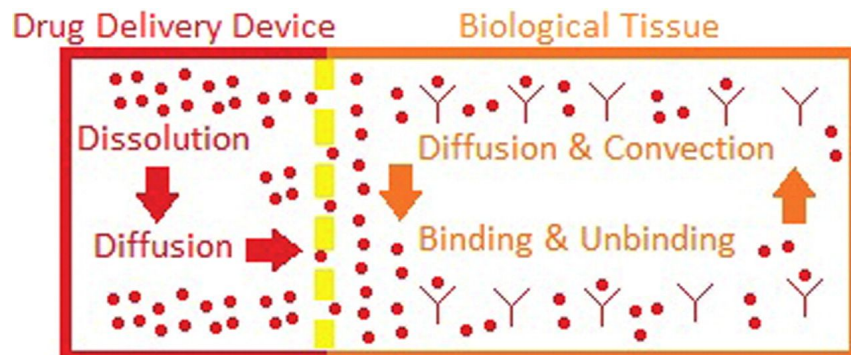
$$\begin{aligned} \frac{dx}{dt} &= \alpha x - \beta xy, \\ \frac{dy}{dt} &= -\gamma y + \delta xy, \end{aligned}$$

The Lotka-Volterra equations modelling predator-prey relationships.

$$\begin{aligned} \frac{dS}{dt} &= -\frac{\beta IS}{N}, \\ \frac{dI}{dt} &= \frac{\beta IS}{N} - \gamma I, \\ \frac{dR}{dt} &= \gamma I \end{aligned}$$

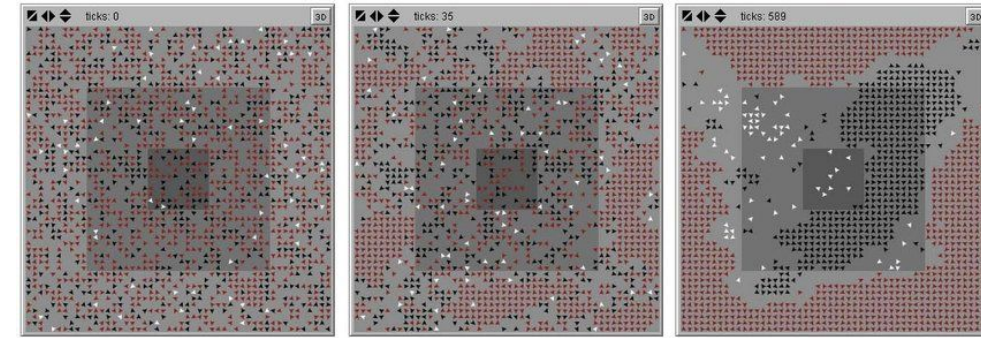
The SIR (S=susceptible, I=infectious, R=removed) model of epidemiology

## Transport models



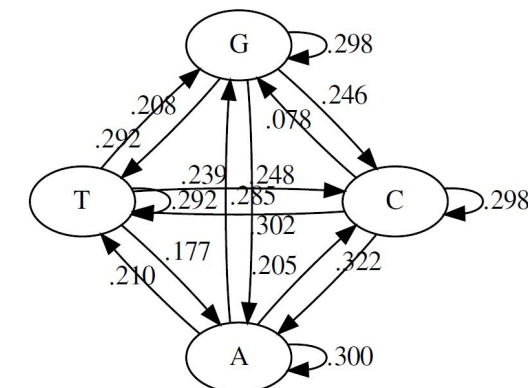
McGinty, Sean, and Giuseppe Pontrelli. 2015. "[A General Model of Coupled Drug Release and Tissue Absorption for Drug Delivery Devices](#)." *Journal of Controlled Release* 217 (November): 327–36.

## Particle models



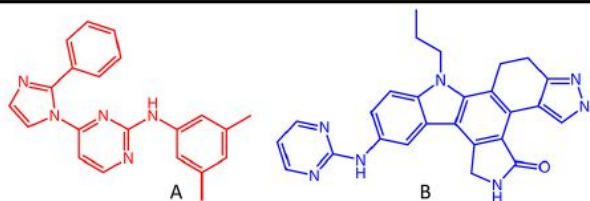
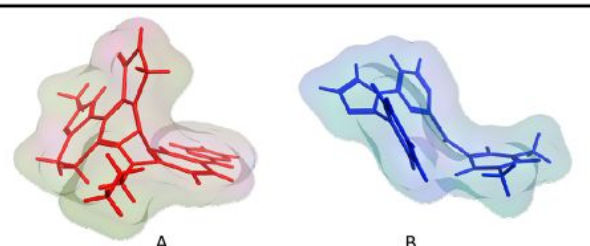
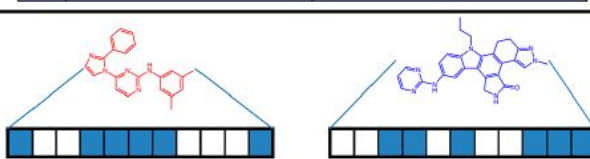
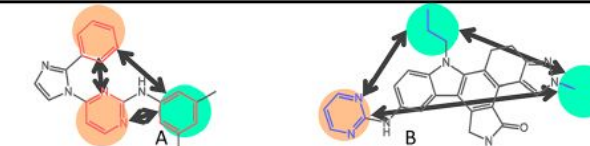
A Study on Socio-spatial Segregation Models Based on Multi-agent Systems by Quadros *et al.* (2012). 10.1109/BWSS.2012.14.

## Finite state models



A finite-state Markov chain modelling DNA sequences

# Molecular similarity and similarity measures

|                       |                                                                                     |                                               |      |                               |                |             |
|-----------------------|-------------------------------------------------------------------------------------|-----------------------------------------------|------|-------------------------------|----------------|-------------|
| Chemical similarity   |                                                                                     | 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 |                                                                                     | 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> | $D_{A,B} = \left[ \sum_{j=1}^n  x_{jA} - x_{jB}  \right] / \left[ \sum_{j=1}^n \max(x_{jA}, x_{jB}) \right]$                                                | $D_{A,B} = 1 - \frac{c}{[a + b - c]}$          |

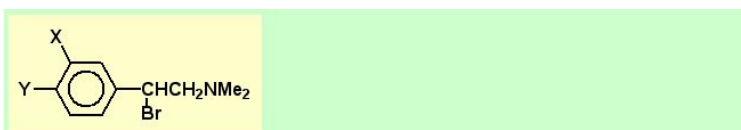
$S$  denotes similarities, while  $D$  denotes distances. The two can be converted to each other by  $similarity = 1/(1 + distance)$ .  $x_{jA}$  means the  $j$ -th feature of molecule A.  $a$  is the number of *on* bits in molecule A,  $b$  is number of *on* bits in molecule B, while  $c$  is the number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

# 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  |   |    |       |    |    |   |    |         |         |

Multivariate regression analysis

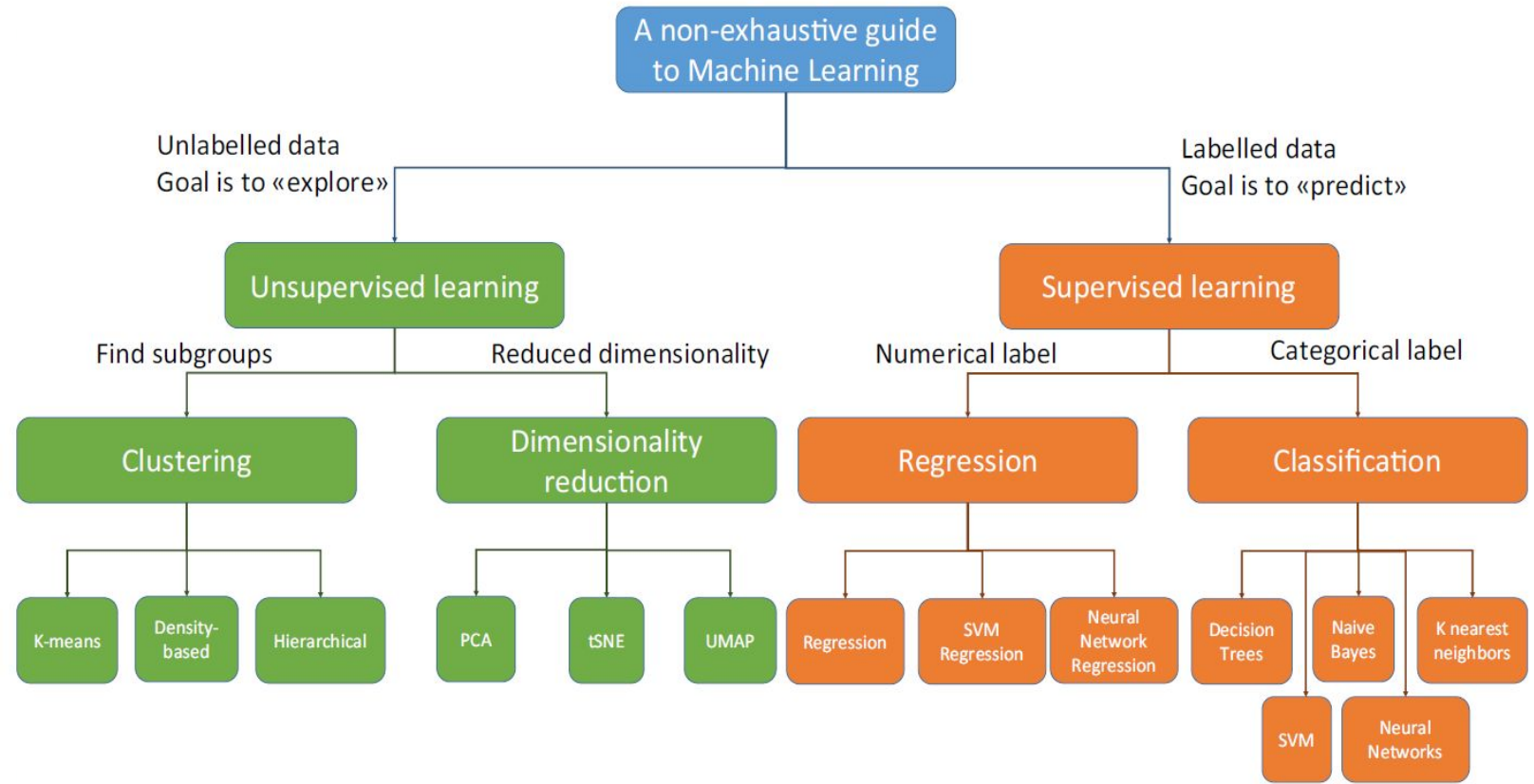
$$\log(1/ED_{50}) = -0.301[m-F] + 0.27[m-Cl] + 0.434[m-Br] + 0.579[m-I] + 0.454[m-Me] + 0.340[p-F] + 0.768[p-Cl] + 1.020[p-Br] + 1.429[p-I] + 1.256[p-Me] + 7.821$$

