

# Хемоинформатика – основные понятия и области применения

**Alexandre Varnek**

*Laboratory of Chemoinformatics, University of Strasbourg*



## Laboratory of Chemoinformatics

**Master in Chemoinformatics (2001)**

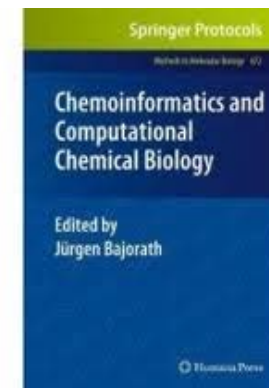
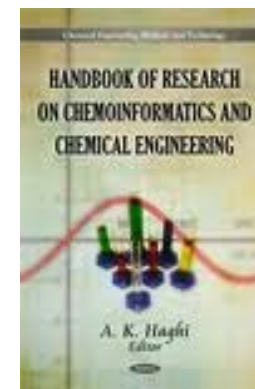
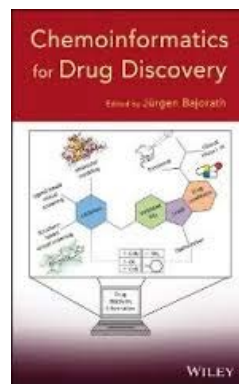
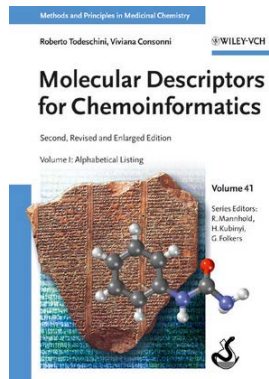
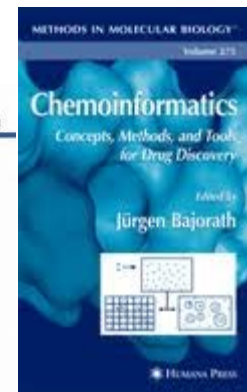
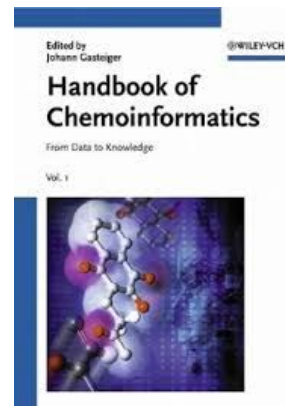
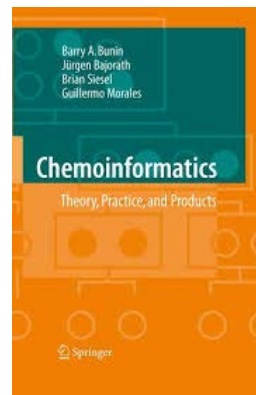
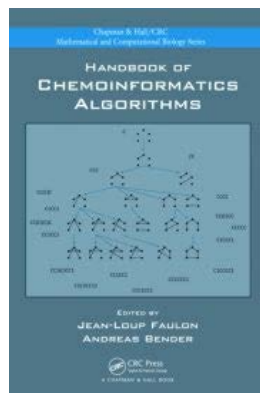
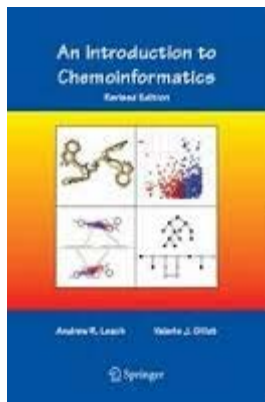
**Master “In Silico Drug Design” (2010)**

Strasbourg - Paris-Diderot - Milan

**Double diploma Strasbourg / Kazan (2013)**



# Selected books in chemoinformatics



# РОССИЙСКИЙ НАУЧНЫЙ ФОНД

ПОДДЕРЖКА И РАЗВИТИЕ

НОВОСТИ

О ФОНДЕ

ДОКУМЕНТЫ

КОНКУРСЫ

КОНТАКТЫ

## Классификатор Фонда

КОД НАИМЕНОВАНИЕ

03-700 Многомасштабное компьютерное моделирование структуры и свойств сложных химических систем и материалов

03-705 Хемоинформатика



**грант РФФ 2014**

# Chemoinformatics:

new discipline combining several „old“ fields

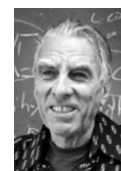
- Chemical databases
- Structure-Activity modeling (QSAR)
- Structure-based drug design
- Computer-aided synthesis design



Michael Lynch



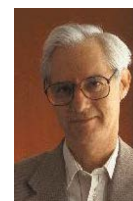
Peter Willett



Corwin Hansch



Johann Gasteiger



Irwin D. Kuntz



Hans-Joachim Böhm



Elias Corey



Ivar Ugi

# OUTLOOK

- **Fundamentals of chemoinformatics**
- **Structure-Property modeling: case studies**
  - **Anti-thrombotics**
  - **Tautomeric equilibria**
  - **Ionic Liquids**
  - **Chemical reactions**
- **Materials design**

# Chemoinformatics: fundamentals

# Review

DOI: 10.1002/minf.201000100

## Cheminformatics as a Theoretical Chemistry Discipline

Alexandre Varnek<sup>\*[a]</sup> and Igor I. Baskin<sup>[b]</sup>

*Mol. Inf.* 2011, 30, 20 – 32

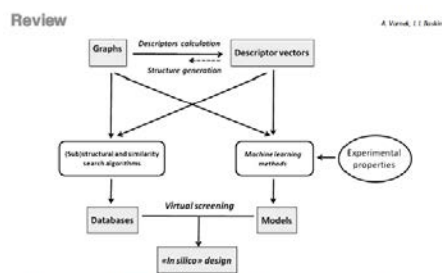


Figure 1. Cheminformatics: from objects to major applications. Note: for each Cheminformatics object (graph, descriptor vector) in the input or in the feature space there exist a associated machine-learning approaches: graph based, vector-based or kernel-based method, respectively.

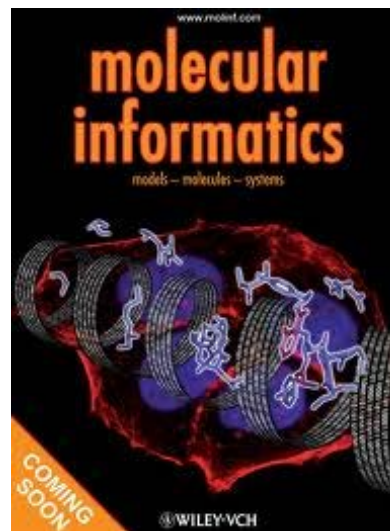
of the kind: molecular model and potential energy equations. Force field methods can be applied to rather large systems containing many thousands of atoms (proteins, solutions, etc.). Cheminformatics considers a molecule as a graph or an ensemble of descriptors generated from this graph. A set of molecules forms a chemical space for which

the relationships between the objects themselves, on one hand, and between their chemical structures and related properties, on the other hand, are established using two main mathematical approaches: graph theory and statistical learning. Due to the rapidity of such calculations, these structure-property relationships can be applied to fast

Alexandre Varnek got his PhD in physical chemistry from the Institute of Inorganic and General Chemistry of the Russian Academy of Sciences, Moscow, in 1989. In 1990, he was a Research Professor in theoretical chemistry at the Moscow Mendeleev University of Chemical Technology. In 1995, Alexandre joined the University of Science and Technology, where he holds the position of a Professor in theoretical chemistry, head of the laboratory of cheminformatics and the director of



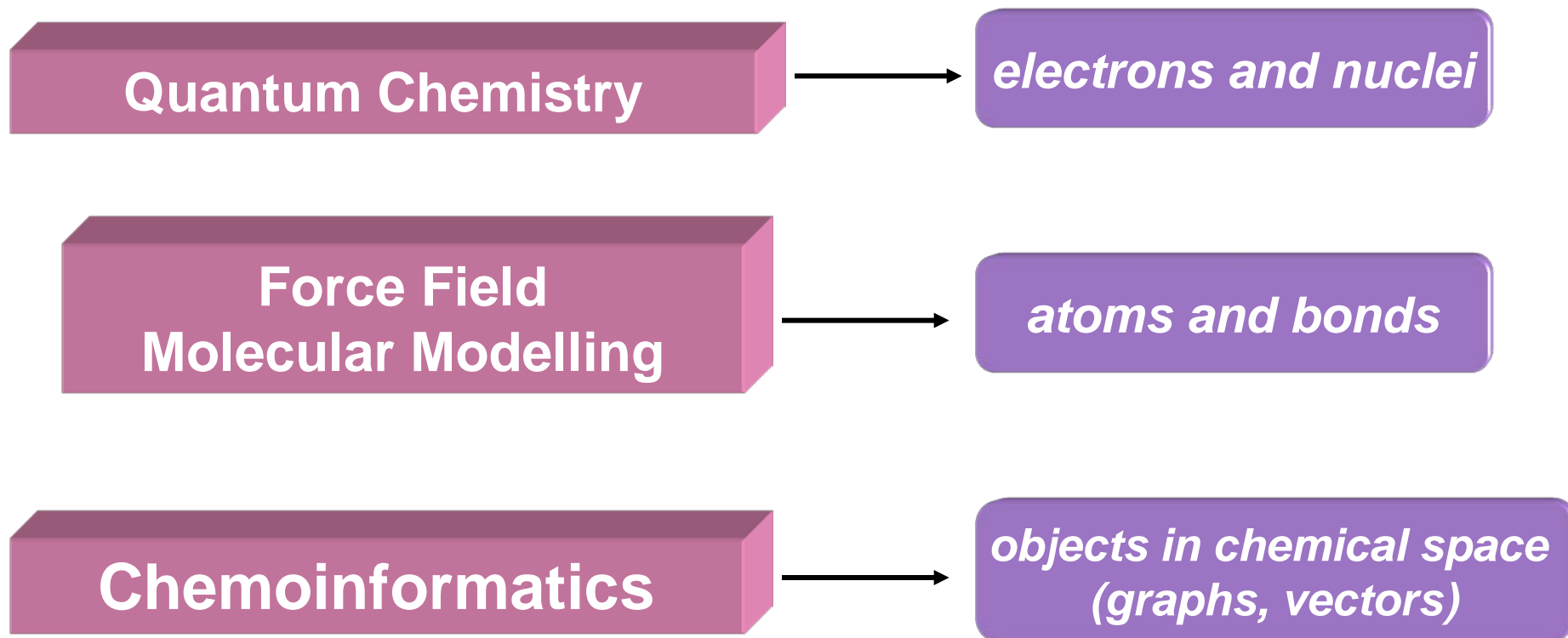
Igor Baskin received his PhD in organic chemistry (1990) and habilitation in mathematical and quantum chemistry (2000) from Lomonosov Moscow State University, Russia. After holding several positions at the Semenov Institute of Chemical Physics and Zelinsky Institute of Organic Chemistry of the Russian Academy of Sciences, Moscow, he joined in 2001 the Chemistry Department of Lomonosov Moscow State University where since 2003 he holds the position of a Leading Scientist. He



**Cheminformatics is defined as individual discipline characterized by its own molecular model, basic concepts, major applications and learning approach**



# *Molecular Model*



# *Molecular Model*

Quantum Chemistry



*electrons and nuclei*

Force Field  
Molecular Modelling



*atoms and bonds*

Chemoinformatics



- *molecular graph*
- *descriptor vector*

# *Basic mathematical approaches*

Quantum Chemistry



*Schrödinger equation,  
HF, DFT, ...*

Force Field  
Molecular Modelling



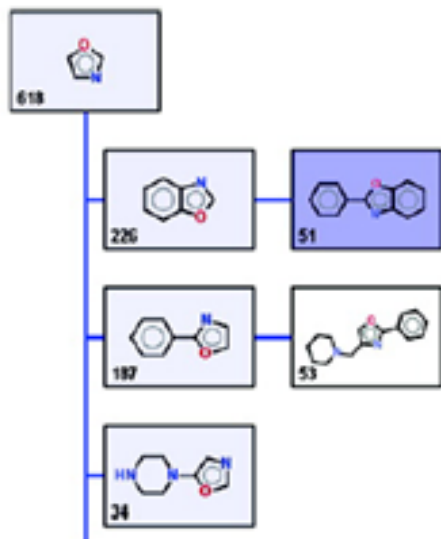
*Classical mechanics  
Statistical mechanics*

Chemoinformatics

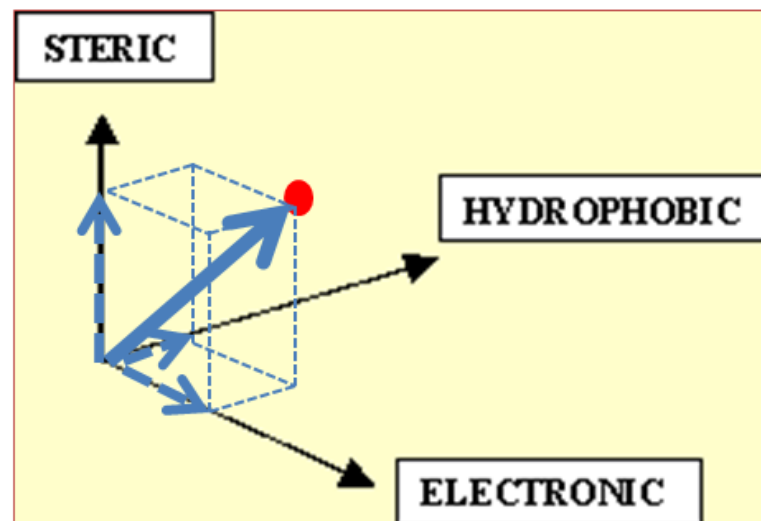


*- Graph theory,  
- Statistical Learning*

# Chemical Space paradigm



graphs-based



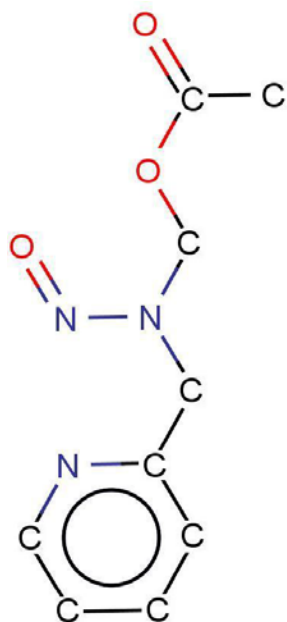
descriptors -based

**SPACE = objects + metric**

# Molecular Descriptors :

ensemble of topological, electronic, geometry parameters calculated directly from molecular structure

Molecular graph



- Topological indices,
- Atomic charges,
- Inductive descriptors,
- Substructural fragments,
- Molecular volume and surface, ...

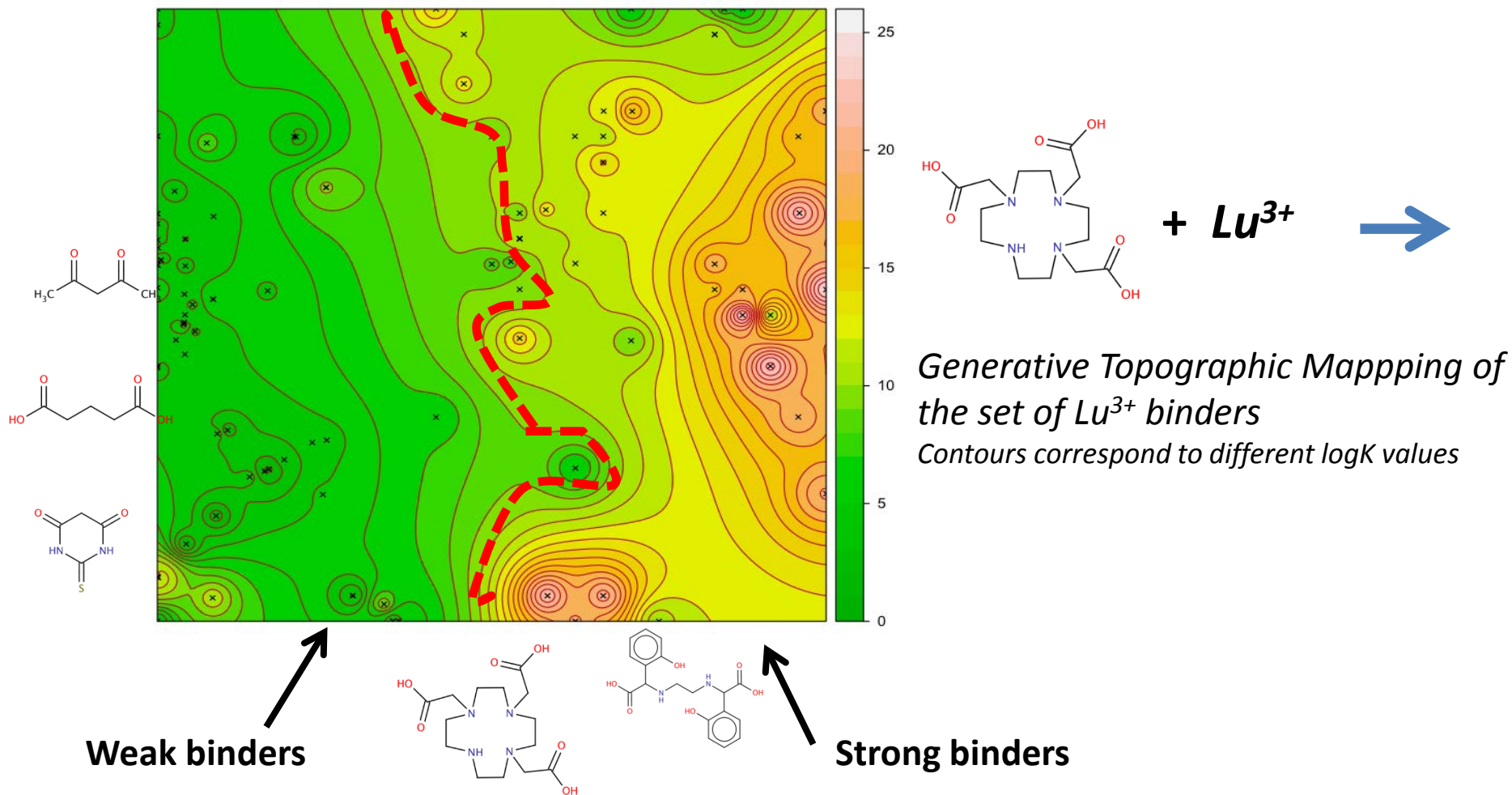
Descriptor vector

Descriptors
$D_1$
$D_2$
...
$D_i$
...