$n = 22, r^2 = 0.94, s = 0.194, F = 17.0$



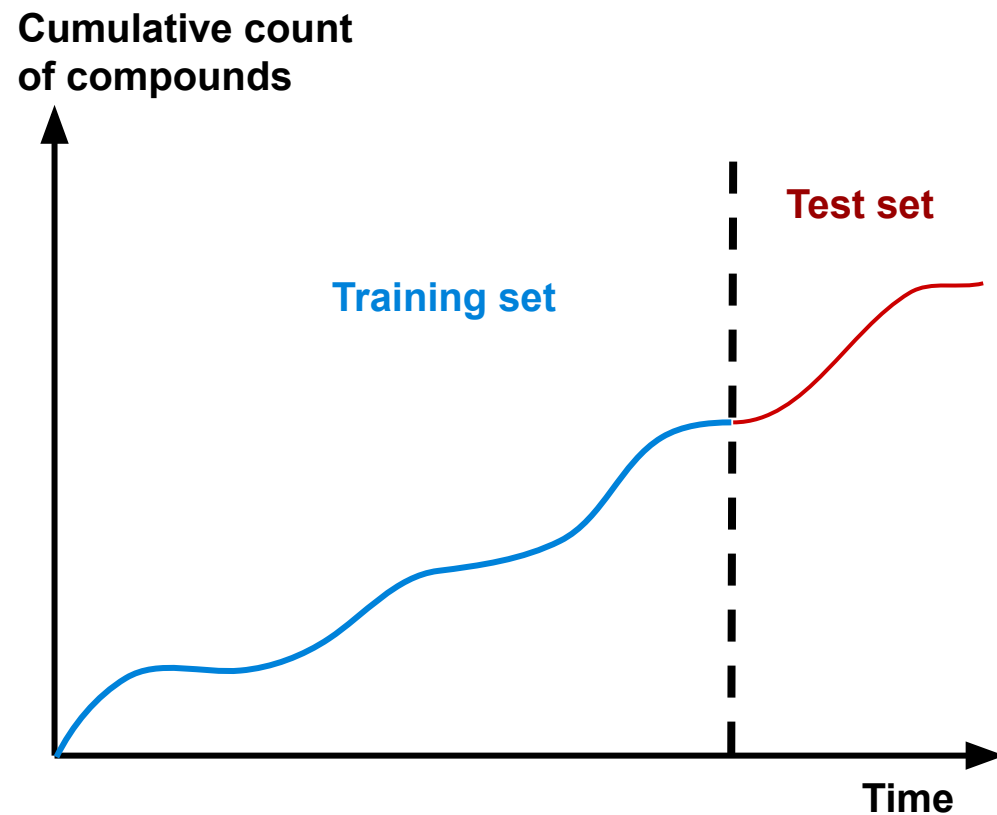
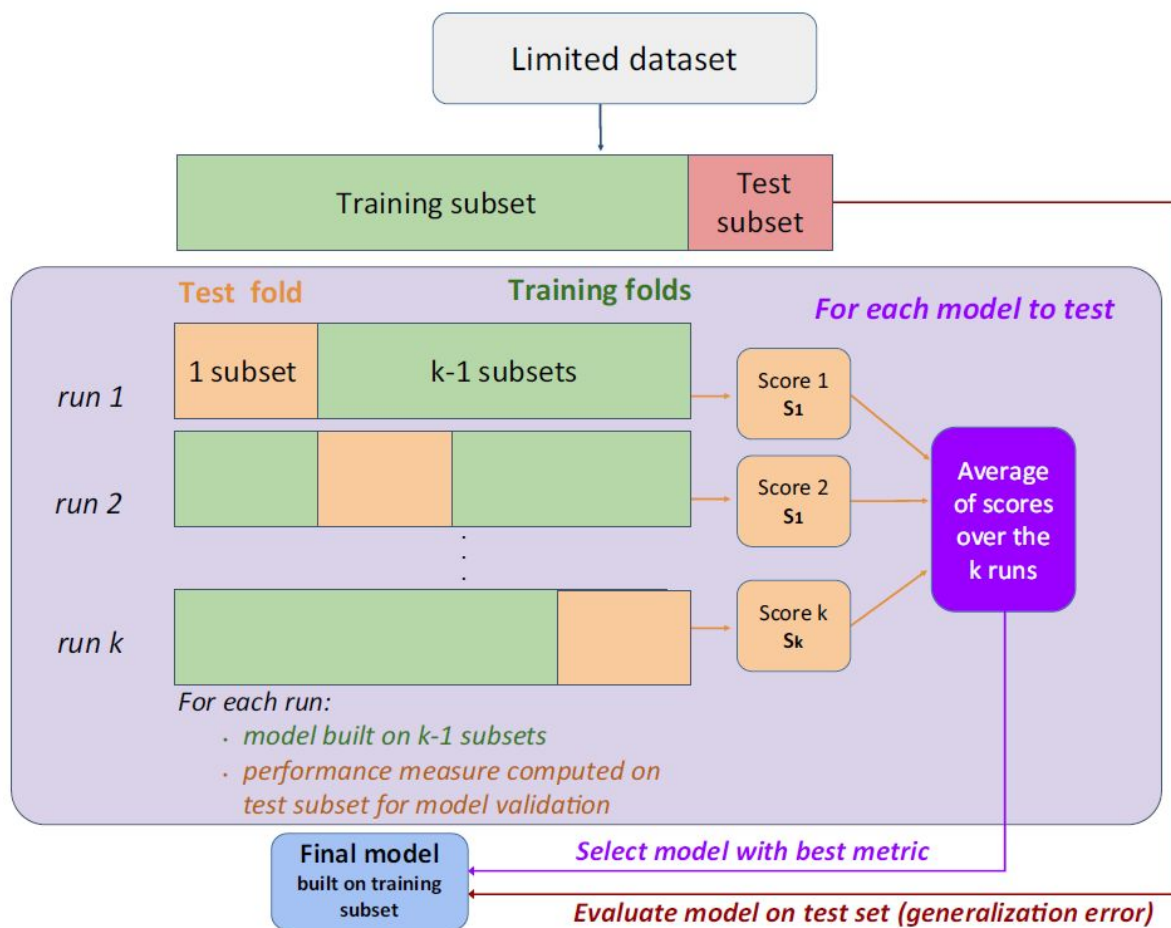
# QSAR models mark the early adoption of statistical modelling and machine learning in drug discovery, the fifth type of mathematical modelling

- 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.