> 5000 types of descriptors are reported

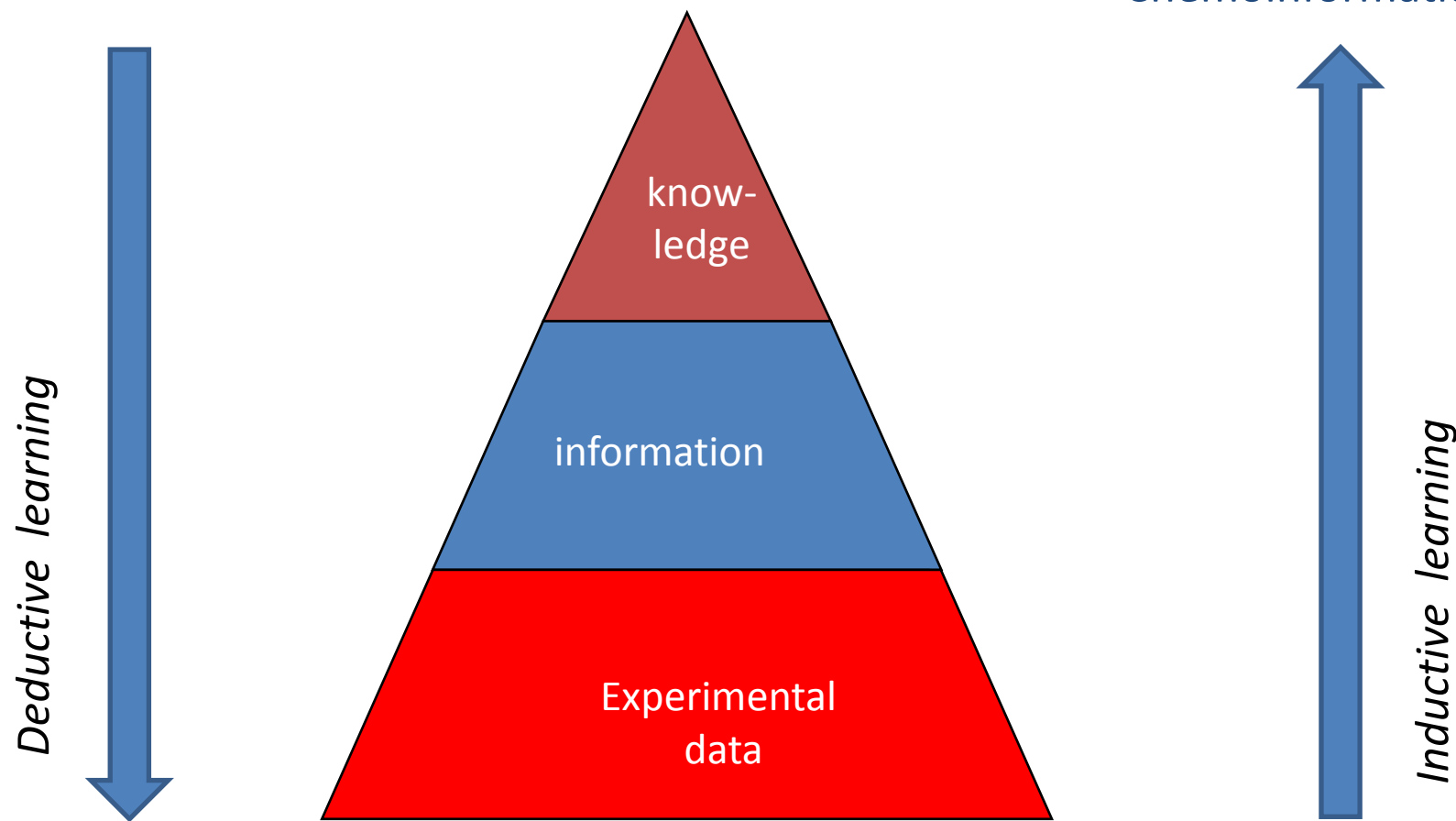
# Chemical space visualization



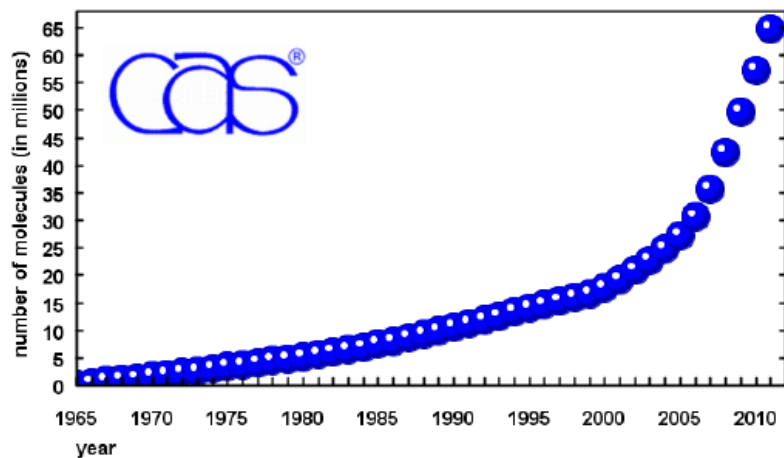
# *Chemoinformatics*: Learning from Data

Quantum  
Mechanics

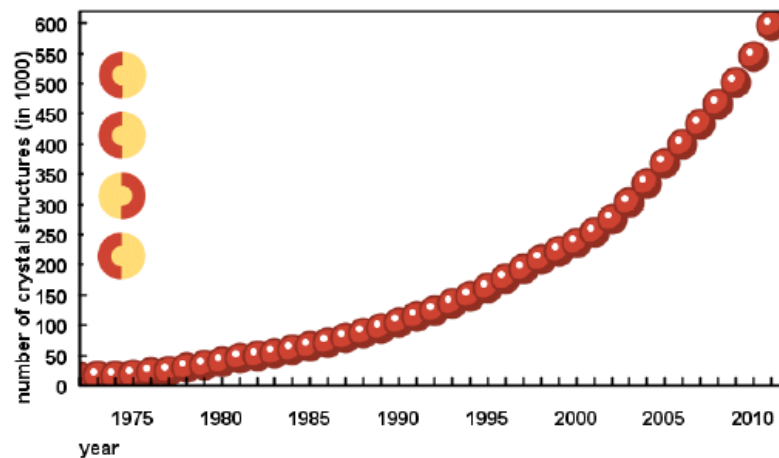
Chemoinformatics



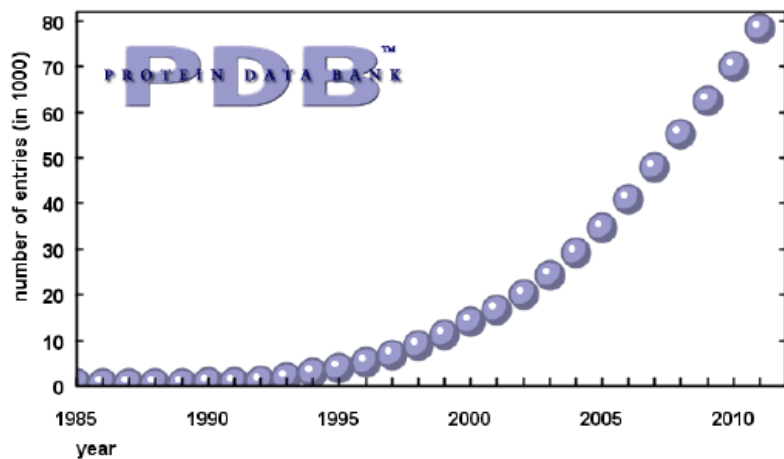
# Data Explosion in Chemistry



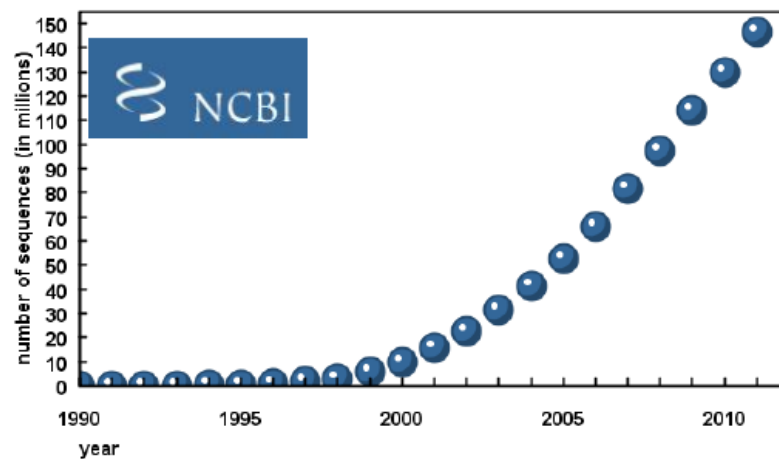
CAS – 65 million molecules



CCDC – 600'000 structures



PDB – 78'000 proteins



GenBank – 145 million sequences



*"Good reactions"™*

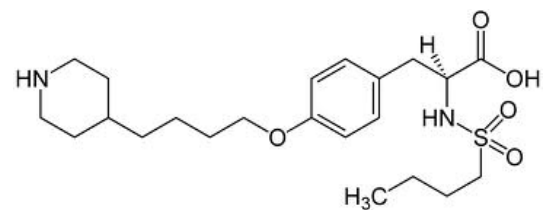
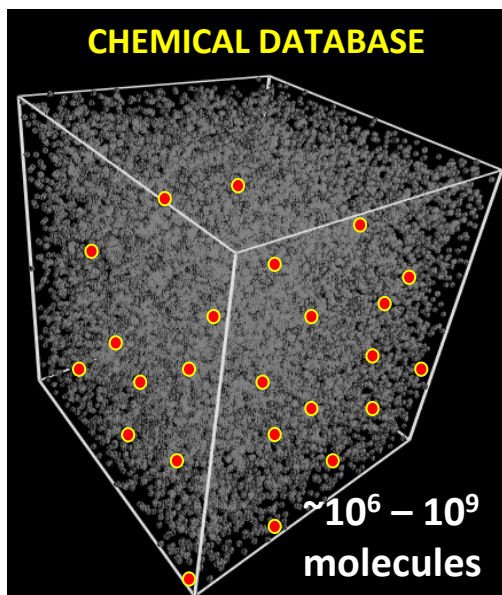


## H-Cube<sup>®</sup> Continuous-flow Hydrogenation Reactor

A revolutionary bench-top standalone hydrogenation reactor combined with automated liquid handler and laboratory automation software

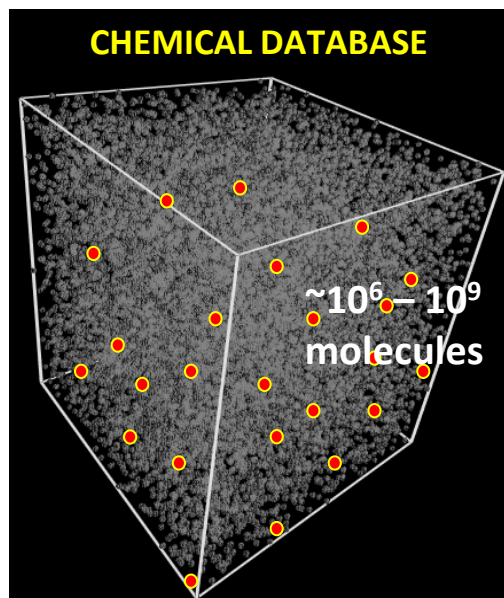
**Up to 100 reactions a day  
≈ 30.000 reactions a year**

# Virtual screening : finding the needle in the haystack

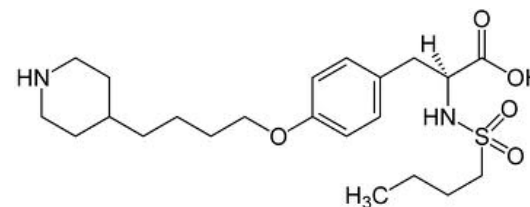


# Chemoinformatics:

## pattern recognition in chemistry

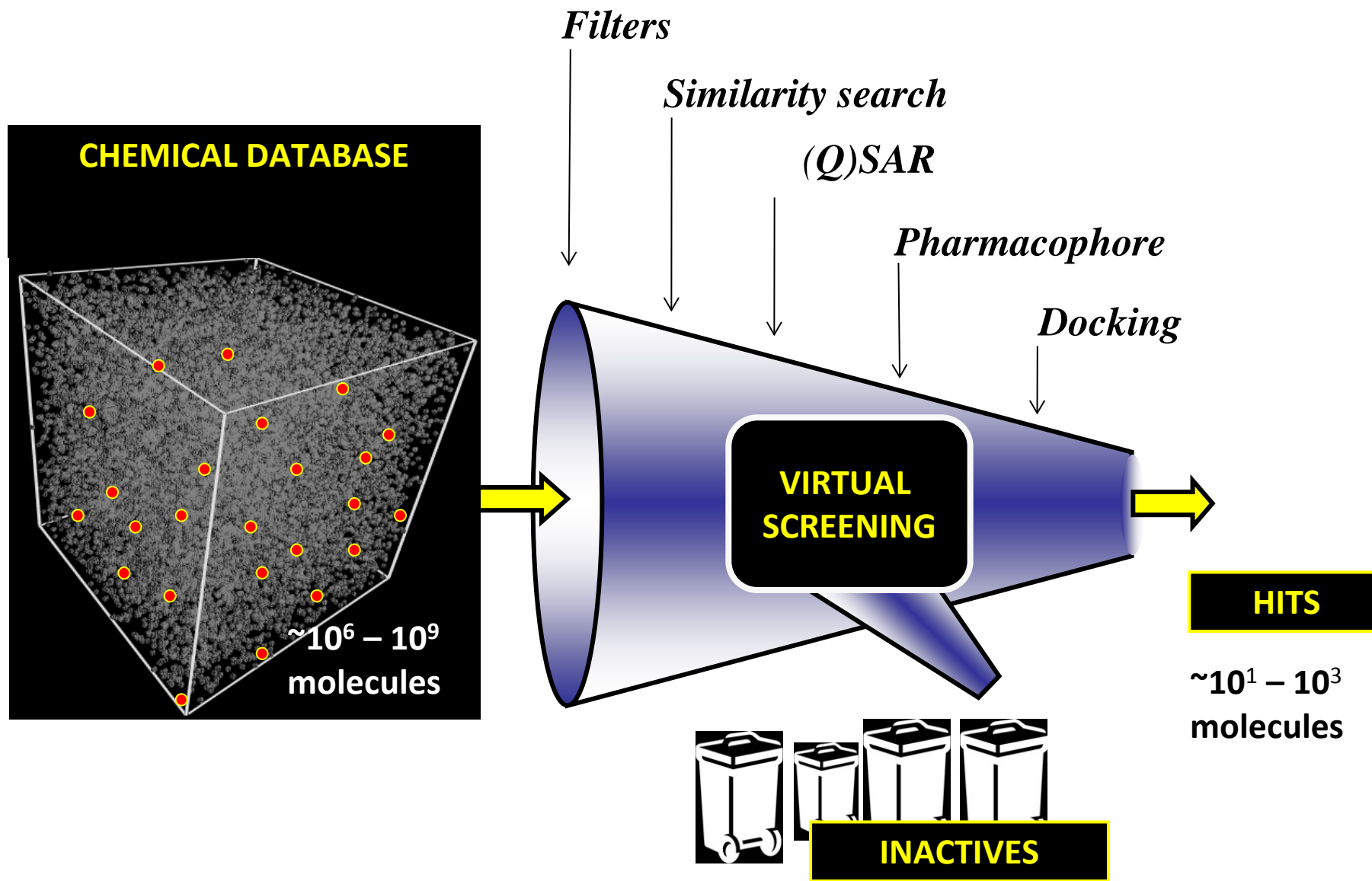


model



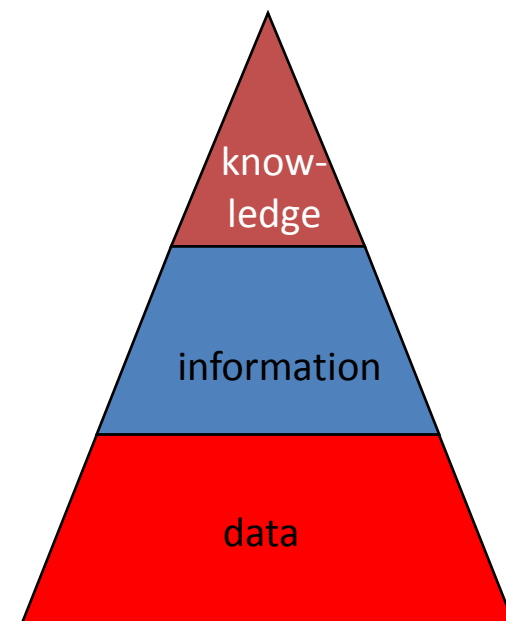
- Specific structural motifs,
  - Selected molecular properties (shape, fields, ...),
  - Interaction patterns,
  - Mathematical equations
- Activity = F (structure)***

# Chemoinformatics: Virtual screening "funnel"



# (Quantitative) Structure-Property Relationships (Q)SPR

$$\begin{aligned} \text{Property} &= \mathbf{F}(\text{structure}) \\ &= \mathbf{F}(\text{descriptors}) \end{aligned}$$



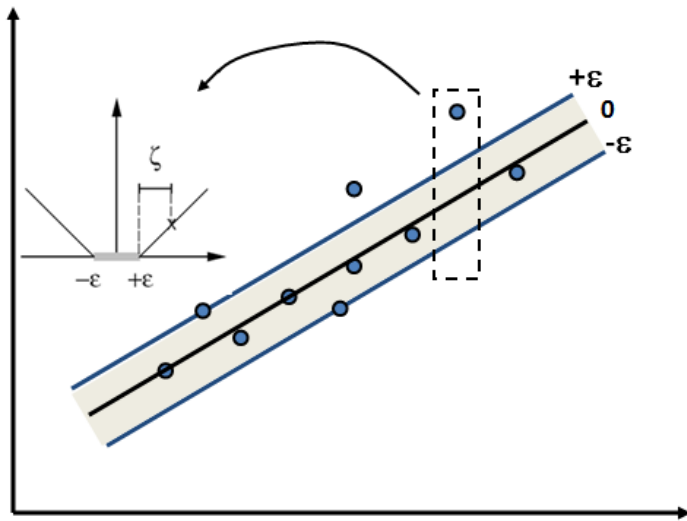
- Two main types of models:
  - **classification** (prediction of classes)
  - **regression** (prediction of numbers)

# Machine learning: *Regression models*

## Multiple Linear Regression (MLR)

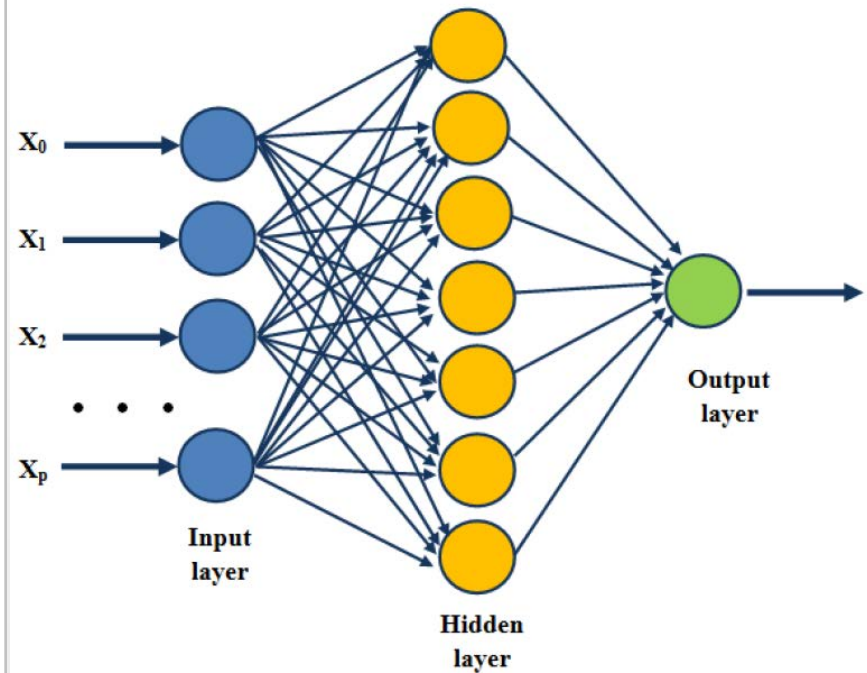
Property =  $a_0 + \sum_{i=1}^k a_i \cdot X_i$

## Support Vector Regression (SVR)



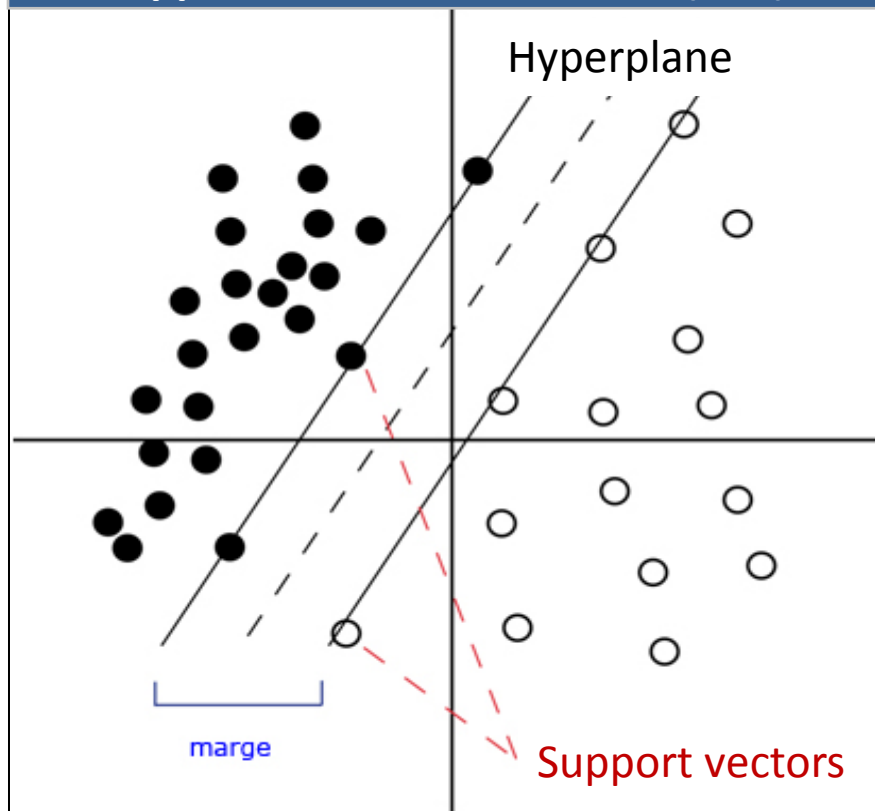
## Associative Neural Networks (ASNN)

### Ensemble of Neural Networks



# Machine learning: *Clasification Models*

## Support Vector Classification (SVC)

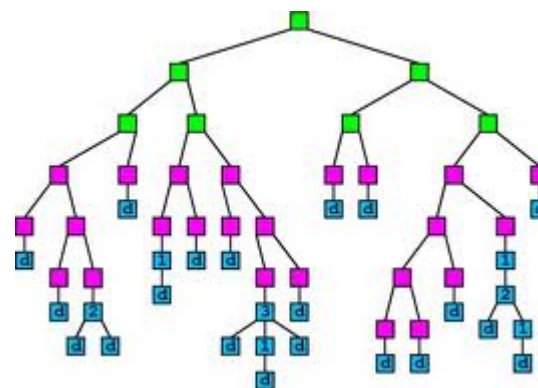


## Naive Bayes (NB)

Probabilistic Method

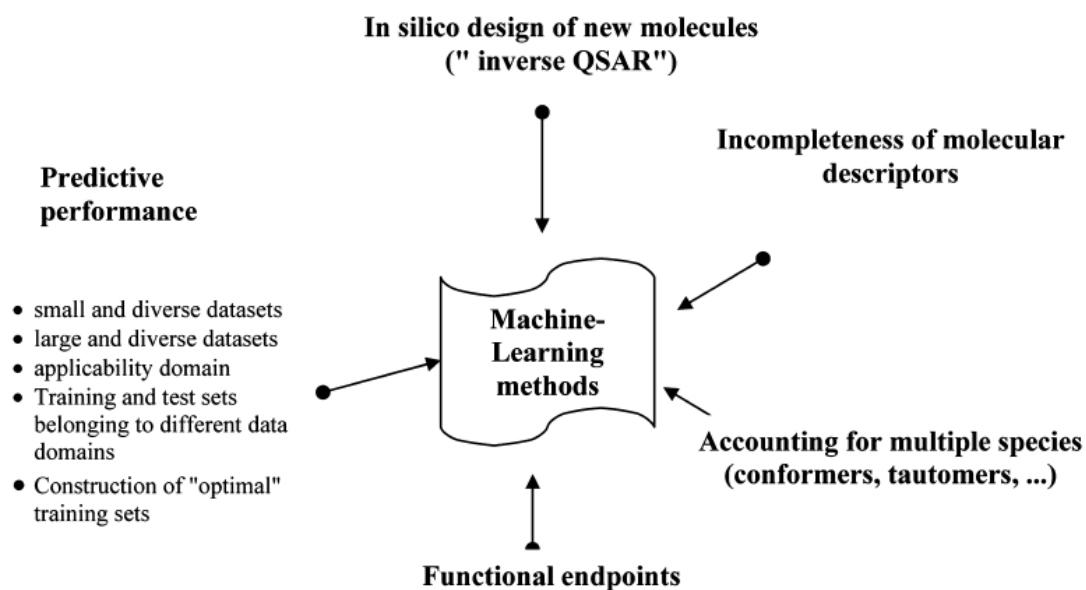
$$P(Y|a,b) = \frac{P(a|Y) \times P(b|Y) \times P(Y)}{P(a) \times P(b)}$$

## Decision Trees



# Machine Learning Methods for Property Prediction in Chemoinformatics: *Quo Vadis?*

Alexandre Varnek<sup>\*,†</sup> and Igor Baskin<sup>†,‡</sup>



**Main challenges of machine learning methods in chemoinformatics**



# Chemoinformatics tools in SciFinder:

The image displays the SciFinder web interface. The top navigation bar includes 'Explore References', 'Explore Substances', and 'Explore Reactions'. The user is logged in as 'Alexandre Varnek'. The main search results area shows two entries for 'Quinoline, labeled with deuterium' (CAS 1161799-94-5) and 'Quinoline' (CAS 91-22-5). Each entry includes a chemical structure, molecular formula (C<sub>9</sub>H<sub>7</sub>N), and a list of available data: References, Reactions, Commercial Sources, and Regulatory Information. The 'Quinoline' entry has approximately 13,164 references.

On the right, the 'Refine by Property Value' panel is open, showing a list of properties categorized into 'Experimental' and 'Predicted'. The 'Predicted' category is circled in red. The 'Predicted' section includes the following properties:

- Boiling Point
- Melting Point
- H Acceptors
- H Donors
- Molecular Weight
- logP
- Freely Rotatable Bonds
- Bioconcentration Factor
- Boiling Point
- Density
- Enthalpy of Vaporization
- Flash Point
- H Acceptor/Donor Sum
- Koc
- logD
- Mass Intrinsic Solubility
- Mass Solubility
- Molar Intrinsic Solubility
- Molar Solubility
- Molar Volume

At the bottom of the panel, there is a checkbox:  Include substances with no value for the

predictions of > 20 physico-chemical properties and NMR spectra for each individual compound

# ISIDA property prediction WEB server

*infochim.u-strasbg.fr/webserv/VSEngine.html*

Prediction of property logP - Page nr. 3 - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://infochim.u-strasbg.fr/userdata/dragos/logP/logPP3.html

Most Visited Personnaliser les liens Windows Media Windows

Courrier :: Boîte de réception Virtual Screening Engine - Laboratoire d... Prediction of property logP - Pag...

Predicted property **logP** for 9677 compounds AS A CONSENSUS OF APPLICABLE LOCAL MODELS

logP	VAR	TRUST	REASON
1.59	0.546	NONE	<ul style="list-style-type: none"><li>- None of the local models have applicability domains covering this compound</li><li>- Individual models failed to reach unanimity - prediction variance exceeds 1.0% of the property range width</li></ul>
3.13	0.127	POOR	<ul style="list-style-type: none"><li>- There are too few (less than 5) local models containing molecule within applicability domain - global consensus is preferred</li><li>- Furthermore, the other local models disagree with the prediction of the minority containing compound inside their applicability domain</li><li>- Individual models failed to reach unanimity - prediction variance exceeds 1.0% of the property range width</li></ul>
2.60	0.105	OPTIMAL	-

Done

# CoMet project

models for stability constants of metal-ligand complexes in water

*in red: the models are available*

H																	He
<b>Li</b>	<b>Be</b>											B	C	N	O	F	Ne
<b>Na</b>	<b>Mg</b>											<b>Al</b>	Si	P	S	Cl	Ar
<b>K</b>	<b>Ca</b>	<b>Sc</b>	Ti	<b>V</b>	<b>Cr</b>	Mn	<b>Fe</b>	<b>Co</b>	<b>Ni</b>	<b>Cu</b>	<b>Zn</b>	<b>Ga</b>	<b>Ge</b>	As	Se	Br	Kr
<b>Rb</b>	<b>Sr</b>	<b>Y</b>	<b>Zr</b>	<b>Nb</b>	<b>Mo</b>	<b>Tc</b>	<b>Ru</b>	<b>Rh</b>	<b>Pd</b>	<b>Ag</b>	<b>Cd</b>	<b>In</b>	<b>Sn</b>	Sb	Te	I	Xe
<b>Cs</b>	<b>Ba</b>	<b>La</b>	Hf	Ta	W	Re	Os	Ir	Pt	Au	<b>Hg</b>	<b>Tl</b>	<b>Pb</b>	<b>Bi</b>	Po	At	Rn
<b>Fr</b>	<b>Ra</b>	<b>Ac</b>															

lanthanides

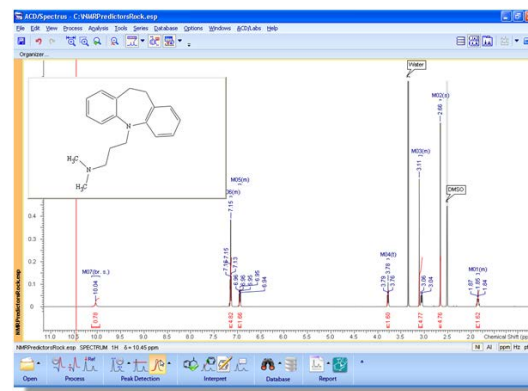
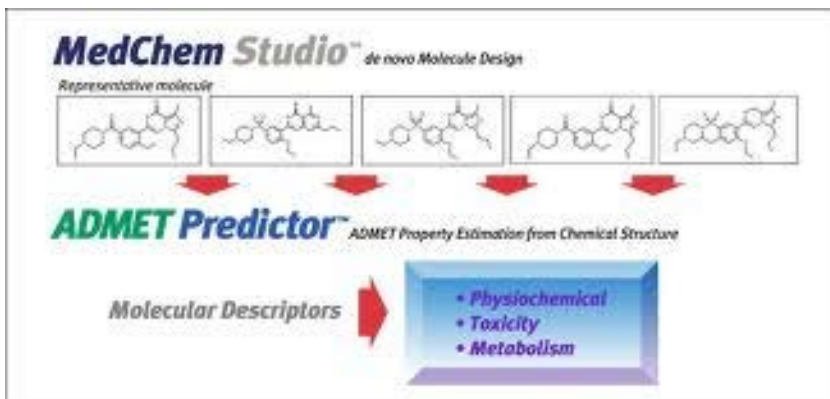
**Ce** **Pr** **Nd** **Pm** **Sm** **Eu** **Gd** **Tb** **Dy** **Ho** **Er** **Tm** **Yb** **Lu**

actinides

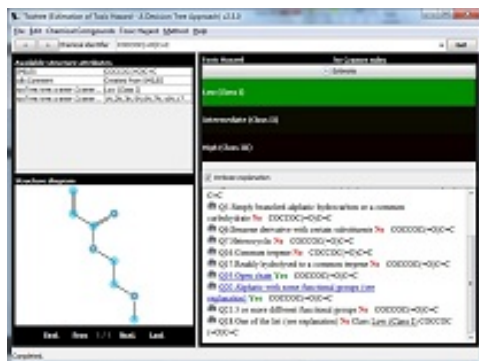
**Th** **Pa** **U** **Np** **Pu** **Am** **Cm** **Bk** **Cf** **Es** **Fm** **Md** **No** **Lr**

# Commercial and Public Software

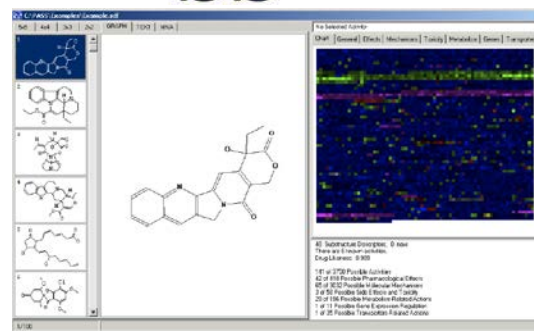
**simulations plus, inc.**  
*integrating science and software*



## Toxtree



## PA SS



# REACH regulation

- The European Union adopted Regulation on the **R**egistration, **E**valuation, **A**uthorisation, and **R**estriction of **C**hemicals (the “REACH Regulation”), which entered into force on June 1, 2007.
- REACH imposes requirements of information of physico-chemical, toxicology and eco-toxicology parameters for the chemicals, production of which exceeds 1 ton.
- More than 30.000 compounds must be tested. Total cost estimated (EU Commission) over a 11 -15 year period is €2.8 - €5.2 bn