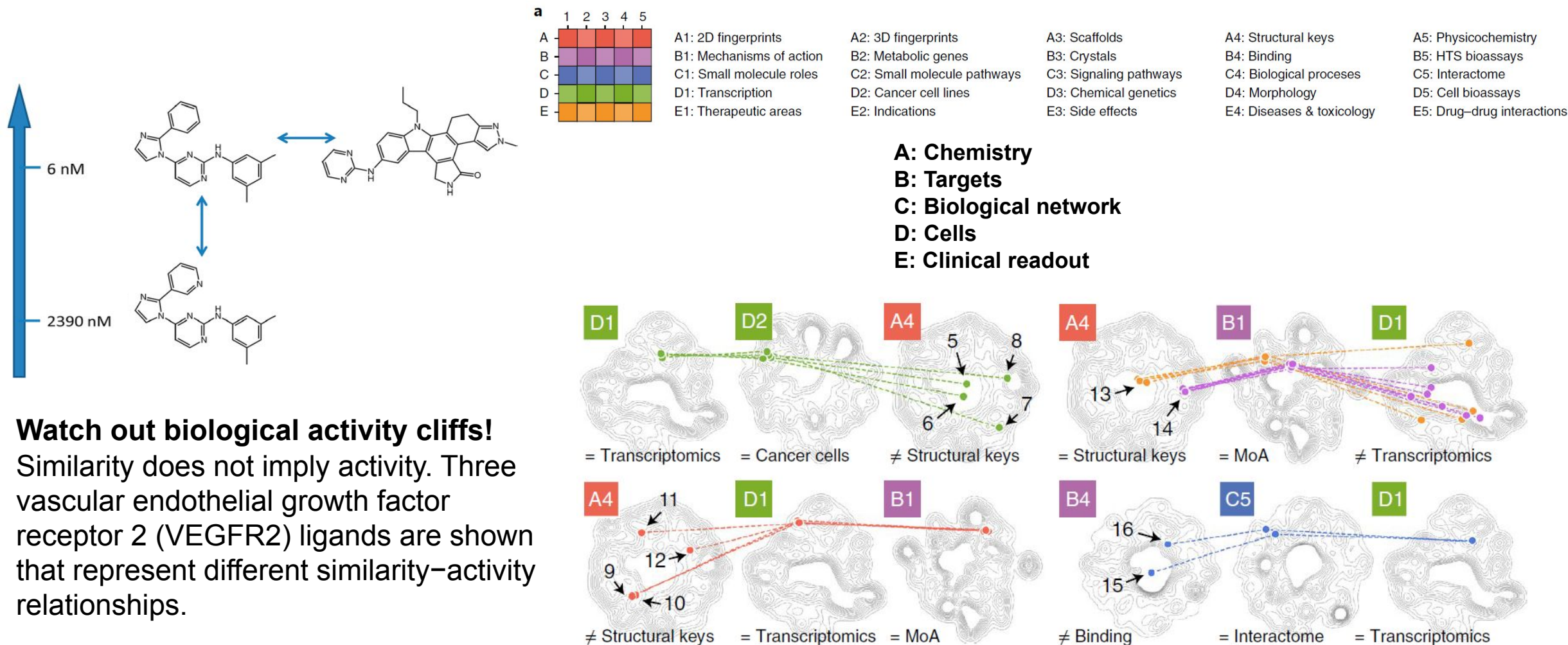


# The general practice of training a supervised learning model



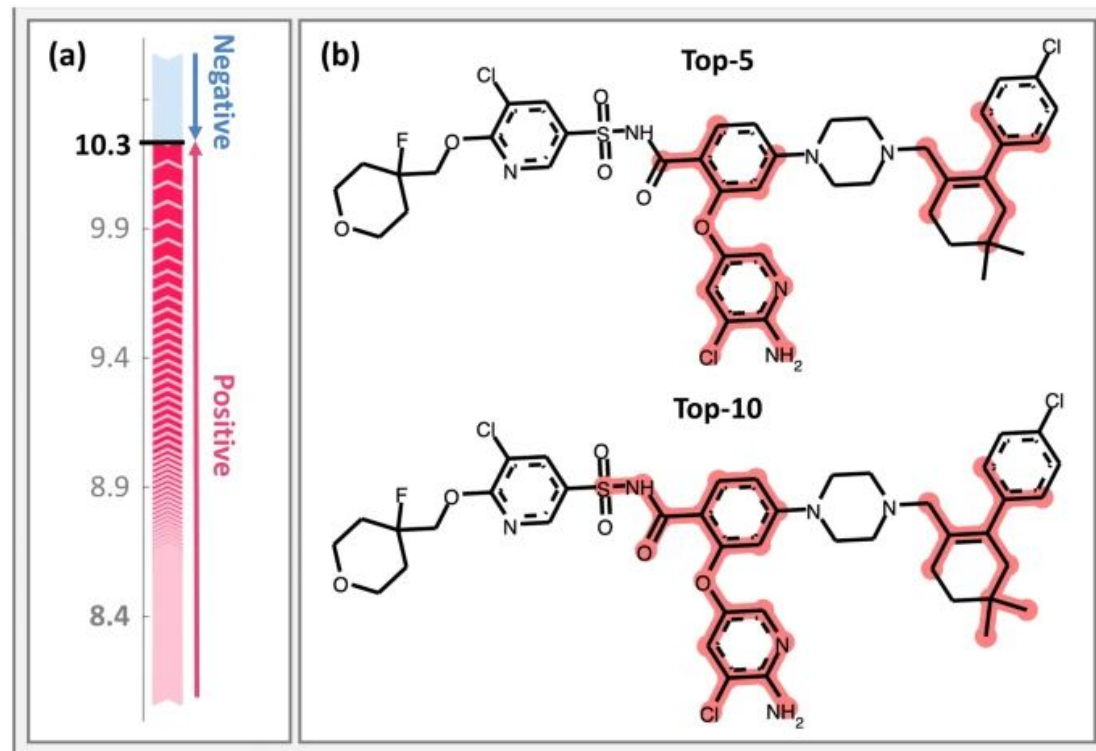
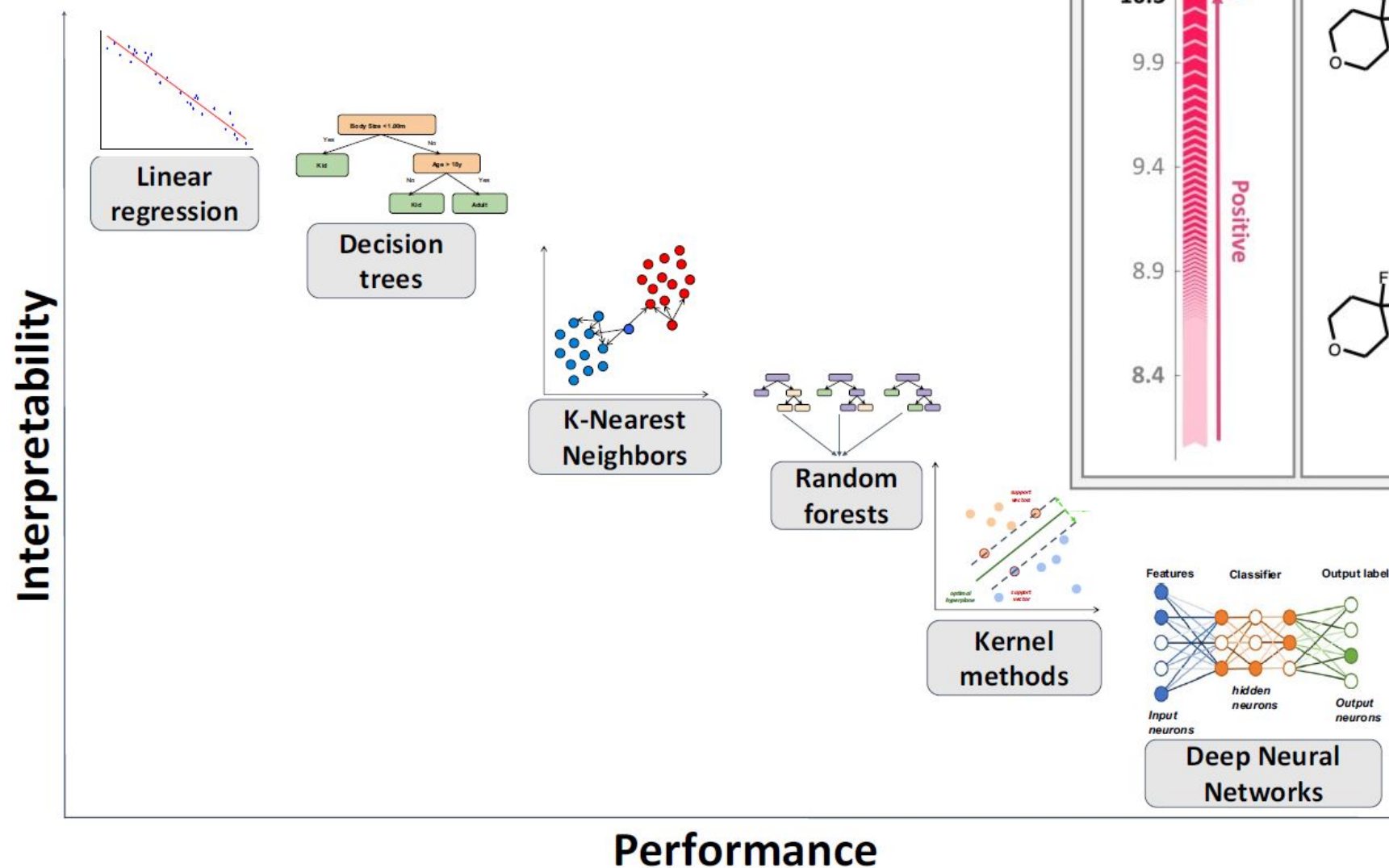
(Left) To assess the generalization ability of a supervised learning algorithm, data are separated into a training subset used for building the model and a test subset used to assess the generalization error (from Badillo *et al.*, 2020) (Right) Temporal validation is especially important for drug discovery,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

# Molecular similarity does not equal biological similarity



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.

# Interpretable and Causal Models will become more important



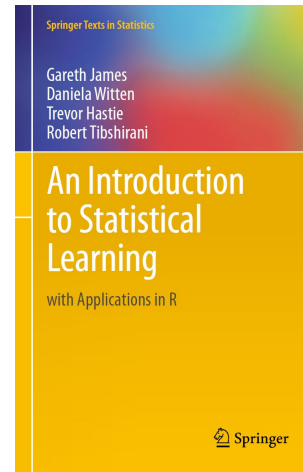
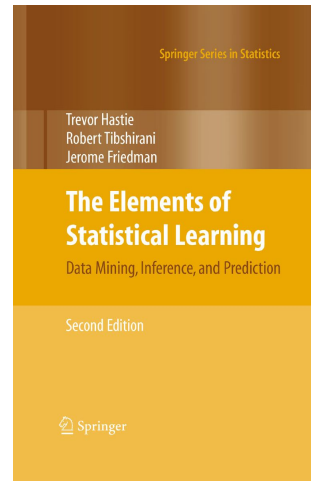
Rodríguez-Pérez, Raquel, and Jürgen Bajorath. "[Interpretation of Machine Learning Models Using Shapley Values: Application to Compound Potency and Multi-Target Activity Predictions.](#)"

Journal of Computer-Aided Molecular Design 34, no. 10 (October 1, 2020): 1013–26..

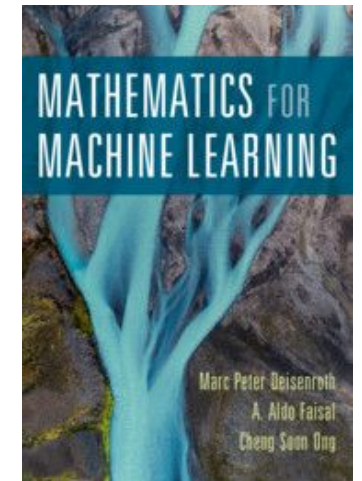


# Resources for learning about machine learning

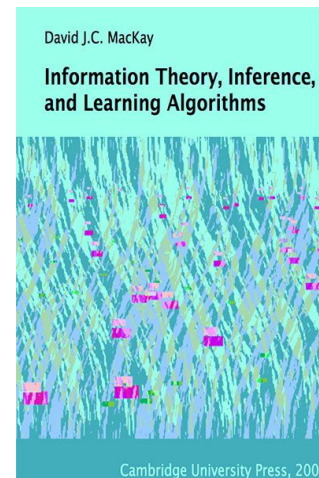
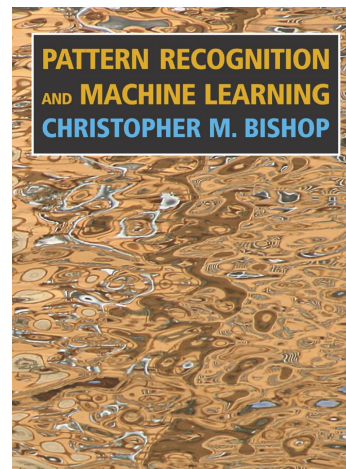
ESL and ISL: From a frequentist view (almost)



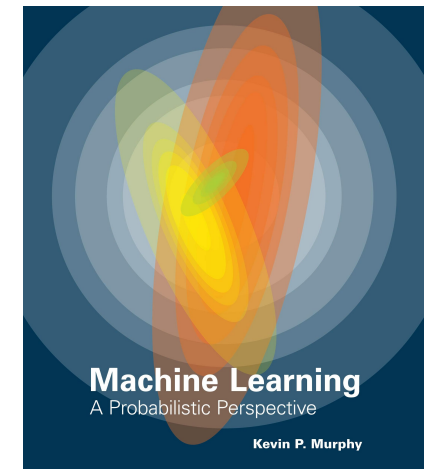
Mathematical foundations



PRML and ITILA: From a Bayesian view



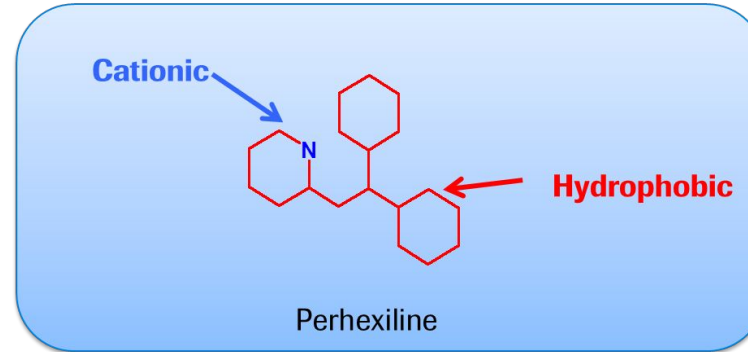
MLaPP: Application oriented, more accessible, and balanced views



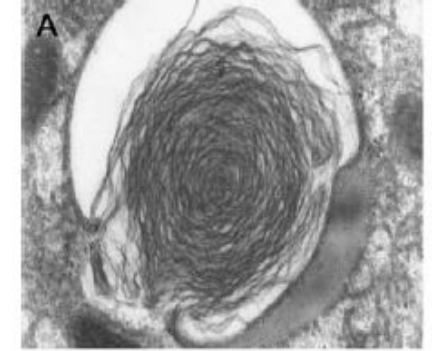


# Drug-induced phospholipidosis is correlated with amphiphilicity

- 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.



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

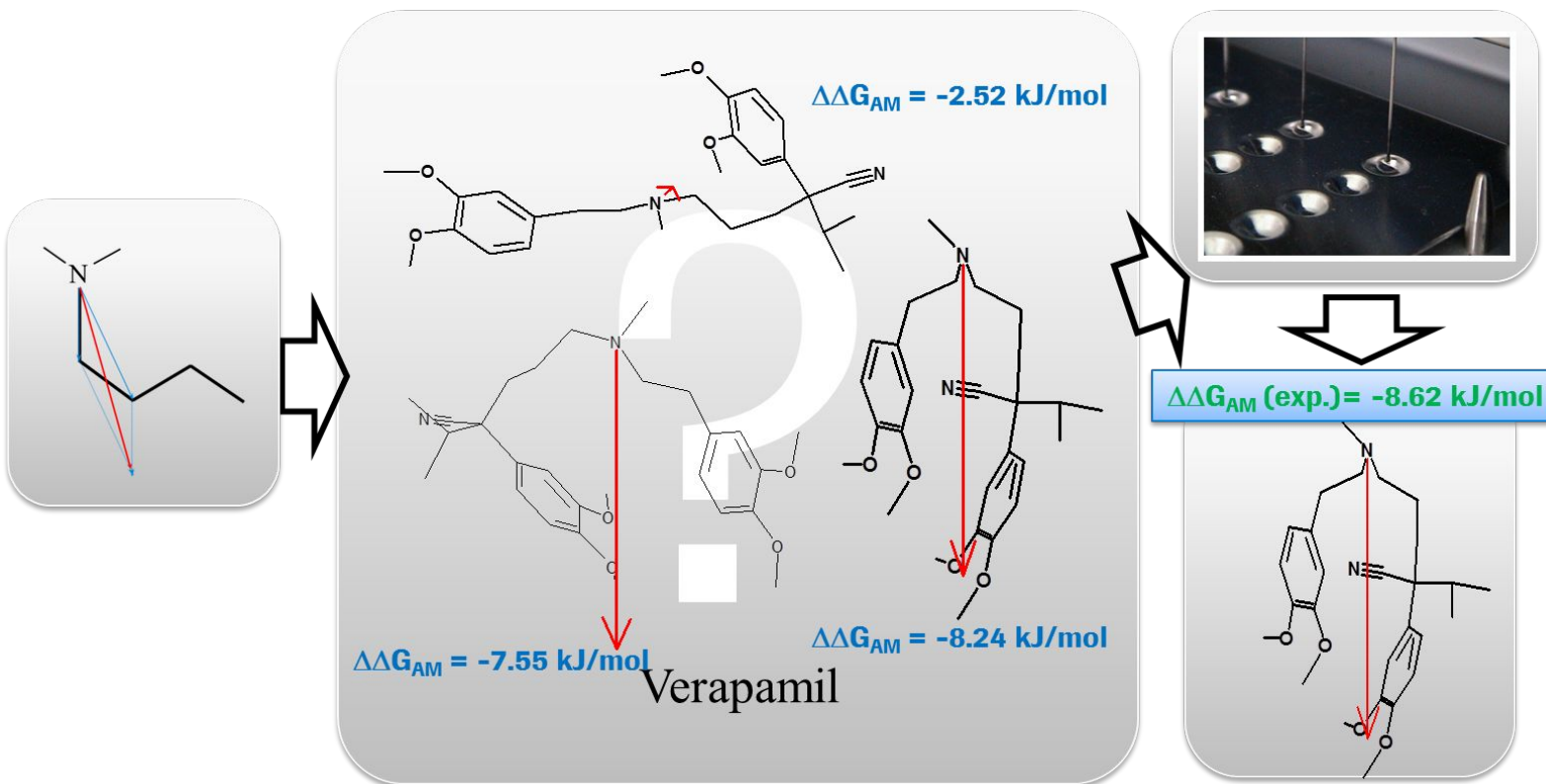
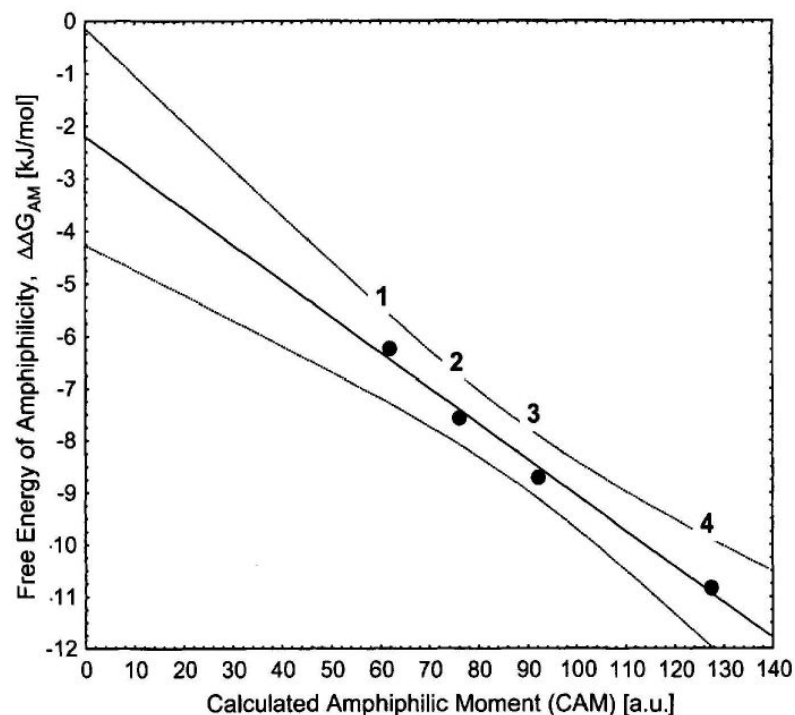
$\vec{\alpha}_i$ : the hydrophobic/hydrophilic contribution of atom/fragment  $i$

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.

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

# In silico prediction of amphiphilicity

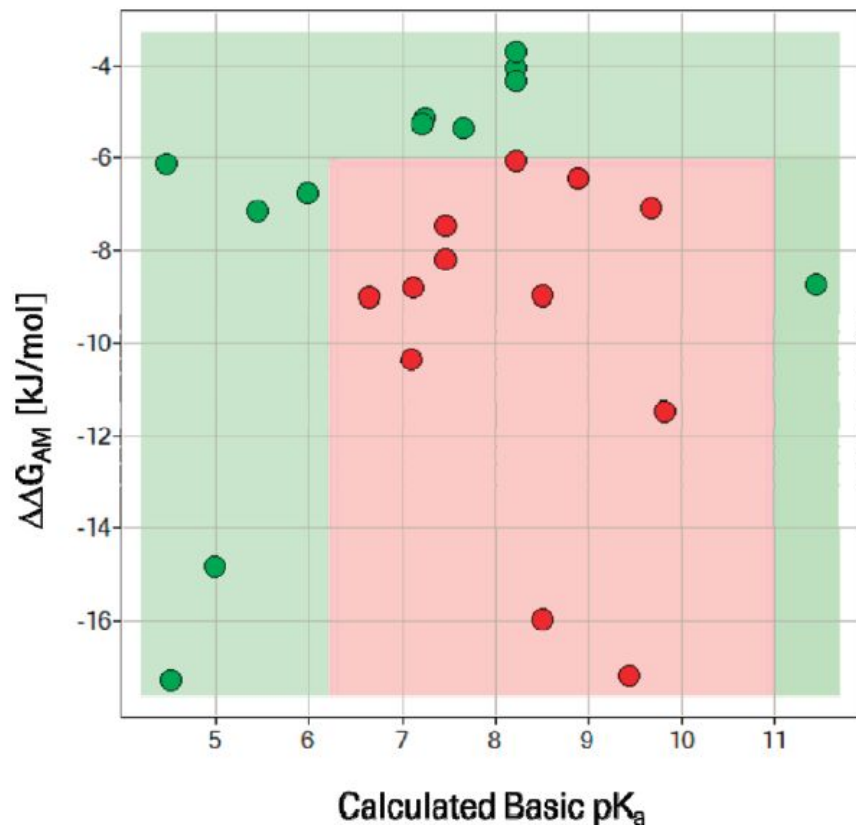
*Development of CAFCA (CAlculated Free energy of amphiphilicity of small Charged Amphiphiles)*



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/<br>in vivo | in silico/<br>in vivo | Exp. PC/<br>in vivo | In silico/<br>in vitro | n=36 |
|----------------------|-----------------------|---------------------|------------------------|------|
| 94%                  | 81%                   | 89%                 | 89%                    |      |

| in vitro/in silico              |                                        |                                        | n=422                     |
|---------------------------------|----------------------------------------|----------------------------------------|---------------------------|
| Accuracy<br>[(TP+TN)/<br>(P+N)] | Sensitivity<br>[True Positive<br>Rate] | Specificity<br>[True Negative<br>Rate] | Precision<br>[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!**