**No Data, No Market!**

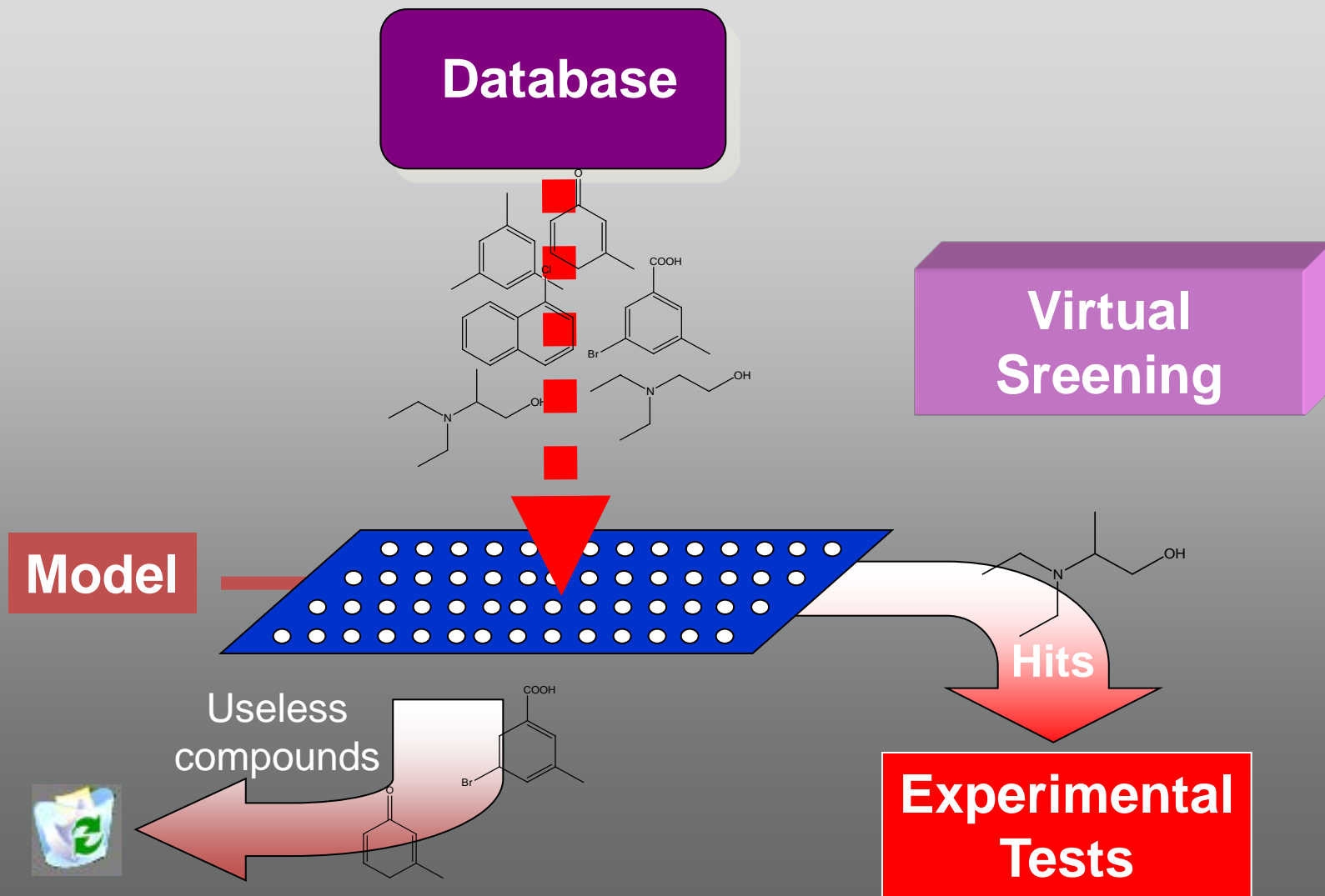
# Chemoinformatics:

## *areas of application*

- Drug design (pharmacodynamics and pharmacokinetics),
- Prediction of physico-chemical properties,
- Materials design,
- Synthesis design,
- Molecular spectra simulations

# **Structure-Property modeling: case studies**

# Screening and hits selection



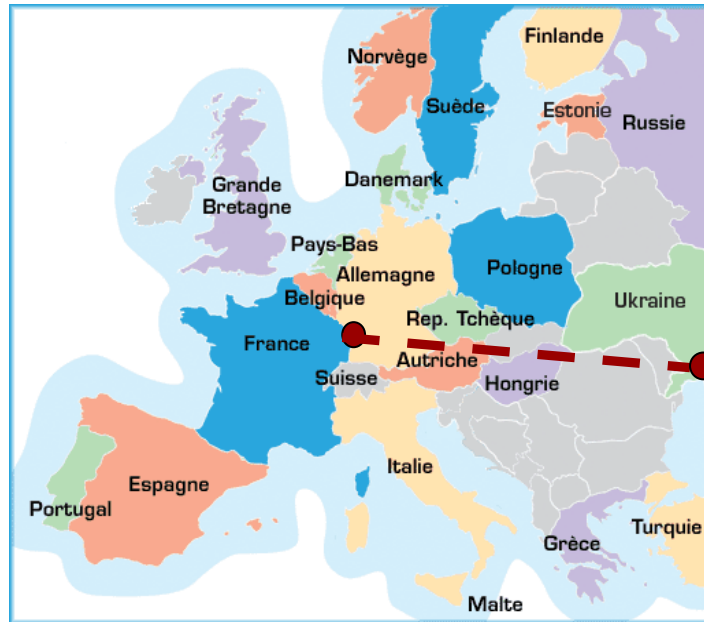


# *In silico* design of new antithrombotics

University of Strasbourg - Bogatsky Institute in Odessa



**A. Varnek**



**S. Andronati**



**V. Kuzmin**



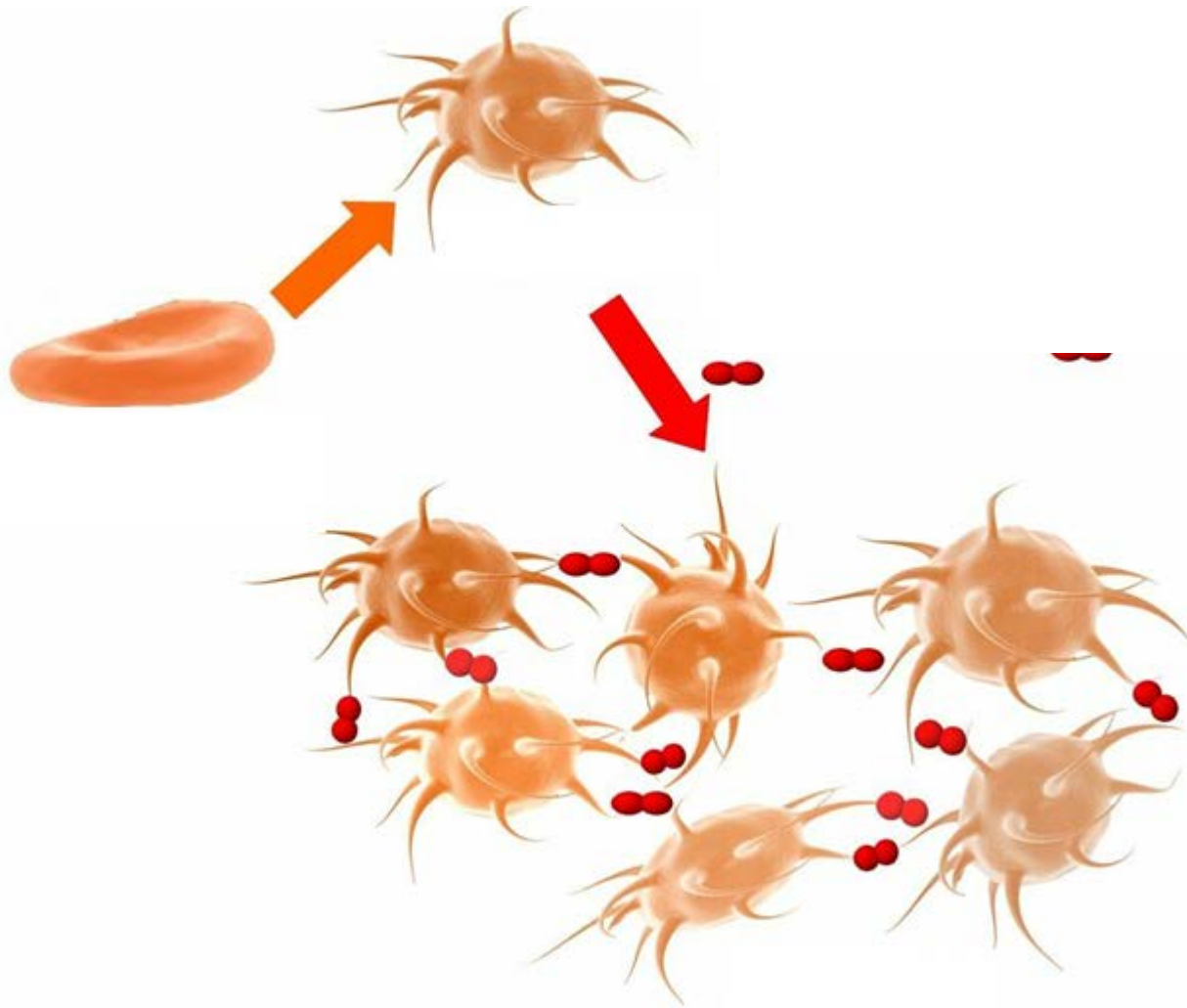
**P. Polishchuk**



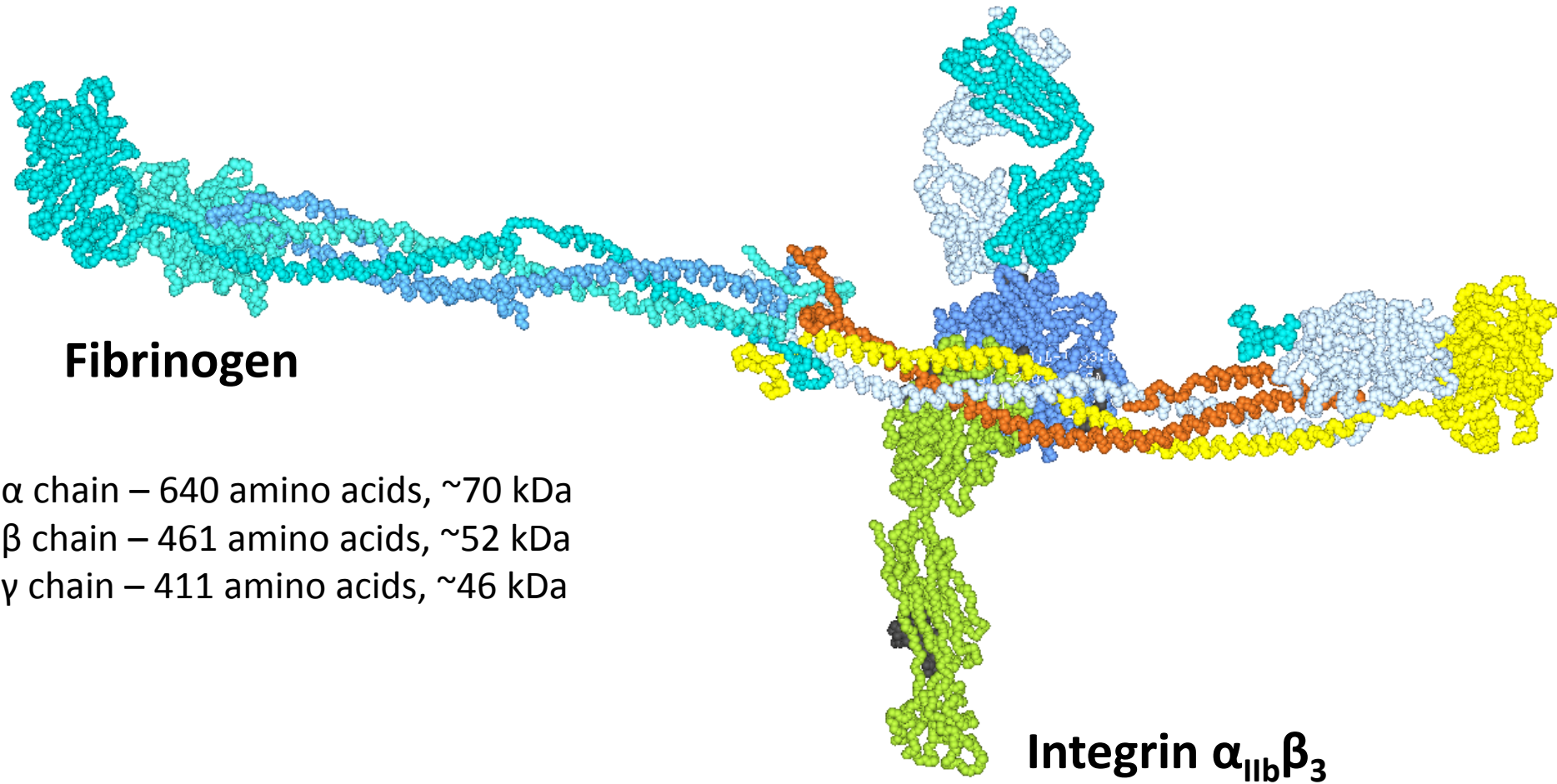
**T. Khristova**

Co-supervised PhD project

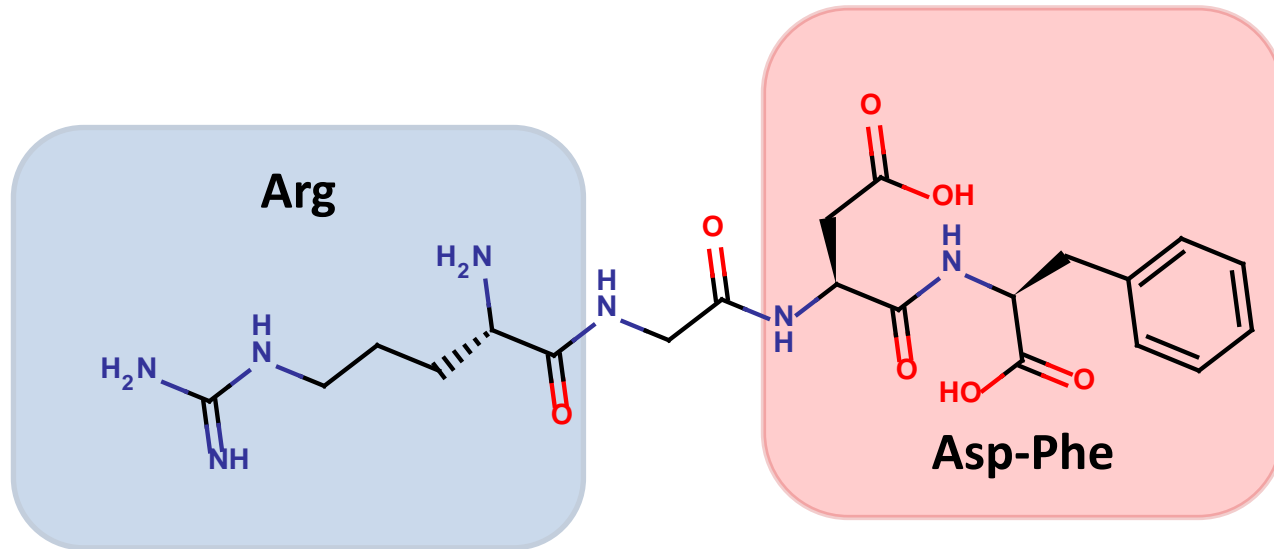
# Platelets aggregation



# Protein-Protein interactions: Complex of *Integrin $\alpha_{IIb}\beta_3$* and *Fibrinogen*



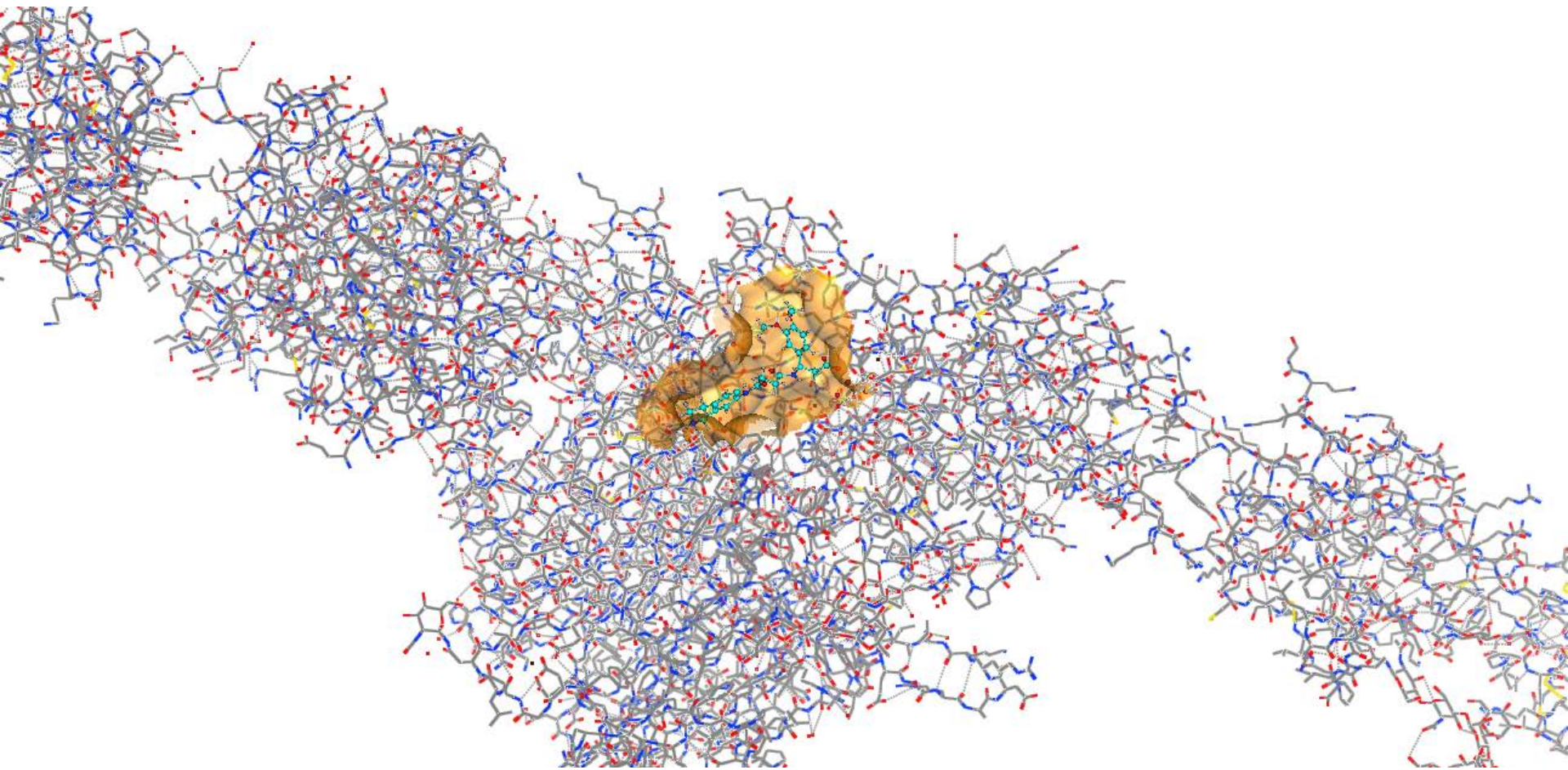
# Residues of fibrinogene interacting with binding site of integrine $\alpha_{IIb}\beta_3$



**Arg-Gly-Asp-Phe  
(RGD)**

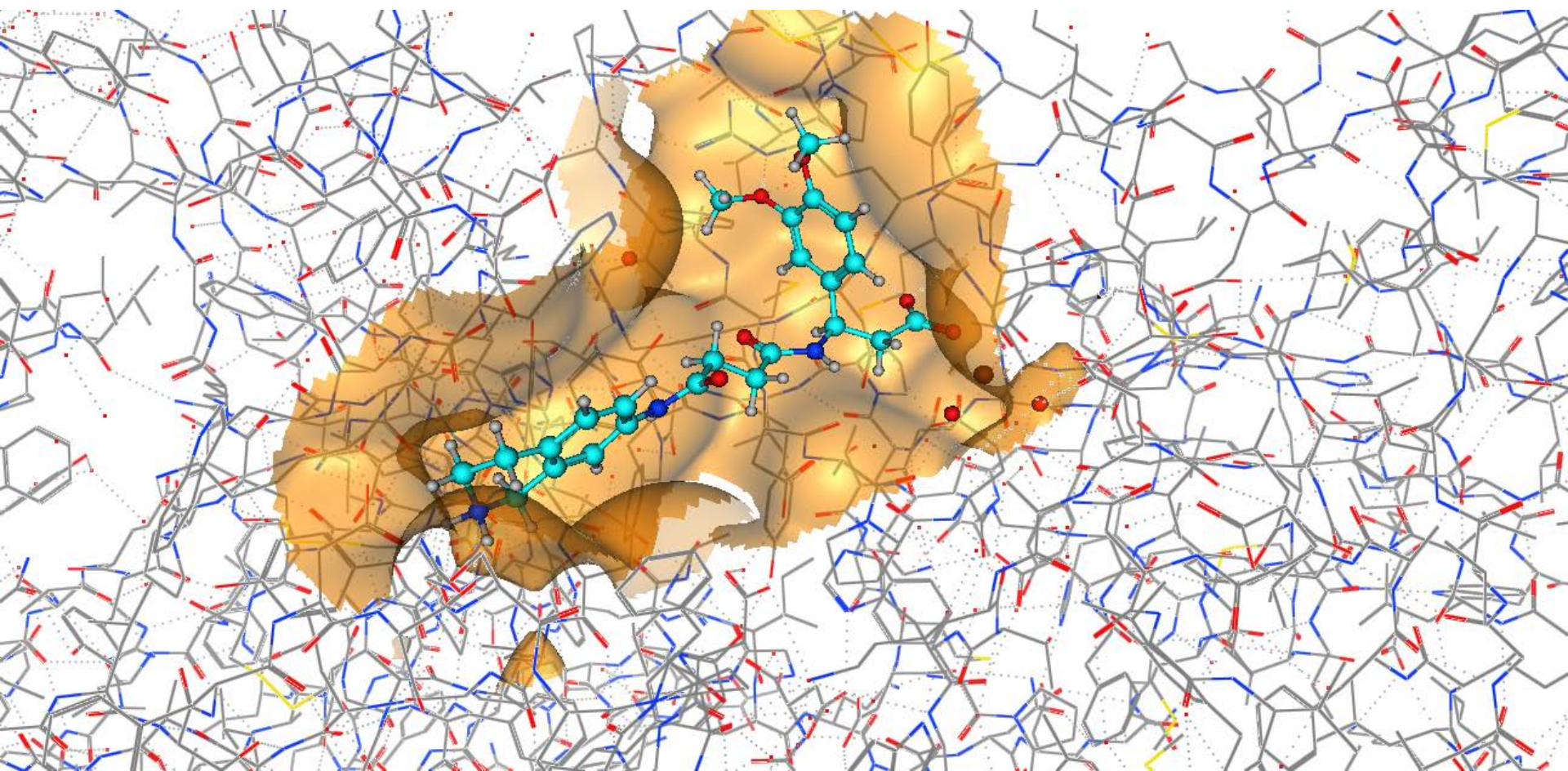
# Protein-Ligand complex:

*Binding of RGD-mimetic to integrin  $\alpha_{IIb}\beta$*



Ligand to protein docking with *MOE*

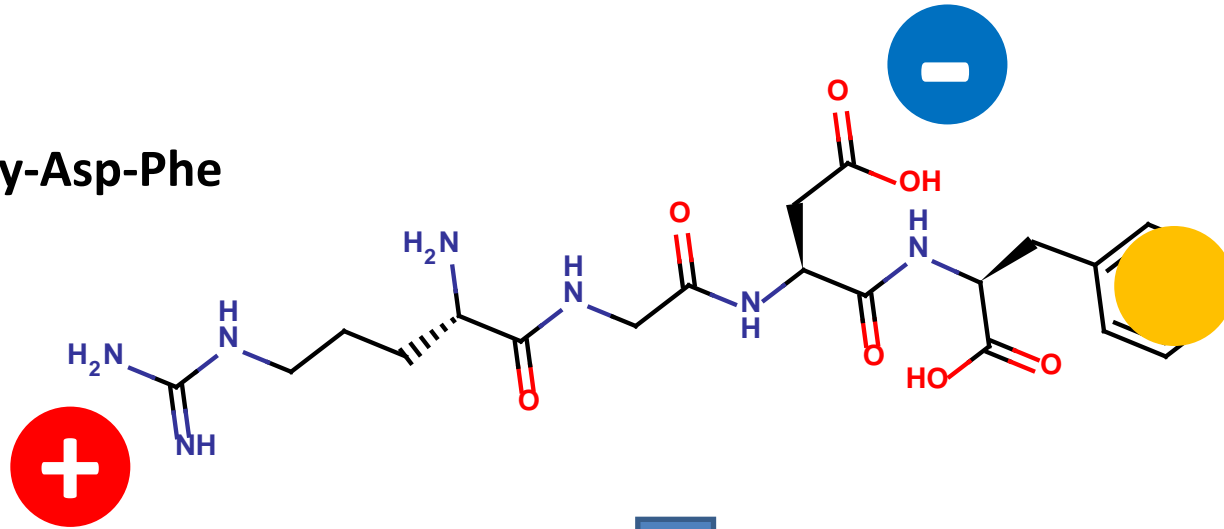
# Protein-Ligand complex: *Binding of RGD-mimetic to integrin $\alpha_{IIb}\beta$*



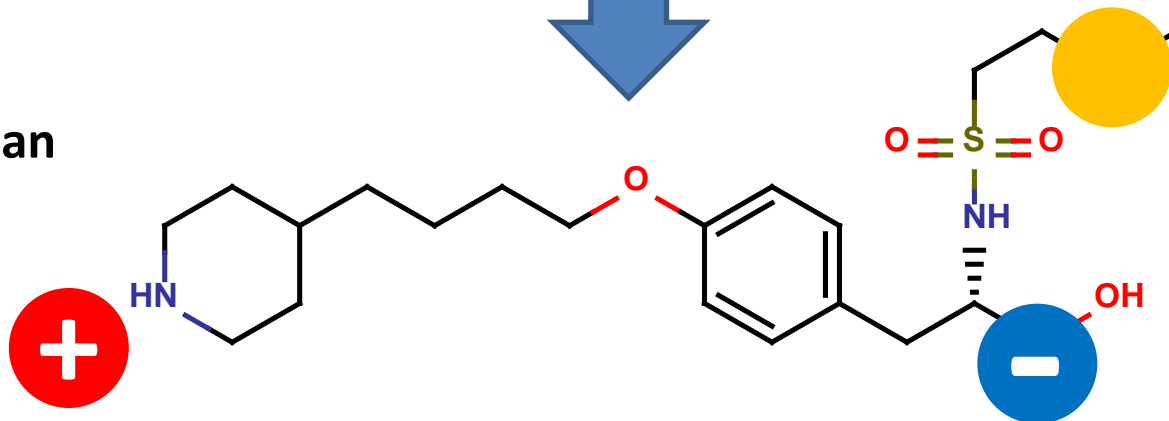
Ligand to protein docking with *MOE*

# What is in common between these two molecules ?

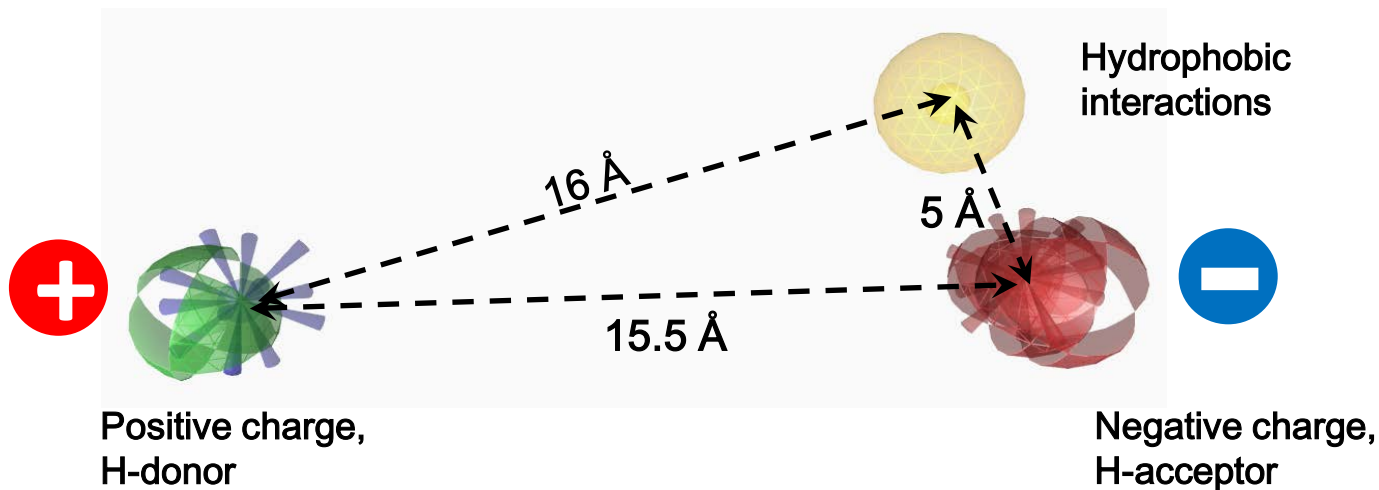
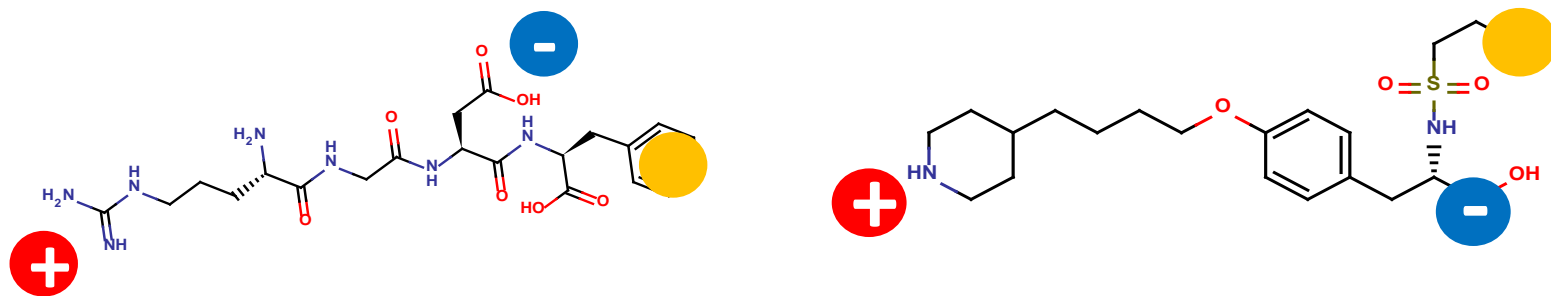
Arg-Gly-Asp-Phe



Tirofiban

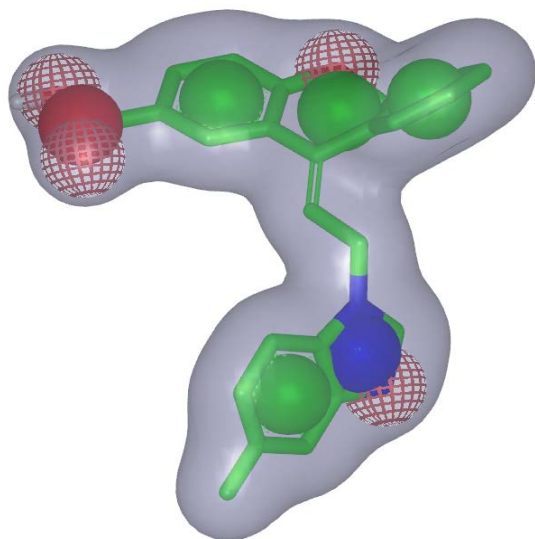


# Pharmacophore model of ligand complementary to integrine $\alpha_{IIb}\beta_3$

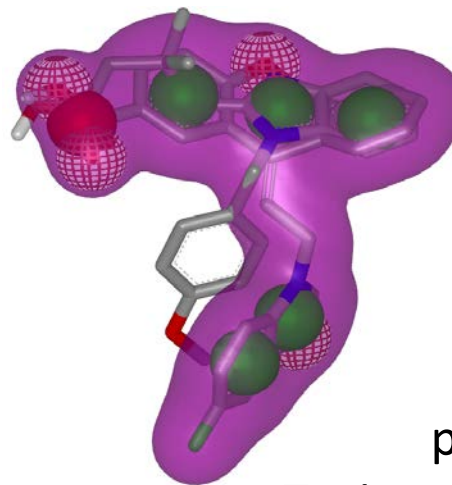
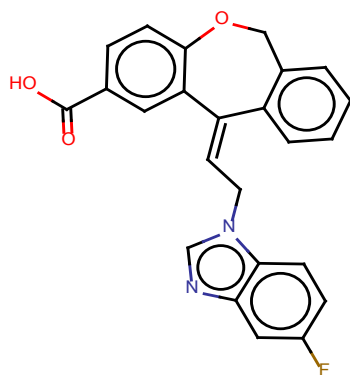




# Molecular Shape similarity analysis

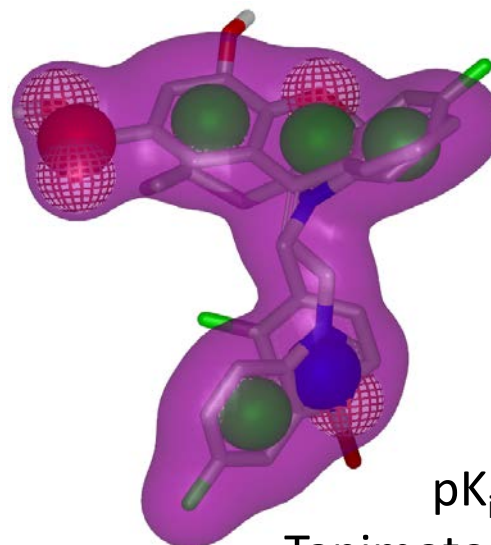
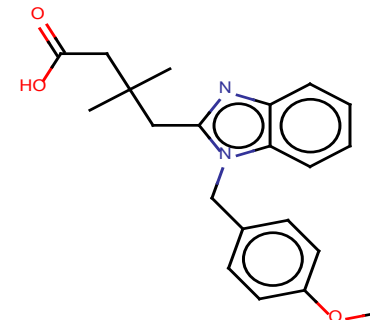