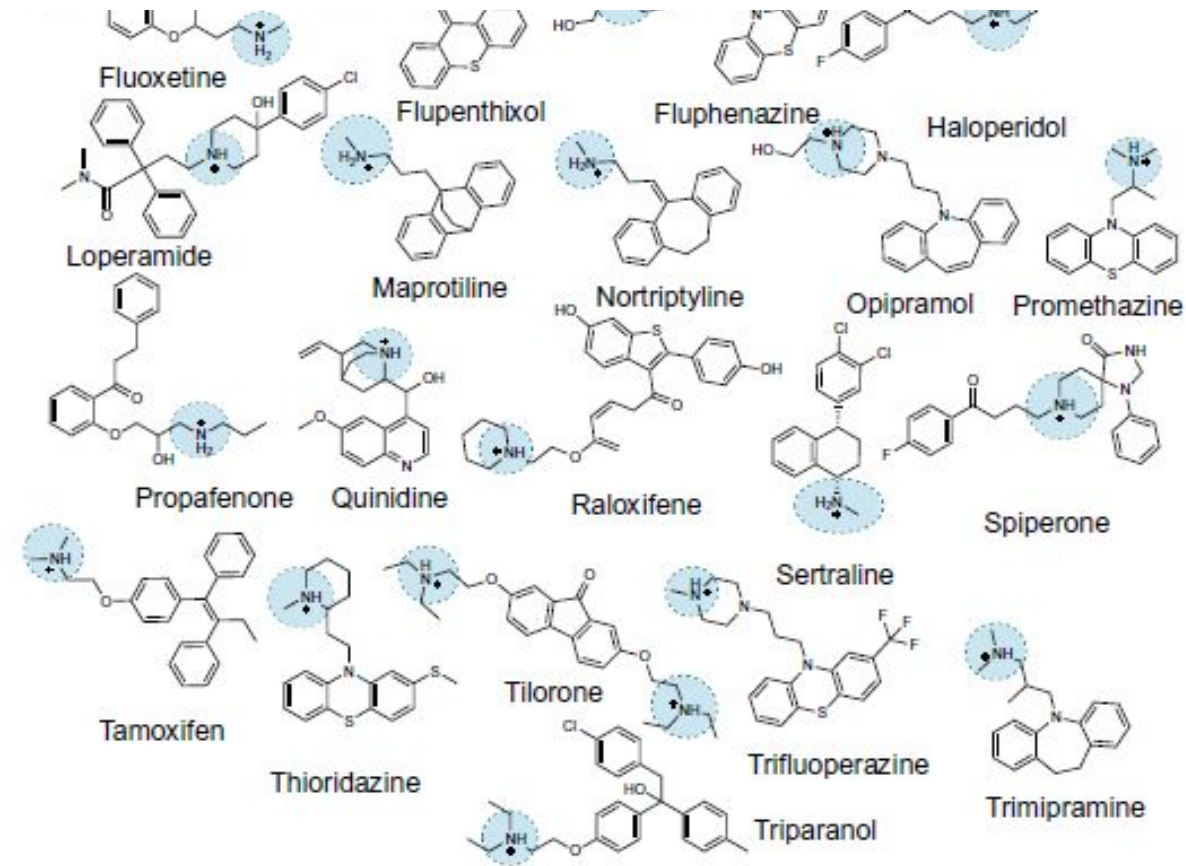
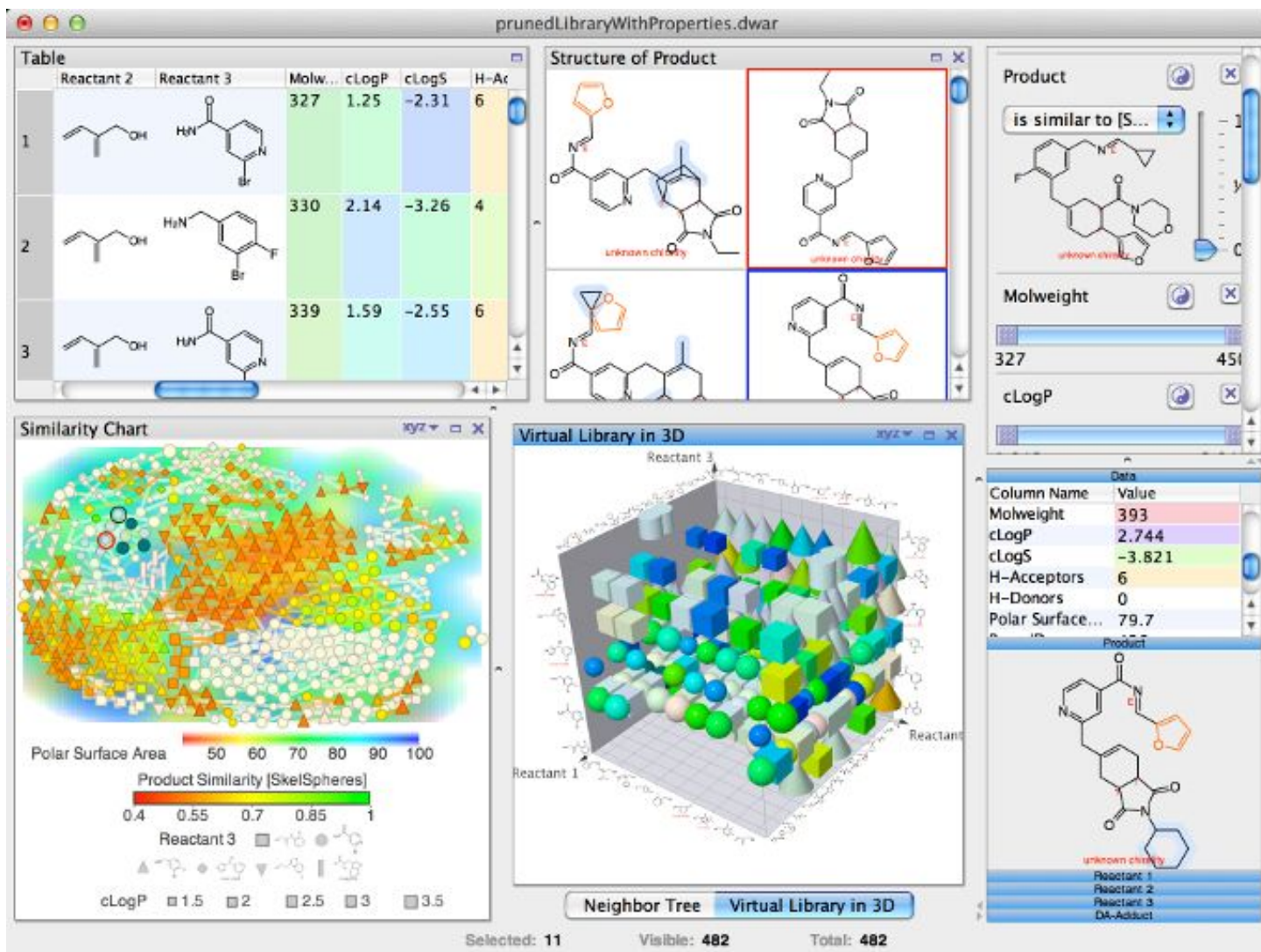


Fig. 1. Representative examples of CADs that are identified in SARS-CoV-2 drug repurposing screens.

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>.



# DataWarrior: an open-source program for data visualization and analysis with chemical intelligence



*DataWarrior* was and still is developed at Actelion/Idorsia Pharmaceuticals Ltd.

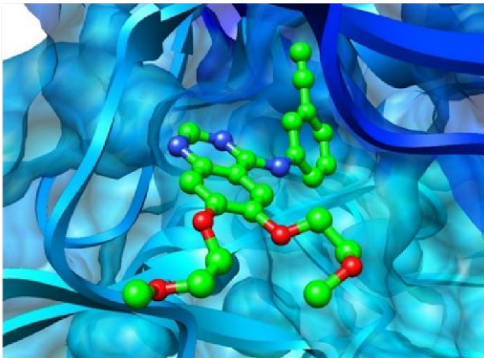
Selected subset of functionalities

- Molecular descriptor calculation
- Similarity calculation
- Compound clustering
- Docking
- ...

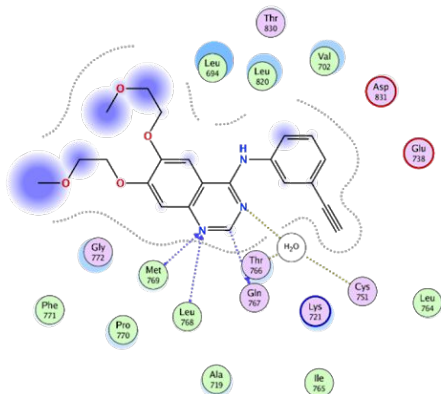
Thomas Sander, Joel Freyss, Modest von Korff, Christian Rufener. *DataWarrior: An Open-Source Program For Chemistry Aware Data Visualization And Analysis*. J Chem Inf Model 2015, 55, 460-473, [doi 10.1021/ci500588j](https://doi.org/10.1021/ci500588j)

# Summary of molecular modelling

## A 3D protein structure-based approaches

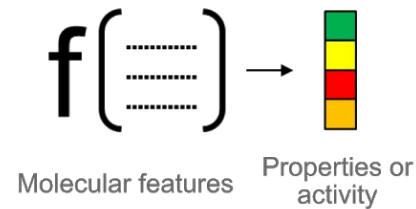
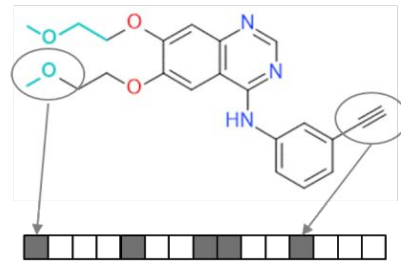


3D model of drug-target complex

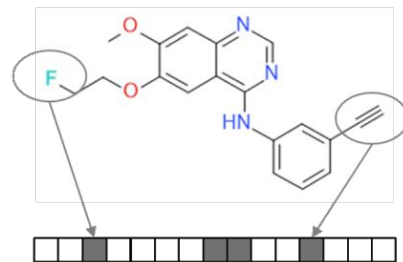


Drug-target interaction map

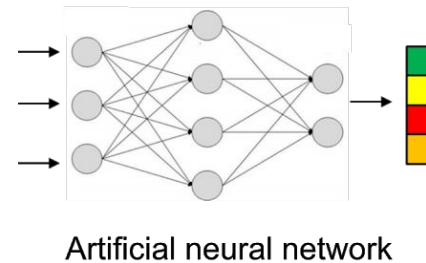
## B Ligand-based approaches



QSAR



Matched molecular pairs and whole-molecule similarity



Artificial neural network

Today we learned ligand-target interaction and **molecular modelling techniques**:

- (A) 3D protein structure-based approaches. An example with docking can be found in the backup slides.
- (B) Ligand-based approaches (similarity search). Another example of amphiphilicity can be found in the backup slides.