$pK_i = 7.82$



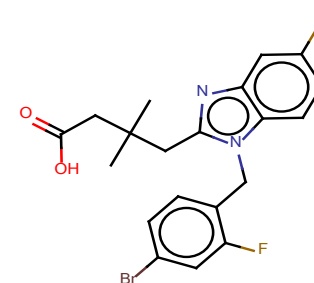
$pK_i = 7.51$

TanimotoCombo = 0.74



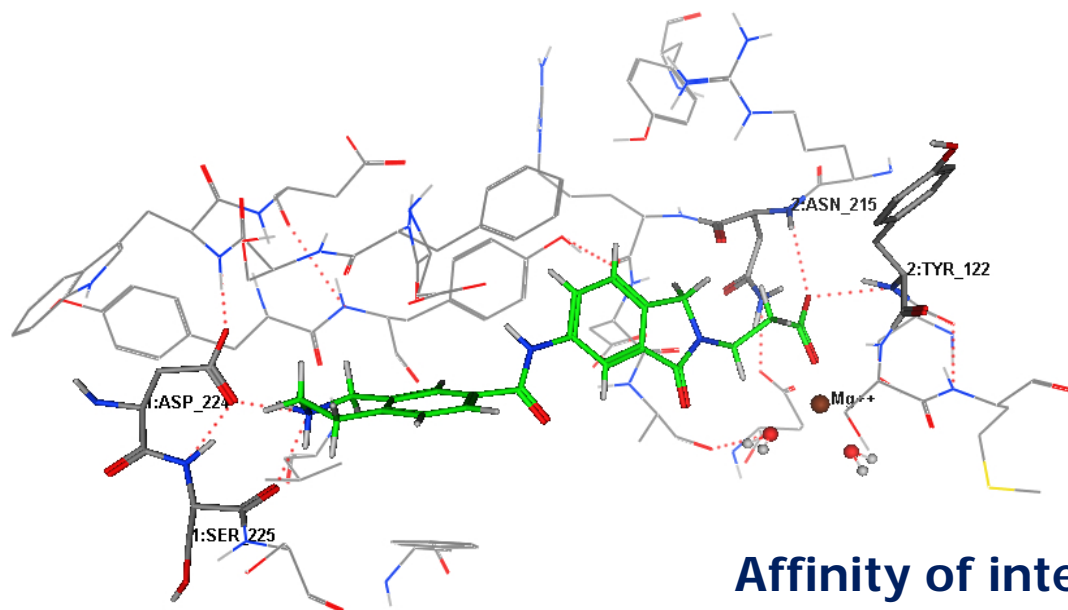
$pK_i = 7.82$

TanimotoCombo = 0.67



## Experimental validation:

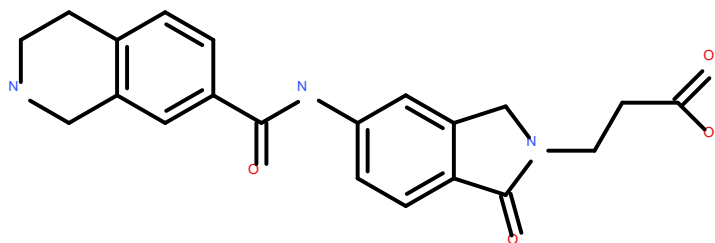
inhibition of platelet aggregation in human platelet-rich plasma



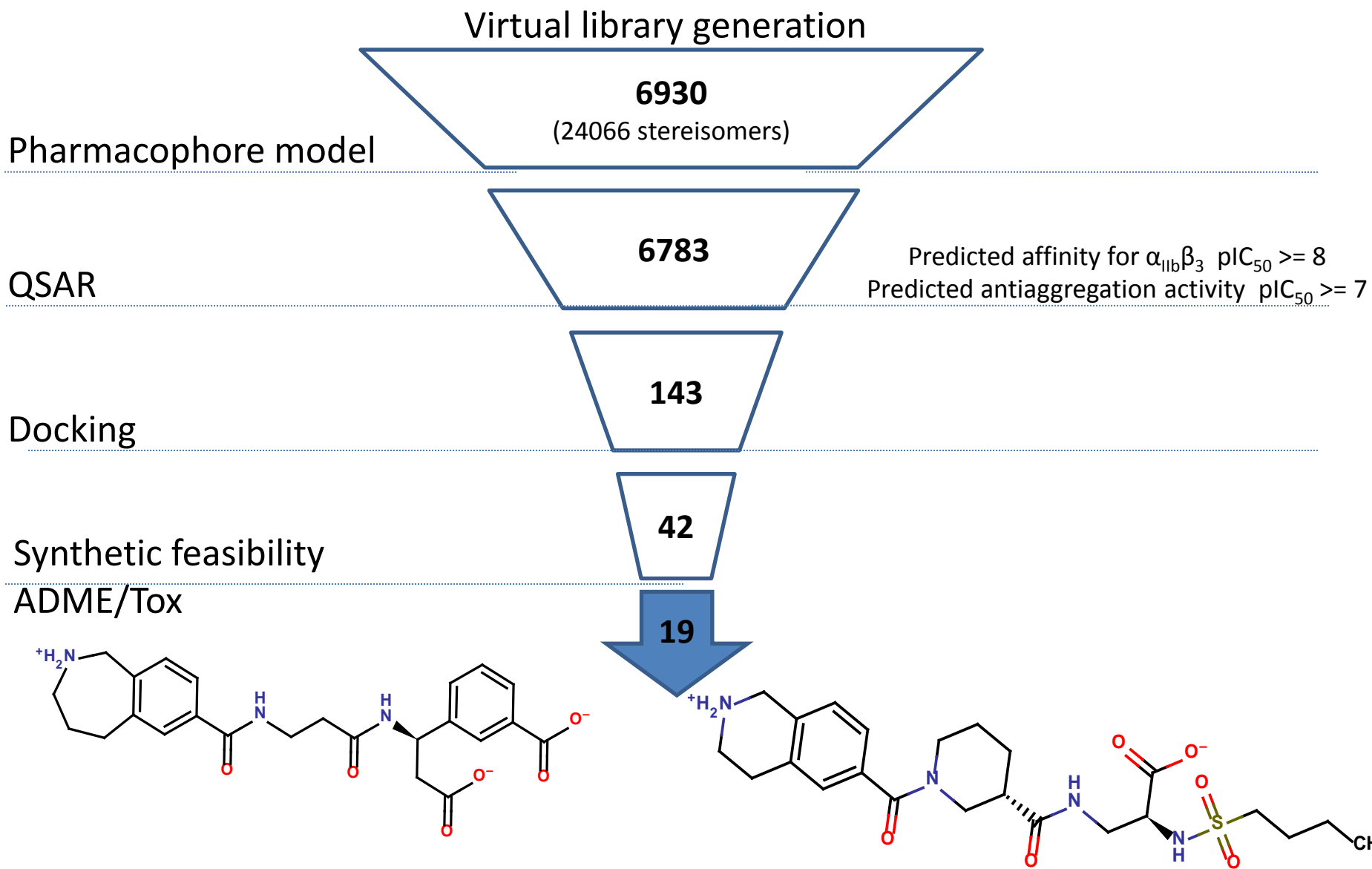
**Affinity of integrin  $\alpha_{11b}\beta_3$  inhibition ( $IC_{50}$ )**

exp = 6.5 nM;

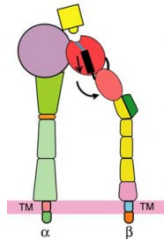
pred = 9.3 nM



# Virtual screening of designed RGD-peptidomimetics

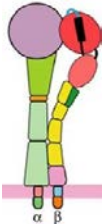


# Eight potent $\alpha$ IIb $\beta$ 3 antagonists with two types of binding modes were developed:

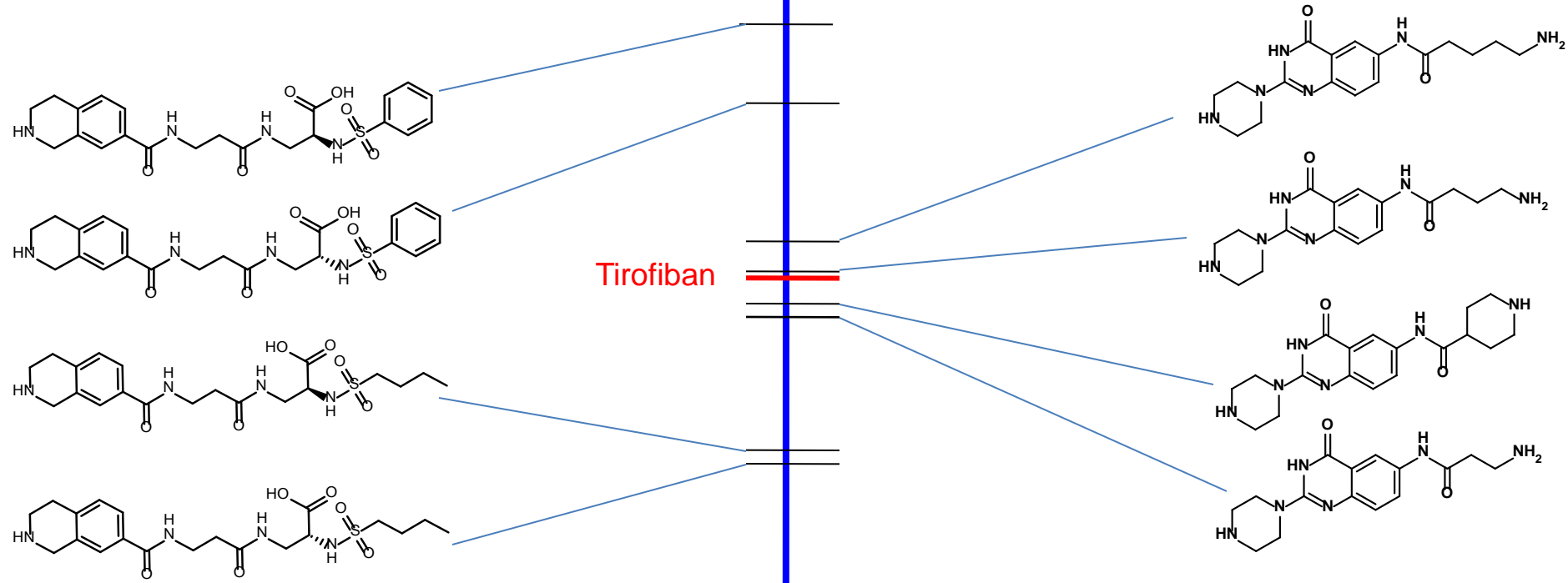


4 “classical”

4 “non classical”



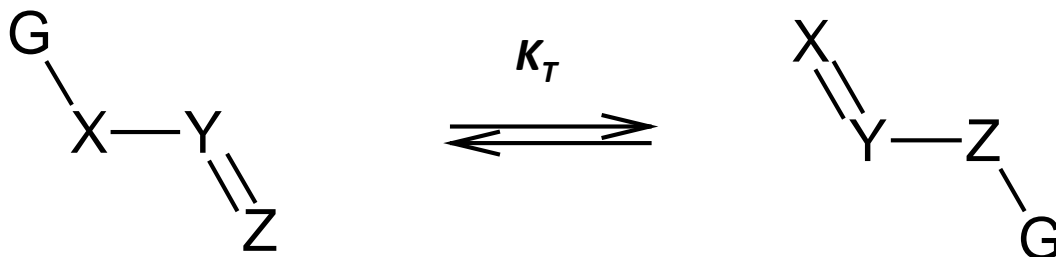
Affinity for  $\alpha$ IIb $\beta$ 3



# Prediction of Tautomer Equilibrium Constants

## IUPAC definition

Isomerism of the general form:



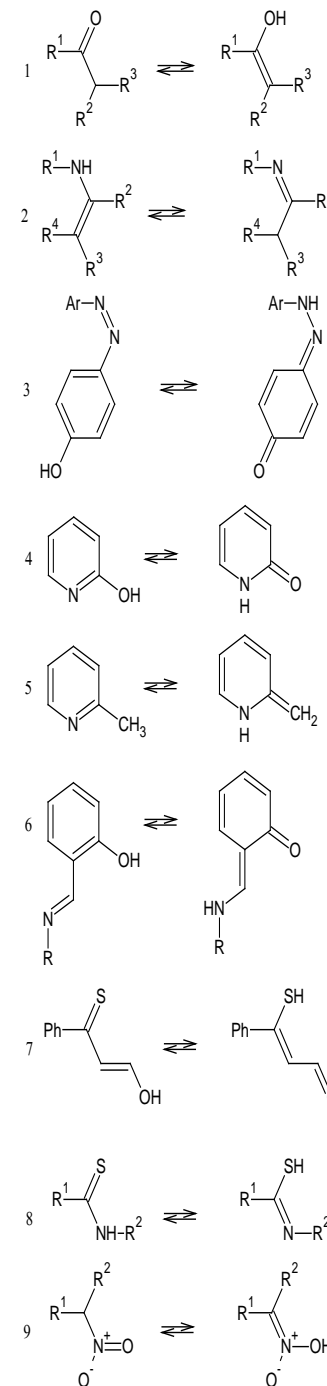
where the isomers (called tautomers) are readily interconvertible

**Goal: to build a model predicting  $\log K_T$**

# CONSIDERED TAUTOMERIC EQUILIBRIA

- Keto-Enol
- Amino-Imino
- Azo-Hydrazine
- Pyridol-Pyridone
- Pyridinoid-Pyridonoid
- Phenol-Imine – Keto-Amine
- Thione-Enol – Keto-Thiol
- Amine-Thione – Imine-Thiol
- Nitro-Acid
- Classical Form -Zwitterion
- Ring-Chain

N
269
105
91
36
6
37
17
82
13
19
50



# Data workflow

Database



Reactions



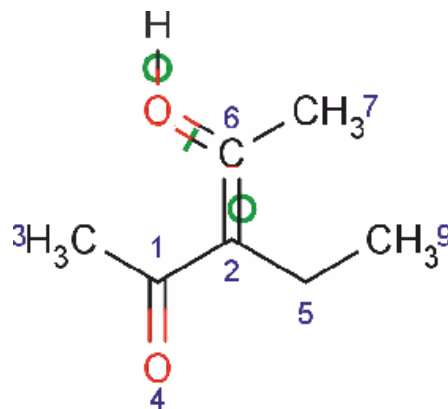
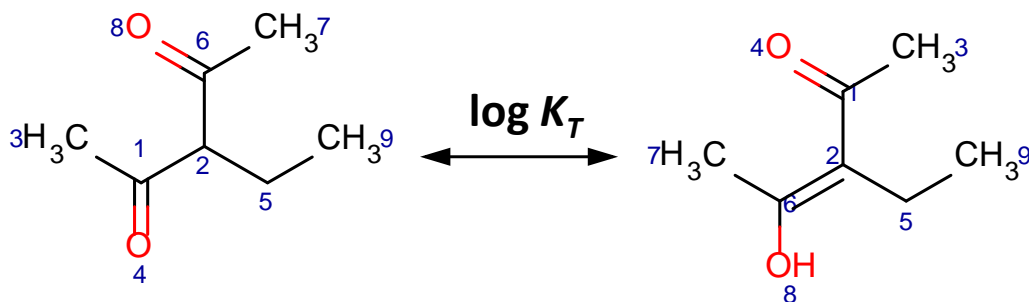
Condensed Graph of Reaction (CGR)

723 tautomeric equilibrium ( $\log K_T$ )

11 types of tautomerism

12 pure solvents

7 water-organic solvent mixtures



—○— formed bond  
—+— broken bond

# Tautomers: Models Performance

	Our model (cross-val)	Our model (ext set)	QM (semi- empirical) <sup>1</sup>	QM (DFT) <sup>2</sup>	ChemAxon
# molecules	723	31	643	31	127
<b>Correct most stable tautomer (%)</b>	<b>85</b>	<b>81</b>	<b>55</b>	<b>45</b>	<b>56</b>
RMSE (log $K_T$ )	0.86	1.34	5.99	8.99	5.38

<sup>1</sup> IEF-PCM/PM6

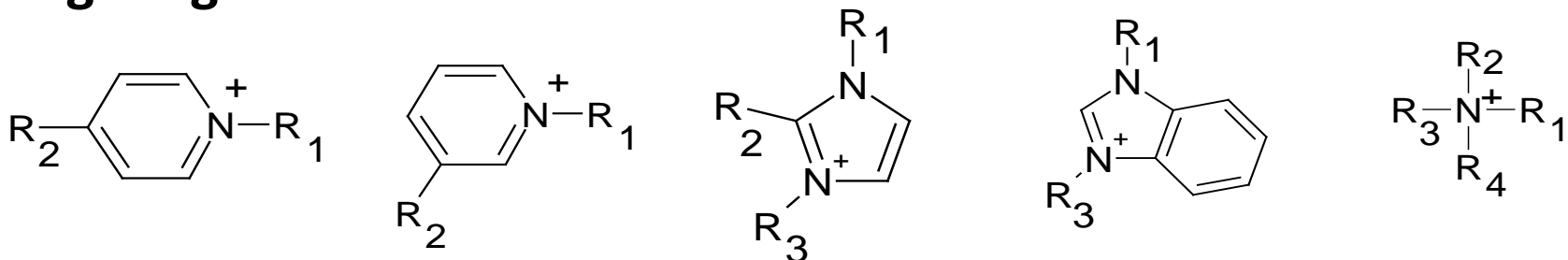
<sup>2</sup> IEF-PCM/B3LYP/6-311++G(d,p)



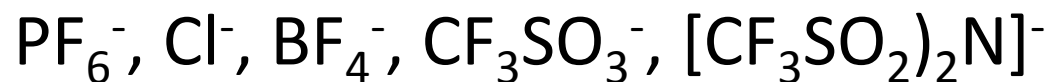
# Ionic Liquids

Ionic Liquids are composed of

large organic cations:

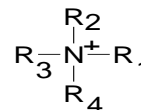
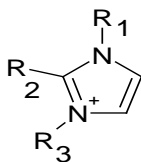
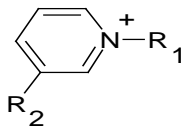
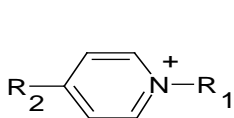


and anions:



# Ionic Liquids

Large organic cations:

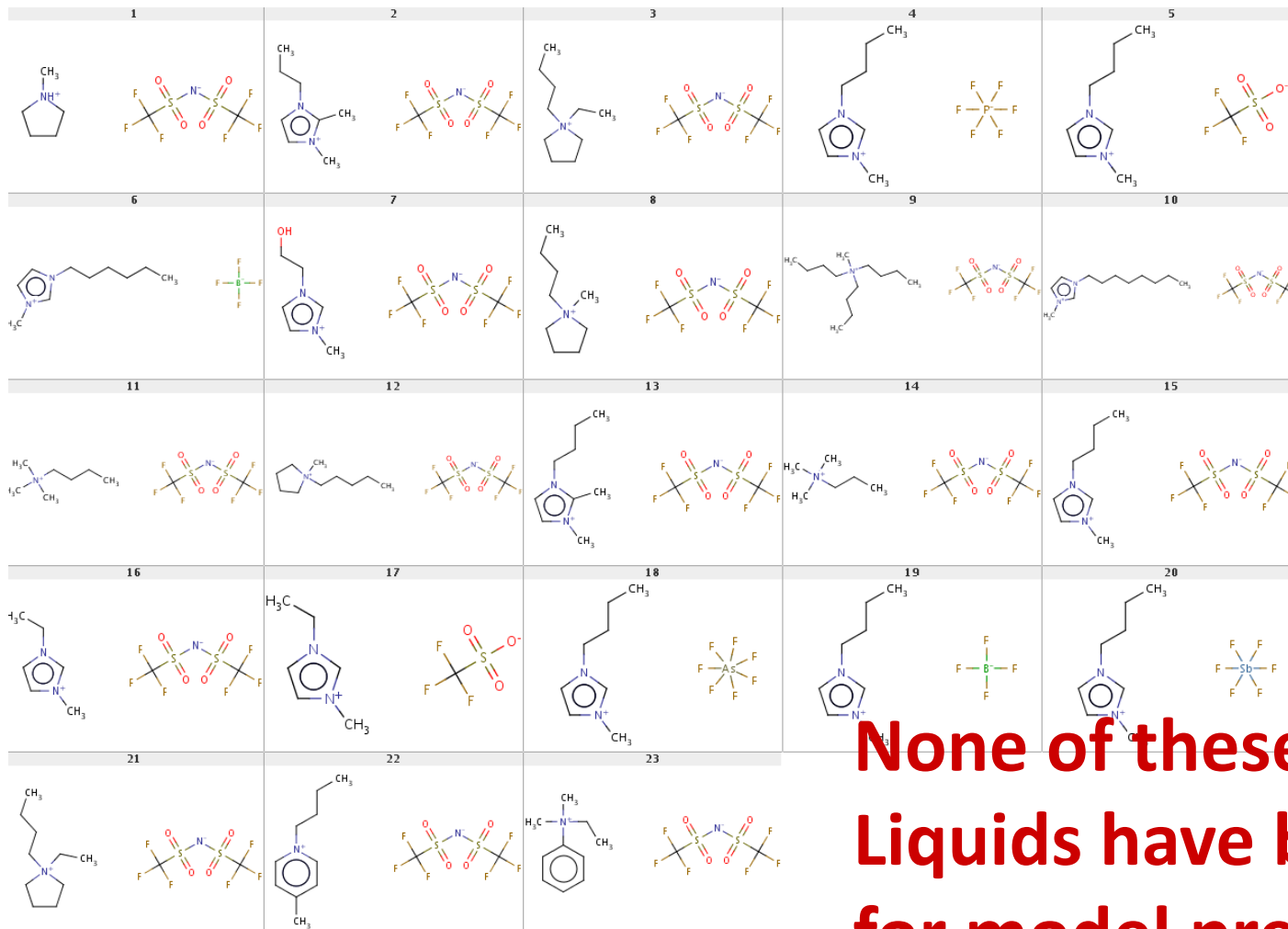


anions:

$\text{PF}_6^-$ ,  $\text{Cl}^-$ ,  $\text{BF}_4^-$ ,  $\text{CF}_3\text{SO}_3^-$ ,  $[\text{CF}_3\text{SO}_2]_2\text{N}^-$

There exist  **$10^{18}$**  combinations of ions  
that could lead to useful ionic liquids

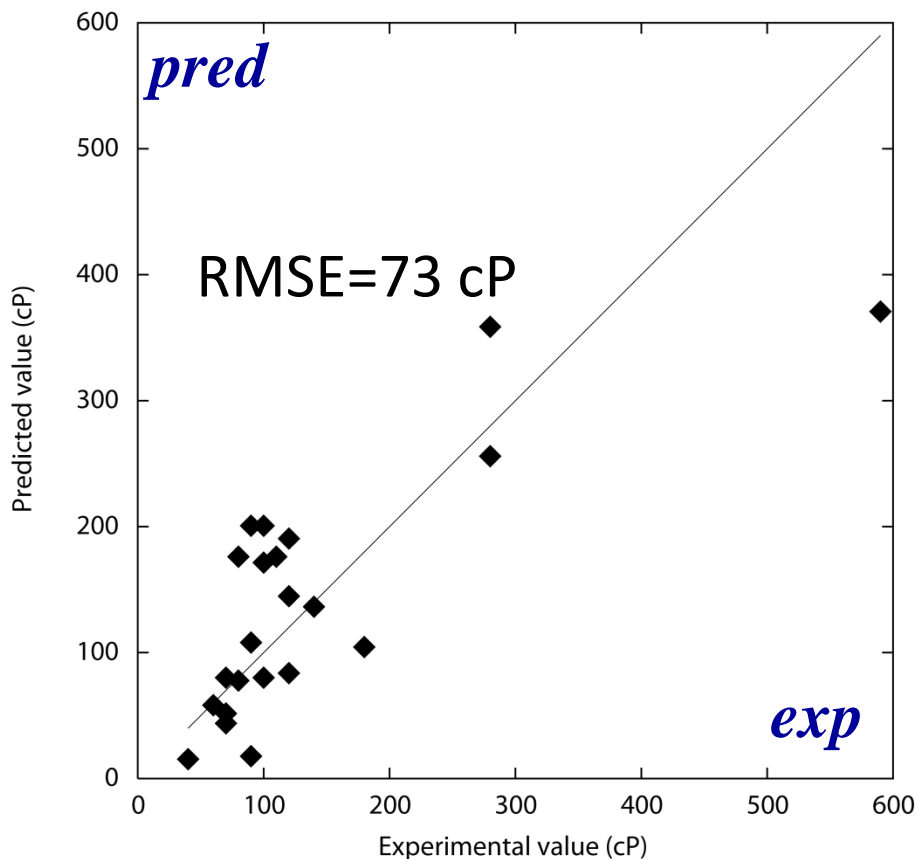
# Viscosity predictions on 23 new ILs



*Solvionics  
company*

**None of these Ionic  
Liquids have been used  
for model preparation**

# Ionic Liquids viscosity: Experimental validation of the Neural Networks models

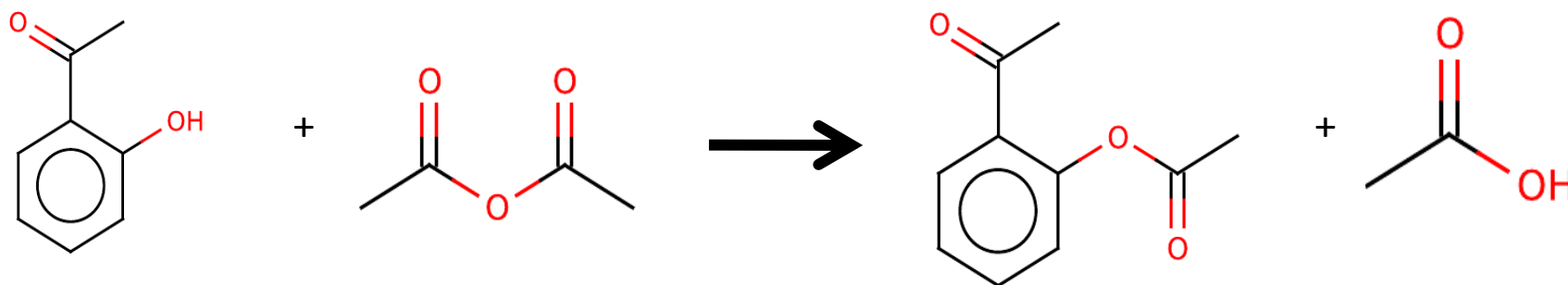


- *prediction error (~70 cP) is similar to the “noise” in the experimental data used for the training of the model*

G. Marcou, I. Billard , A. Ouadi and A. Varnek,  
*J. Phys. Chem. B*, 2011, **115** (1), 93–98

# Prediction of optimal reaction conditions

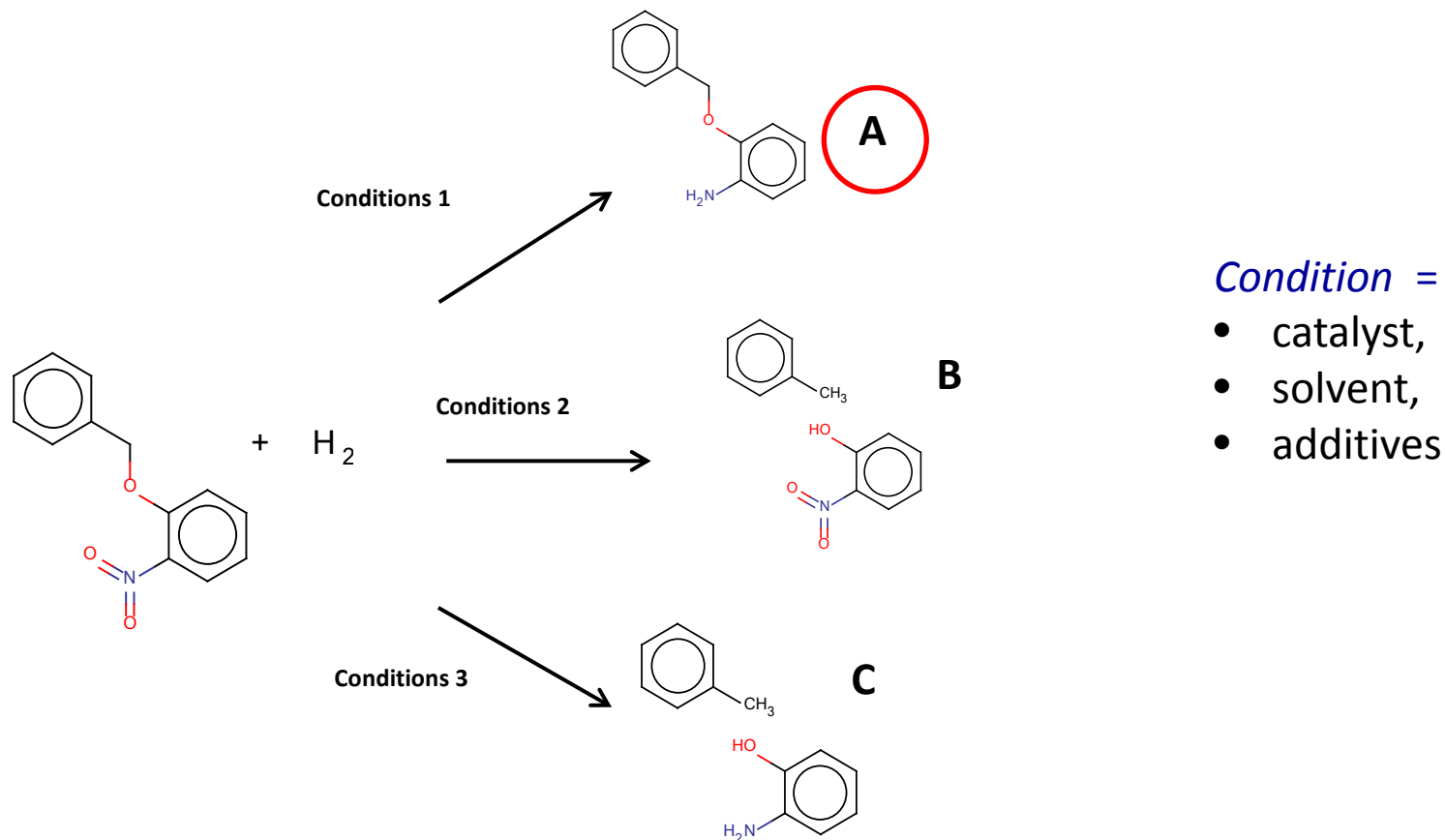
# What is a reasonable way to encode a chemical reaction ?



Chemical reactions are difficult objects :

- many species;
- two types of species: reactants and products;
- multi-step reactions,
- dependent on experimental conditions

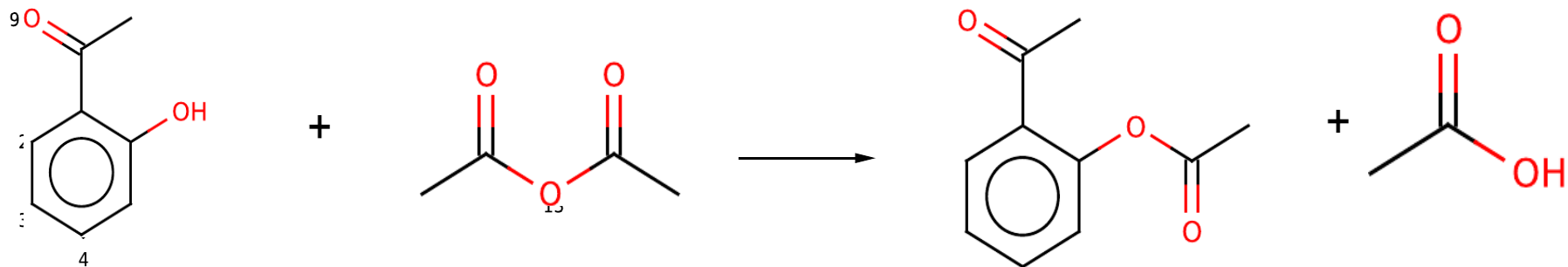
# Selective hydrogenation



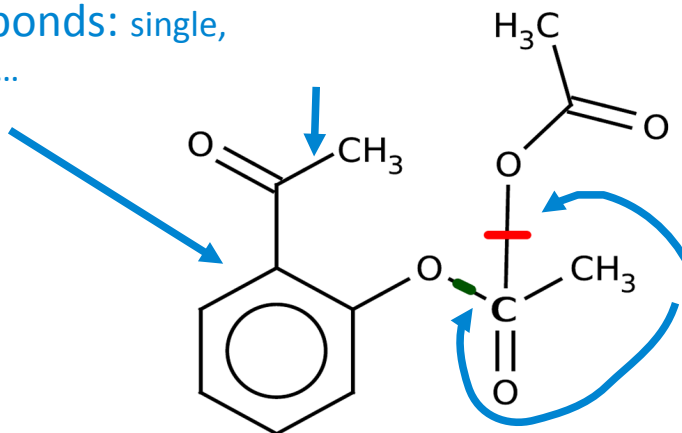
Example reaction leading to products A, B and C under specific conditions.

The goal is to predict reaction conditions leading selectively to product **A**.

# Reactions representation: Condensed Graphs of Reactions



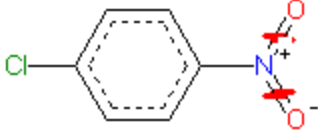
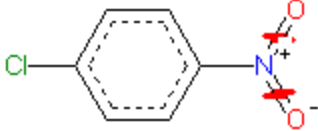
Conventional bonds: single, double, aromatic, ...



Dynamical bonds: created single, broken single, ...