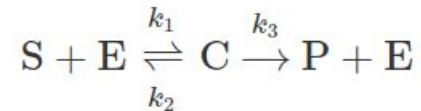
Zhang, Jitao David, Lisa Sach-Peltason, Christian Kramer, Ken Wang, and Martin Ebeling. 2020. “[Multiscale Modelling of Drug Mechanism and Safety](#).” *Drug Discovery Today* 25 (3): 519–34.

# **Why modelling molecules is not enough for drug discovery?**

***The importance of networks***

# Simulation of biological networks with ordinary differential expression: the simplest case

Given the reaction



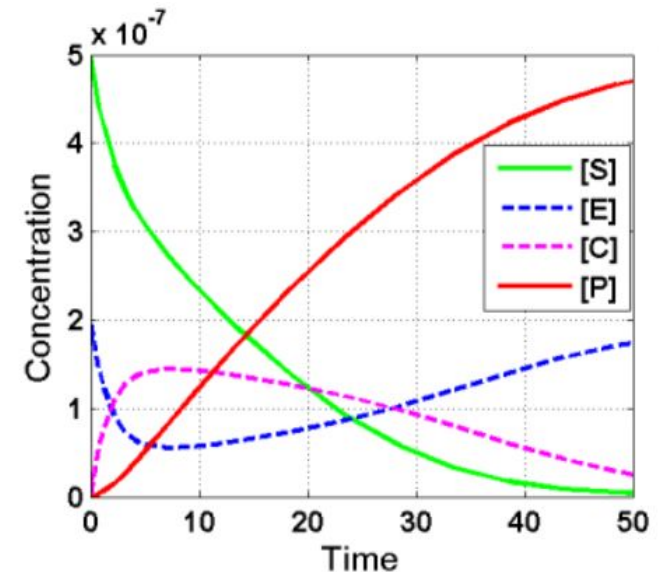
Given the initial values and rate constants

- $S(0) = 5e^{-7}$
- $E(0) = 2e^{-7}$
- $C(0) = P(0) = 0$
- $k_1 = 1e^6$
- $k_2 = 1e^{-4}$
- $k_3 = 0.1$

According to the law of mass action

$$\begin{aligned}\frac{d[S]}{dt} &= -k_1[E][S] + k_2[C], \\ \frac{d[E]}{dt} &= -k_1[E][S] + (k_2 + k_3)[C], \\ \frac{d[C]}{dt} &= k_1[E][S] - (k_2 + k_3)[C], \\ \frac{d[P]}{dt} &= k_3[C],\end{aligned}$$

It is possible to simulate the concentration changes by time *deterministically*.



See [Systems Engineering Wiki \(tue.nl\)](https://www.systems-engineering.wiki/tue.nl) for MATLAB/COPASI codes and *Stochastic Modelling for Systems Biology* by Darren J. Wilkinson



# Chemical Master Equations (CME): a particle model of chemical reaction

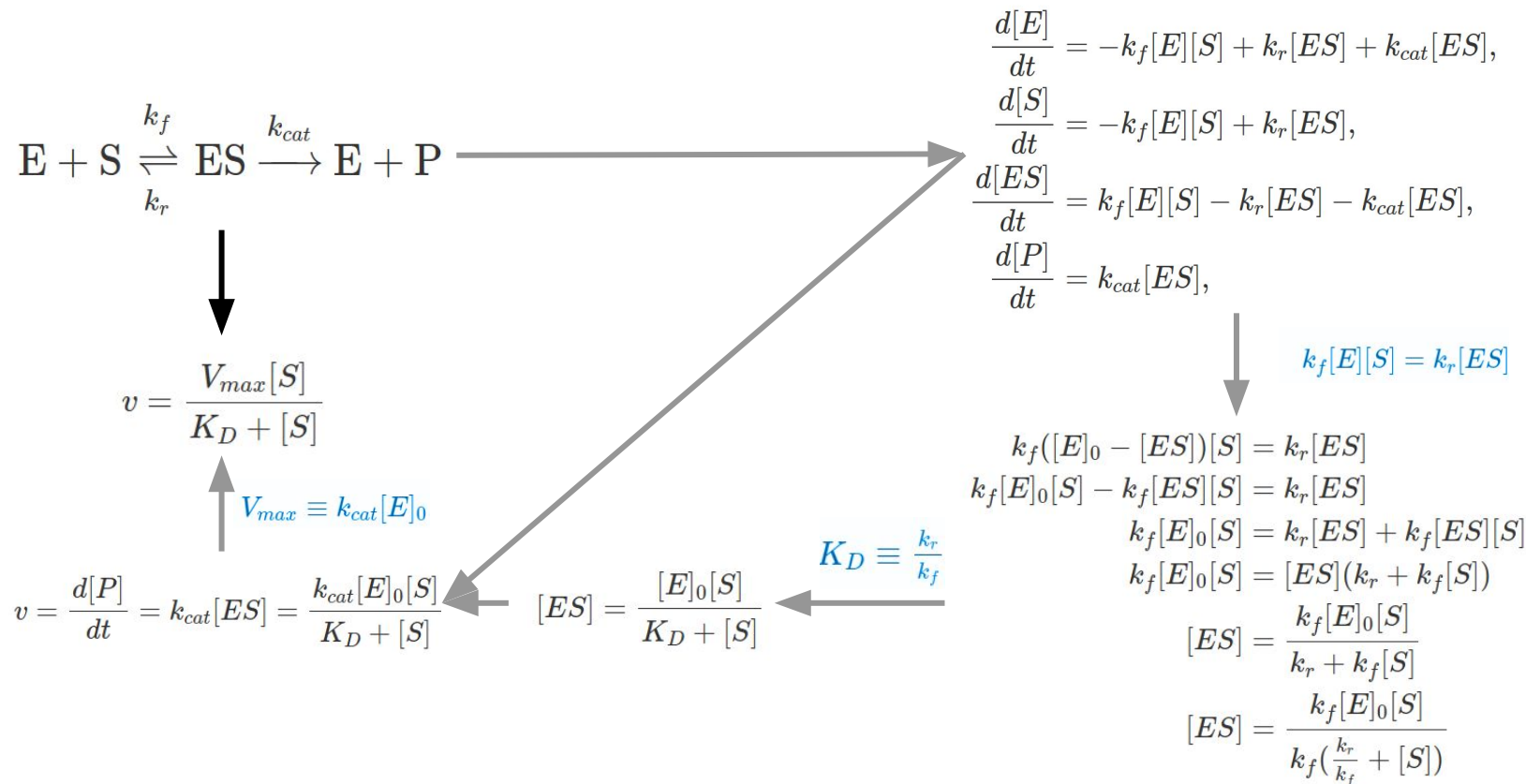
Given the reaction  $A + B \xrightleftharpoons[k_2]{k_1} C + D$  and the initial condition  $X(0) = \begin{bmatrix} K \\ K \\ 0 \\ 0 \end{bmatrix}$  ( $K$  molecules of species A and of species B respectively)

The state vector  $X(t)$  can take at any time point *one* of the values  $\begin{bmatrix} K \\ K \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K-1 \\ K-1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} K-2 \\ K-2 \\ 2 \\ 2 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ K \\ K \end{bmatrix},$

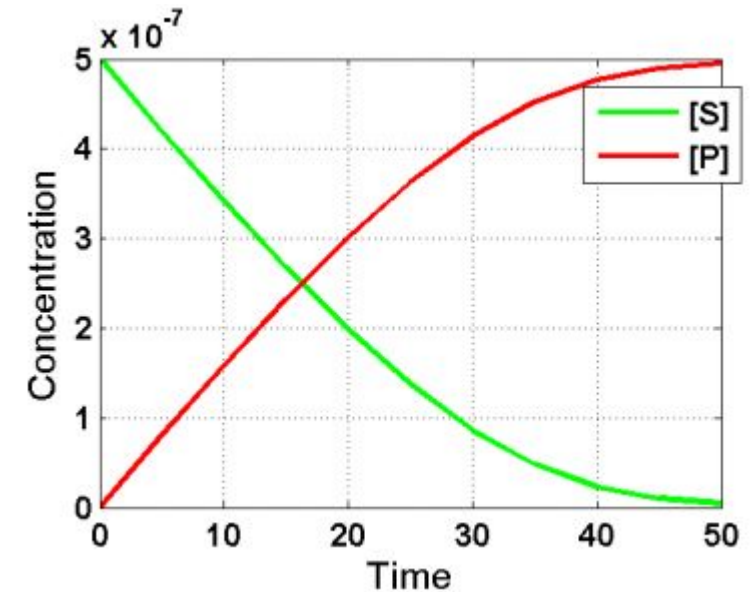
Theoretically we can build an ODE system with  $K+1$  equations to model *every state of the reaction*, down to every particle. In reality, the dimension is so high so that a simulation is not feasible.

**CME is a set of ODEs, with each ODE representing one possible state of the system. Solution of the  $k$ th equation at time  $t$  is a real number giving the probability of system being in that particular state at that time.**

# Reaction Rate Equations (RRE): a compartment model



**RRE simulation of the Michaelis-Menten model**



Source: [Systems Engineering Wiki \(tue.nl\)](https://www.tue.nl/systems-engineering/wiki/)

**RRE is a set of ODEs, with each ODE representing one chemical species. Solution of the  $j$ th equation at time  $t$  is a real number representing the concentration of species  $j$  at time  $t$ .**

# The Gillespie's algorithm and the chemical Langevin equation allow stochastic simulation of biological networks

- The *stochastic simulation algorithm* (exact SSA), also called *Gillespie's algorithm*, allows stochastic simulation of a reaction.
- It is performed in four steps
  - **Initialize** the system with initial conditions
  - Given a state at time  $t$ , we can define a probability  $p$  that reaction  $j$  takes place in the time interval  $[t+\tau, t+\tau+d\tau)$ . It is the product of two density functions of two random variables: the probability of reaction  $j$  happens (proportional to the number of substrate molecules), multiplied by the time until next reaction, which is exponentially distributed. This is known as the **Monte Carlo** step.
  - Let the randomly selected reaction happen and **update** the time.
  - **Iterate** until substrates are exhausted or simulation time is over.
- Further computation tricks such as 'tau-leaping' by lumping together reactions are possible. The chemical Langevin equation (CLE) replaces further accelerates stochastic simulation by approximating the Poisson distribution with the normal distribution.