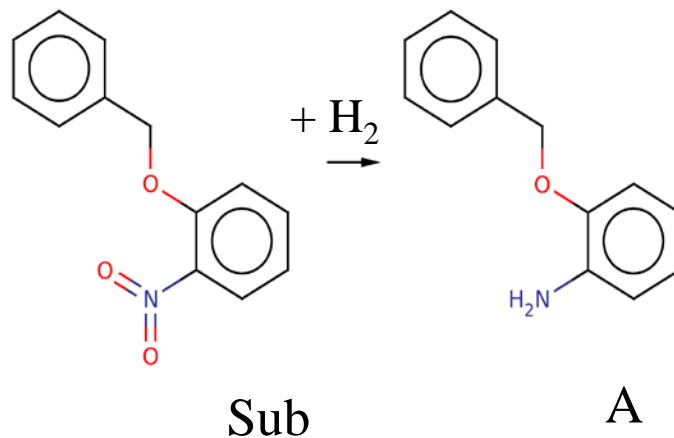
Notice: atom-to-atom mapping is required



3745	Ir/CaCO3 (5%)	DMF	None	25	1	4-Chloro-phenylamine	75	101650_1259		0.513
	Ir/CaCO3 (5%)	DMF	None					101650_582	c:c-C-O>>c:c-C-O	99
	Ir/CaCO3 (5%)	DMF	None					101650_873	c:c-C-O>>c:c-C-O	100
	Ir/CaCO3 (5%)	DMF	None					101650_873	C-O-c:c>>C-O-c:c	100
	Ir/CaCO3 (5%)	DMF	None					101650_873	c:c-c-O>>c:c-c-O	100
	Ir/CaCO3 (5%)	DMF	None					101650_1298	c:c-c-O>>c:c-c-O	100
3697	Pt/C 10%	THF	None	25	1	4-Chloro-phenylamine	78	101650_1243		0.513
	Pt/C 10%	THF	None					101650_566	c:c-C-O>>c:c-C-O	100
	Pt/C 10%	THF	None					101650_707	c:c-C-O>>c:c-C-O	58

Search details

- [-] DATABASE
  - Name: C:\v
  - CpdNb: 388
  - DescNb: 39
- [-] QUERY
  - Name: C:\v
  - CpdNb: 1
  - DescNb: 17
- [-] METRIC
  - Type: Tanim
  - CutOff: 49
- [-] FILTER
  - Property: Y
  - Value: 70
- [-] FRAGMENTATION
  - Type: I(AB,
- [-] WEIGHT
  - Type: redu



### Conditions suggested by the program

	<i><b>catalyst</b></i>	<i><b>solvent</b></i>	<i><b>additif</b></i>
1	Pt/C (10%)	THF	None
2	Pt/C (10%)	DMF	None
3	Ir/CaCO <sub>3</sub> (5%)	EtOH	NEt <sub>3</sub> (5 %)
4	Ir/CaCO <sub>3</sub> (5%)	Hexane	None
5	Ir/CaCO <sub>3</sub> (5%)	DMF	None

# Materials design



# Why material informatics is difficult

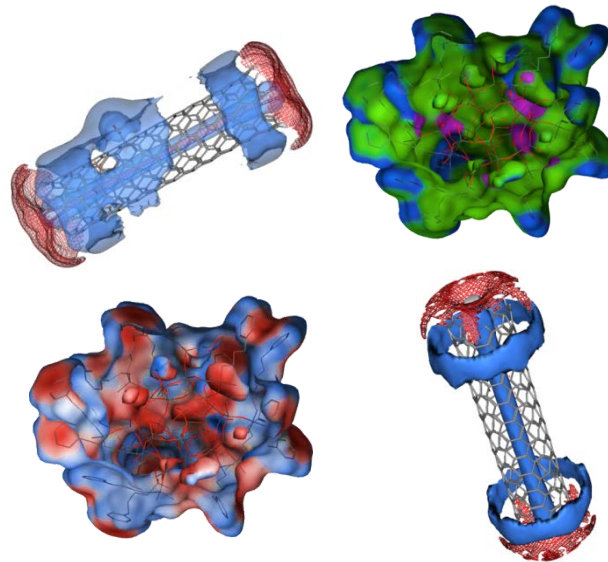
- **Materials data is more inconsistent** - different measurements and different conditions are used by different labs.
- **Materials data is more unpredictable** - many properties can not be calculated by higher level methods (DFT), especially for amorphous materials. MD is not appropriate for electronic properties.
- **Materials data is more process-dependent** - material properties depend not only on their chemical composition but also on how they are fabricated

Descriptors can be used to quantify the controllable parts of the factors involved, but not the effects of defect and impurities

# Material Quantitative Structure-Property Relationship (MQSPR)

$$\text{Material Property} = f(\text{descriptors})$$

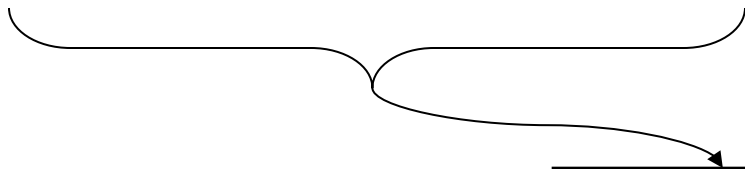
Molecular Descriptors  
Physicochemical Descriptors  
Macroscopic properties  
Morphology information



Descriptors

Model

Property



# Representing Potential-Energy Surfaces by Artificial Neural Networks

Eur. Phys. J. B (2014) 87: 152  
DOI: 10.1140/epjb/e2014-50070-0

THE EUROPEAN  
PHYSICAL JOURNAL B

Colloquium

## Next generation interatomic potentials for condensed systems

Christopher Michael Handley and Jörg Behler<sup>a</sup>

THE JOURNAL OF CHEMICAL PHYSICS 139, 054112 (2013)



## Permutation invariant polynomial neural network approach to fitting potential energy surfaces

Bin Jiang and Hua Guo<sup>a)</sup>

## Representing Potential Energy Surfaces with Neural Networks and High Dimensional Model Representations

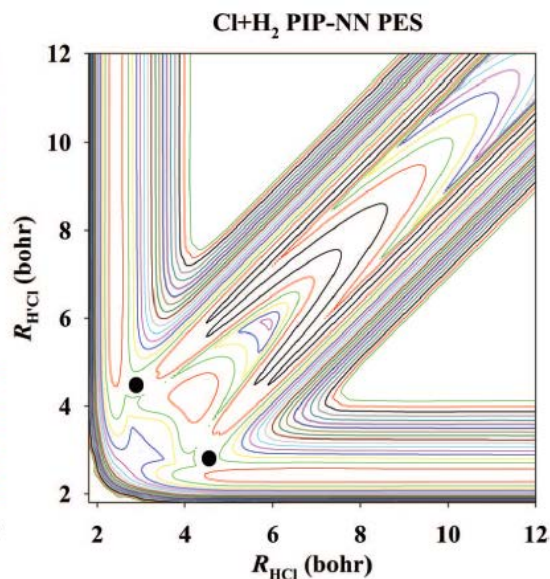
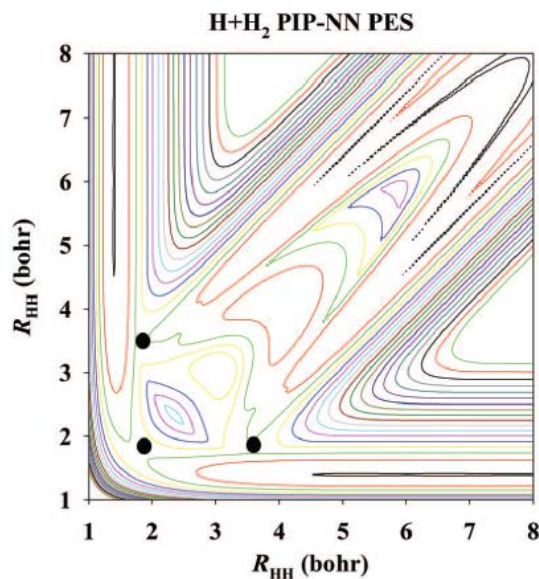
Sergei Manzhos and Tucker Carrington

THE JOURNAL OF CHEMICAL PHYSICS 122, 084104 (2005)

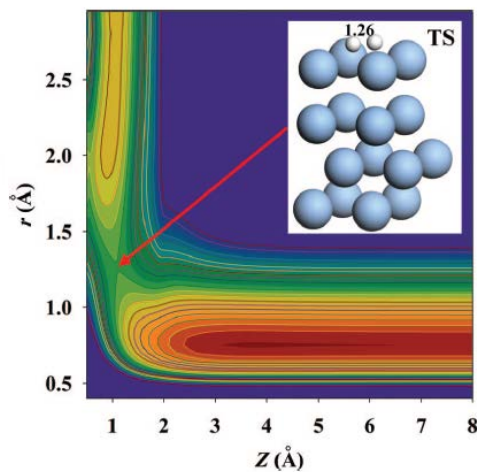
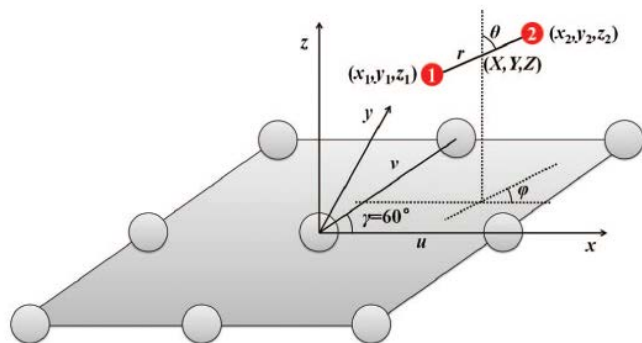
## *Ab initio* potential-energy surfaces for complex, multichannel systems using modified novelty sampling and feedforward neural networks

L. M. Raff

# Guo's Permutation-Invariant Neural Network for Approximating Potential-Energy Surfaces



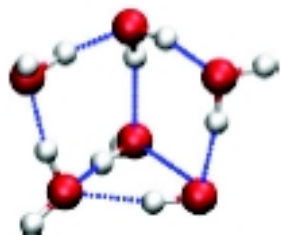
Three- and four-atom reaction systems (*J. Chem. Phys.*, **2013**, *139*, 054112; *J. Chem. Phys.*, **2013**, *139*, 204103; *J. Chem. Phys.*, **2014**, *140*, 044327)



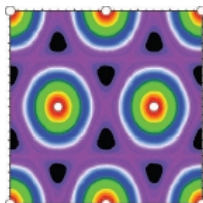
Molecule-surface interactions H<sub>2</sub> + Ag(111) and H<sub>2</sub> + Cu(111) (*J. Chem. Phys.*, **2014**, *141*, 034109)

# Behler's Permutation-Invariant Symmetry-Adapted Neural Network for Approximating Potential-Energy Surfaces.

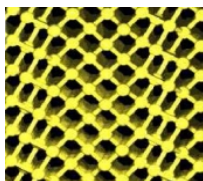
## Applications



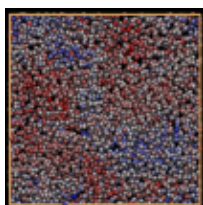
Water clusters (*J. Chem. Phys.*, **2012**, *136*, 064103; *J. Phys. Chem. A*, **2013**, *117*, 7356; *Z. Phys. Chem.*, **2013**, *227*, 1559)



Molecule-surface interactions (*J. Chem. Phys.*, **2007**, *127*, 014705)



Metadynamics simulations of the high-pressure phases of silicon (*Phys. Rev. Lett.*, **2008**, *100*, 185501)

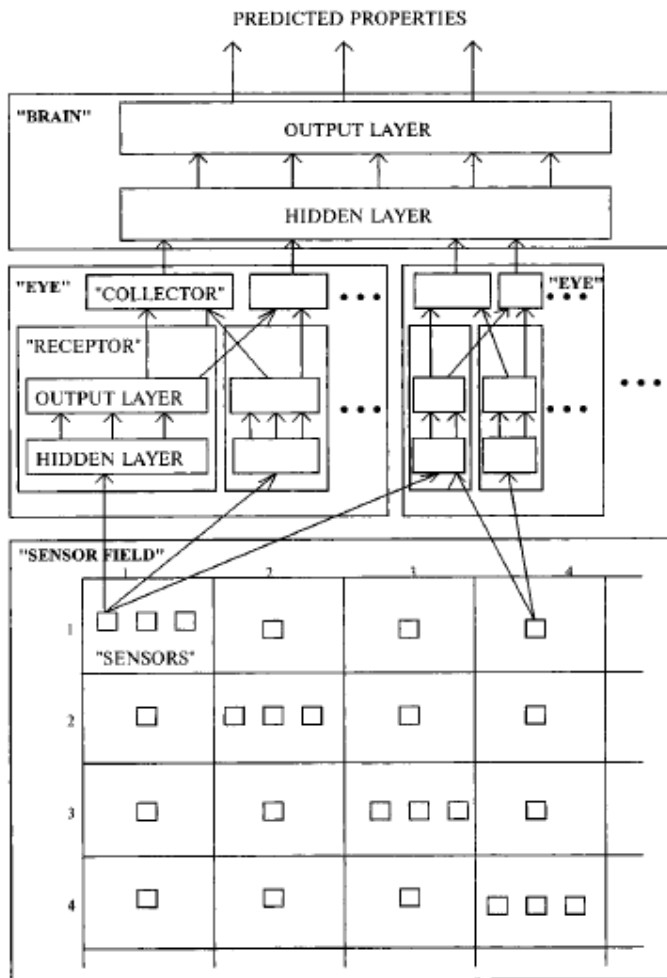


Large-scale MD simulations of phase change materials GeTe (*J. Phys. Chem. Lett.*, **2013**, *4*, 4241)



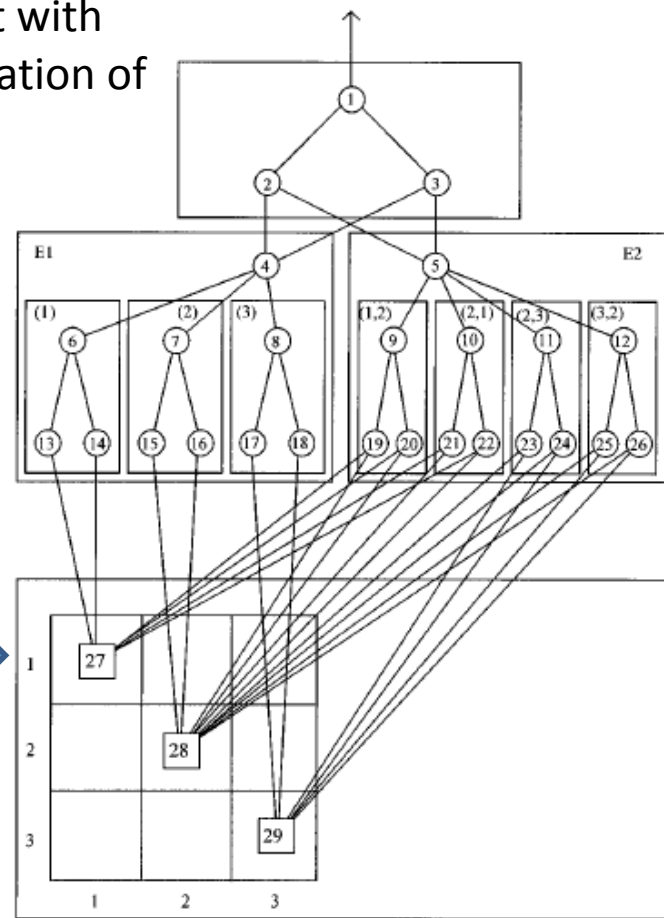
# The First Permutation-Invariant Neural Network

Initially developed for chemoinformatics purposes (QSPR modeling)



**General architecture**

Output is invariant with respect to permutation of identical atoms



Properties of atoms (on diagonal)

Distances or bond orders (off-diagonal)

**Application to 3-atom system**

### Neural Network Potentials:

- now applicable to condensed systems (solids, large clusters, liquids)
- NN reproduces total energies very accurately (also in value!)
- provide analytic derivatives (forces, stress)
- fast and linear scaling with system size
- can be constructed using any electronic structure method

### Limitations of NN Potentials:

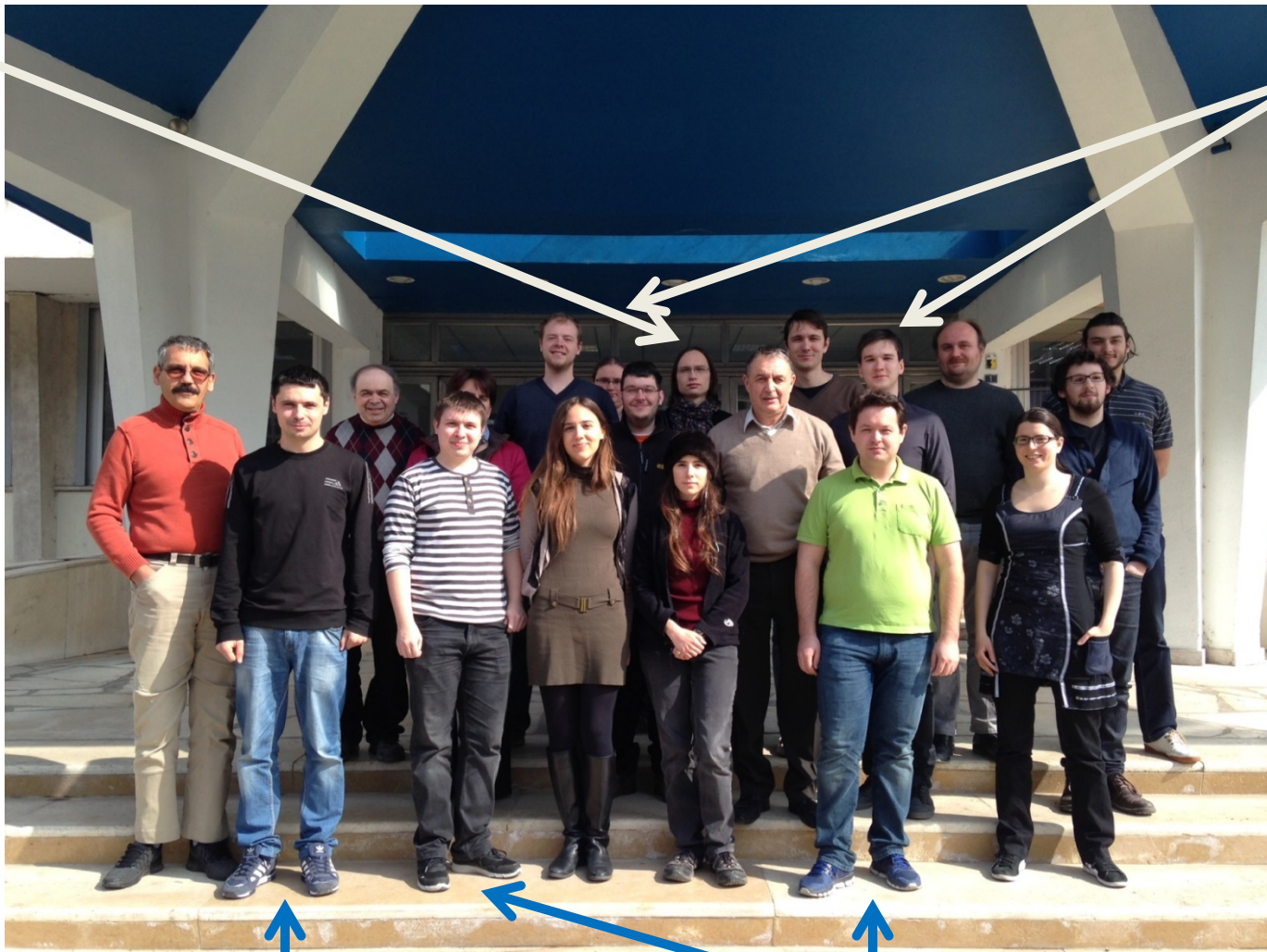
- need many reference calculations
- limited transferability, no extrapolation
- limited number of chemical elements

### Outlook:

- useful tool to speed up extended simulations

Moscow

Odessa



Kiev

Kazan