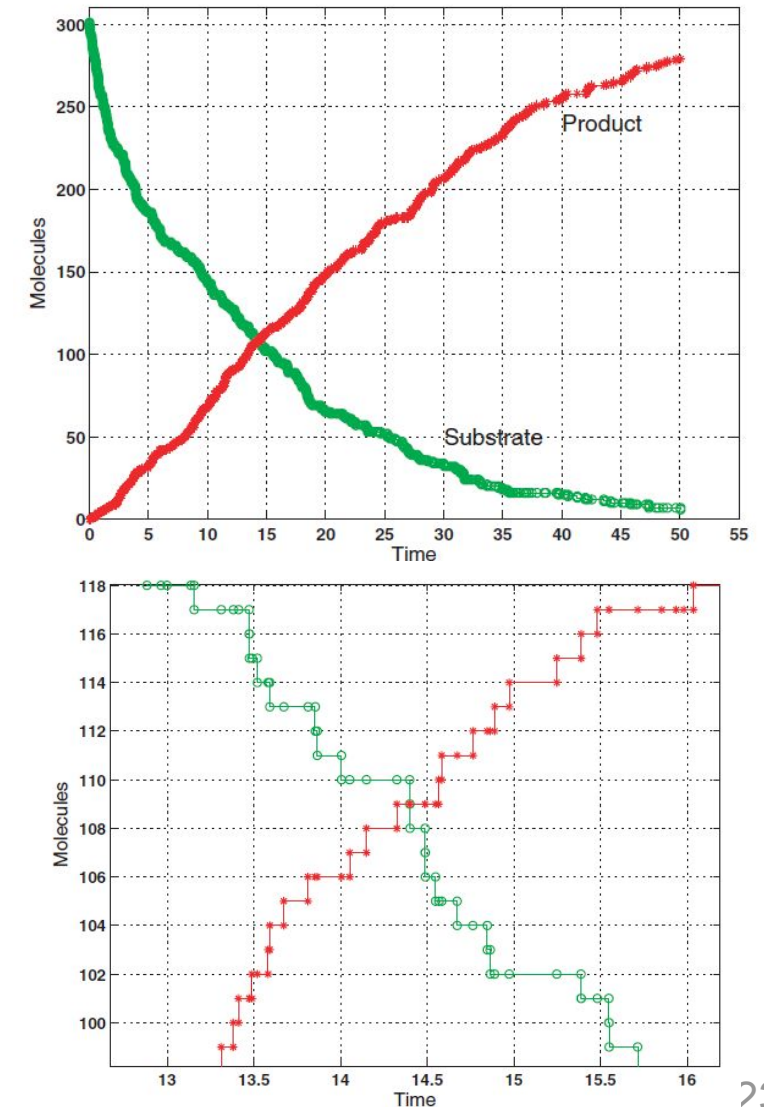
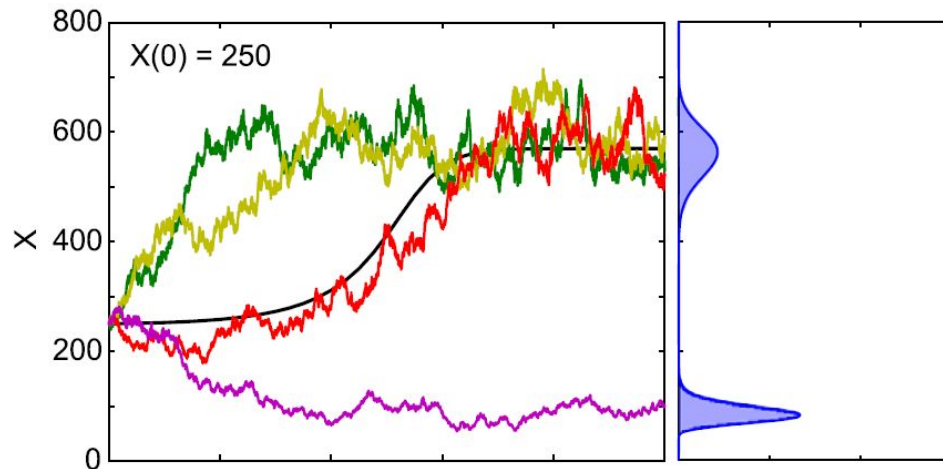


Figure source and further reading: Higham, Desmond J. 2008. "Modeling and Simulating Chemical Reactions." *SIAM Review* 50 (2): 347–68. <https://doi.org/10.1137/060666457>.

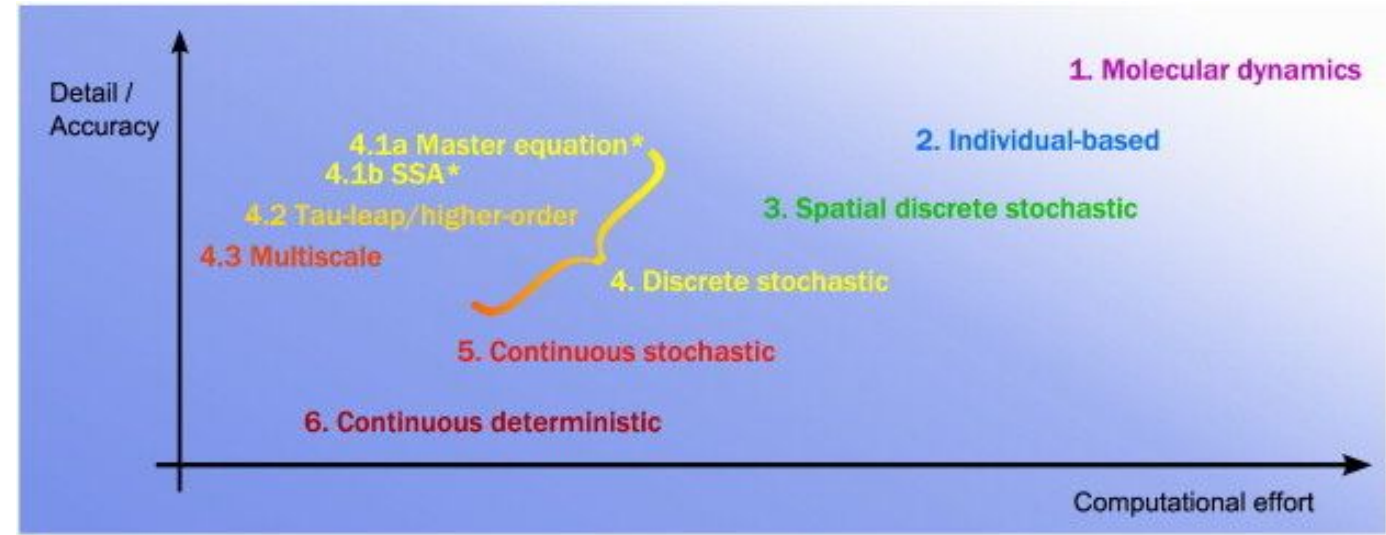
# Why stochastic modelling?



- Stochastic modelling can reveal individual trajectories that are otherwise ‘averaged’ by ODE models.
- Small systems and single-molecule studies show stochastic behaviour.
- It is possible to consider both extrinsic and intrinsic factors and take them into the model.

Székely and Burrage. 2014. “[Stochastic Simulation in Systems Biology](#).” *Computational and Structural Biotechnology Journal* 12 (20–21): 14–25.

Also see *Stochastic Modelling for Systems Biology* by Darren J. Wilkinson.



Advantages and disadvantages of several modelling/simulation methods.

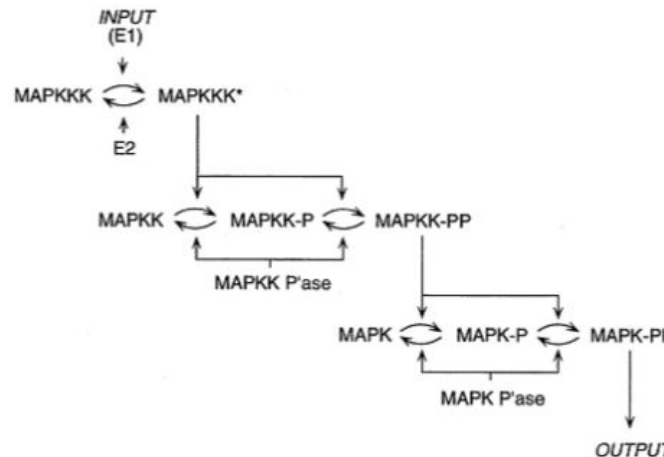
| Simulation method | Cat. | Advantages                                              | Disadvantages                                                                | References            | Software                                                        |
|-------------------|------|---------------------------------------------------------|------------------------------------------------------------------------------|-----------------------|-----------------------------------------------------------------|
| Master equation   | 4    | Exact                                                   | Very computationally intensive                                               | [85,143]              |                                                                 |
| SSA               | 4    | Statistically exact                                     | Very computationally intensive                                               | [82,109]              | COPASI [144]<br>StochKit [145]<br>STOCKS [146]<br>BioNetS [147] |
| Tau-leap          | 4    | Relatively fast                                         | Approximate; too slow for large systems or frequent/multiscale reactions     | [83,113,118]          | StochKit [145]                                                  |
| Higher-order      | 4    | Relatively fast; accurate                               | Approximate; too slow for large systems or frequent/multiscale reactions     | [83,121,122,124,125]  |                                                                 |
| Multiscale/hybrid | 4    | Fast; good for systems with disparate reaction scales   | Approximate; problems with coupling different scales                         | [131,132,137,139,148] | COPASI [144]<br>BioNetS [147]                                   |
| Brownian dynamics | 2    | Tracks individual molecules                             | Slow; molecule size must be artificially added                               | [149,150]             | Smoldyn [149,151]<br>MCell [152]                                |
| Compartment-based | 3    | Accounts for diffusion between homogeneous compartments | Slow; compartment size must be set manually; each compartment is homogeneous | [150,153,154]         | MesoRD [153]<br>URDME [155]                                     |
| SDE               | 5    | Fast                                                    | Continuous; Gaussian noise                                                   | [76]                  | BioNetS [147]                                                   |
| PDE (R-D)         | 6    | Very fast; spatial                                      | Continuous; no noise                                                         | [156]                 |                                                                 |
| ODE               | 6    | Very fast                                               | Continuous; no noise                                                         | [157]                 |                                                                 |

Cat. represents Category from Fig. 2. Abbreviations: SSA, stochastic simulation algorithm; SDE, stochastic differential equation; PDE (R-D), partial differential equation (classical reaction-diffusion equations); ODE, ordinary differential equation.

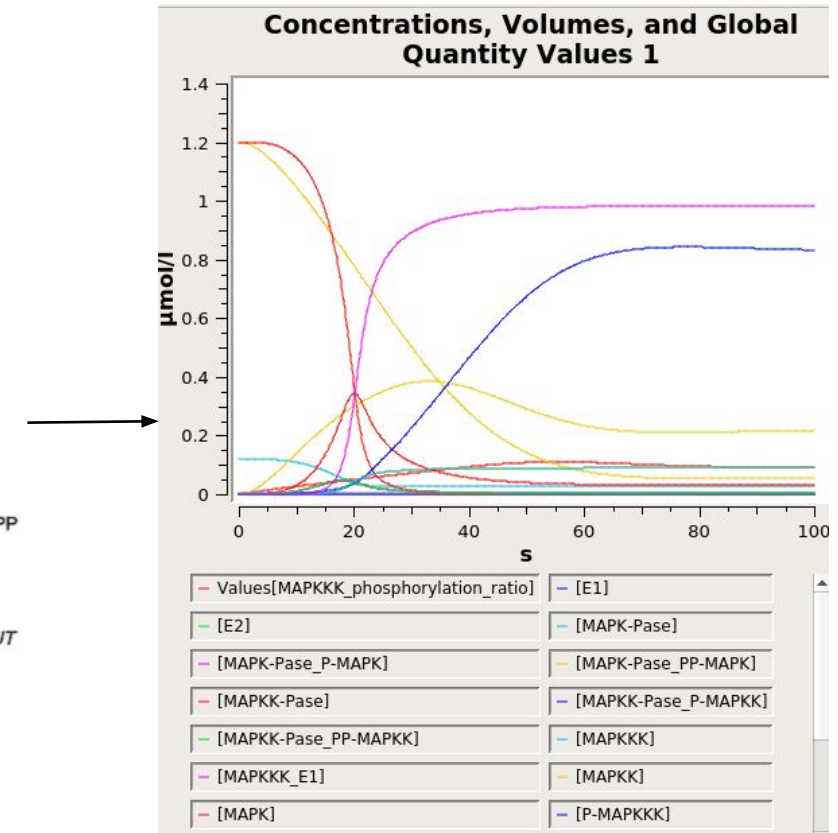


# Biochemical system simulator COPASI

- Freely available at <http://COPASI.org/>
- COPASI supports two types of simulations: (1) **ordinary differential equation (ODE)** based simulation, (2) **stochastic kinetic simulation**, among others using the [stochastic Runge–Kutta method](#) (RI5) and [Gillespie's algorithm](#)
  - Resources to learn more about stochastic modelling: [MIT OpenCourseWare](#) by Jeff Gore, and [Stochastic Processes: An Introduction, Third Edition](#) by Jones and Smith
- Tutorials also available on [the website of European Bioinformatics Institute \(EBI\)](#)
- The mathematical concept and software tools are important for detailed analysis of enzymatic reactions, especially in the presence of drugs and/or disease-relevant mutation



Huang and Ferrell, PNAS, 2006

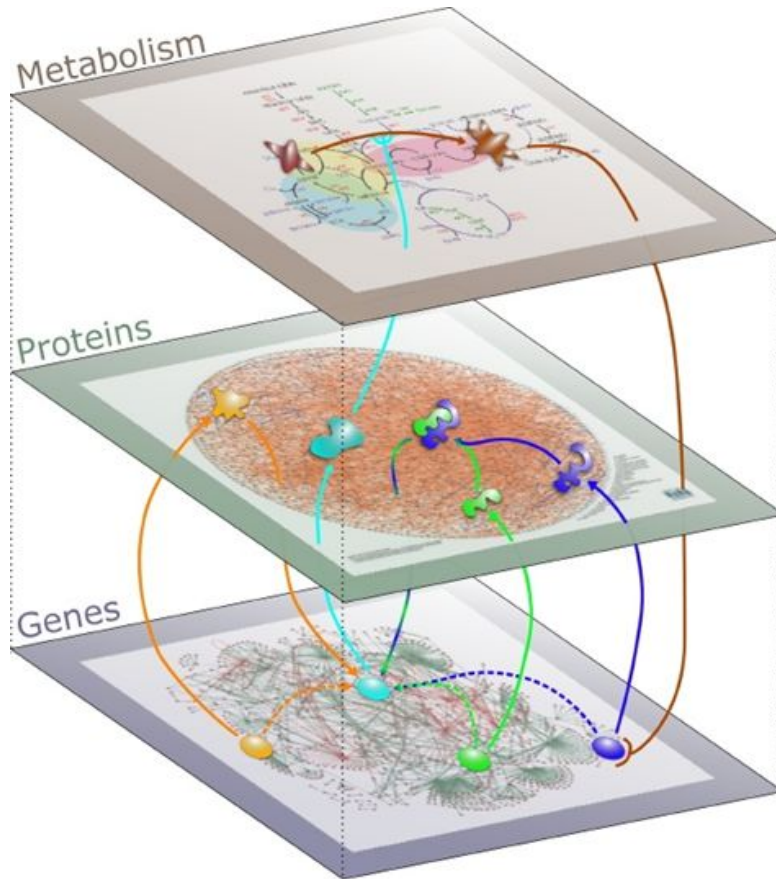


ODE-based simulation of dynamics

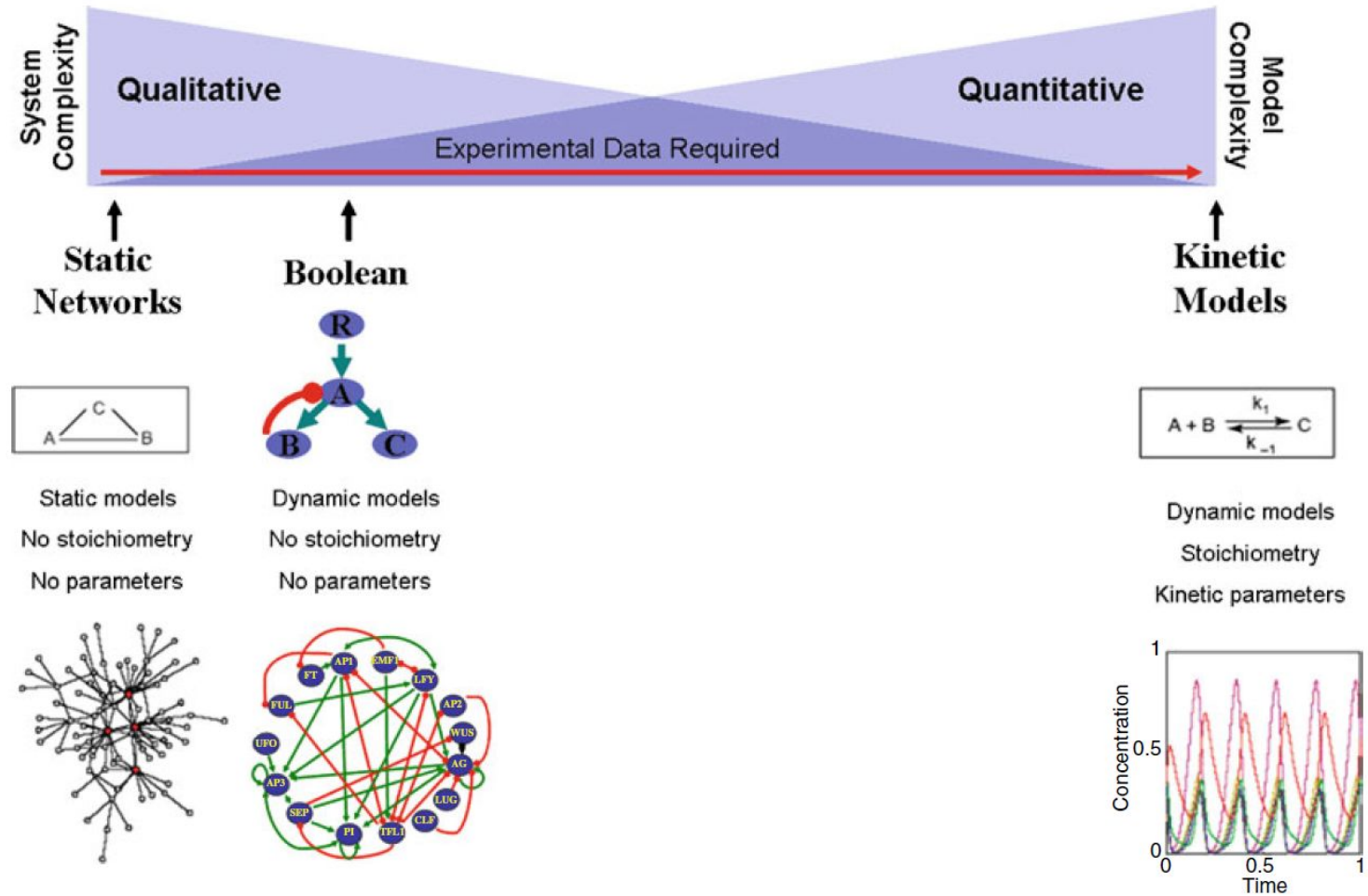
# Summary

- **QSAR and machine learning models in drug discovery.**
- **Machine learning should be guided by chemical and biological models to improve human understanding.**
- **ODE-based compartment models and stochastic models can be used to model small to moderate biochemical networks.**
- **The network nature of biology requires models beyond the molecular level.**

# Modelling biological networks



Stéphane CHÉDIN & Jean LABARRE, [www-dsv.cefa.fr](http://www-dsv.cefa.fr)



Garg, Abhishek, Kartik Mohanram, Giovanni De Micheli, and Ioannis Xenarios. 2012. "[Implicit Methods for Qualitative Modeling of Gene Regulatory Networks](#)." In *Gene Regulatory Networks: Methods and Protocols*, edited by Bart Deplancke and Nele Gheldof, 397–443. Methods in Molecular Biology. Totowa, NJ: Humana Press.