

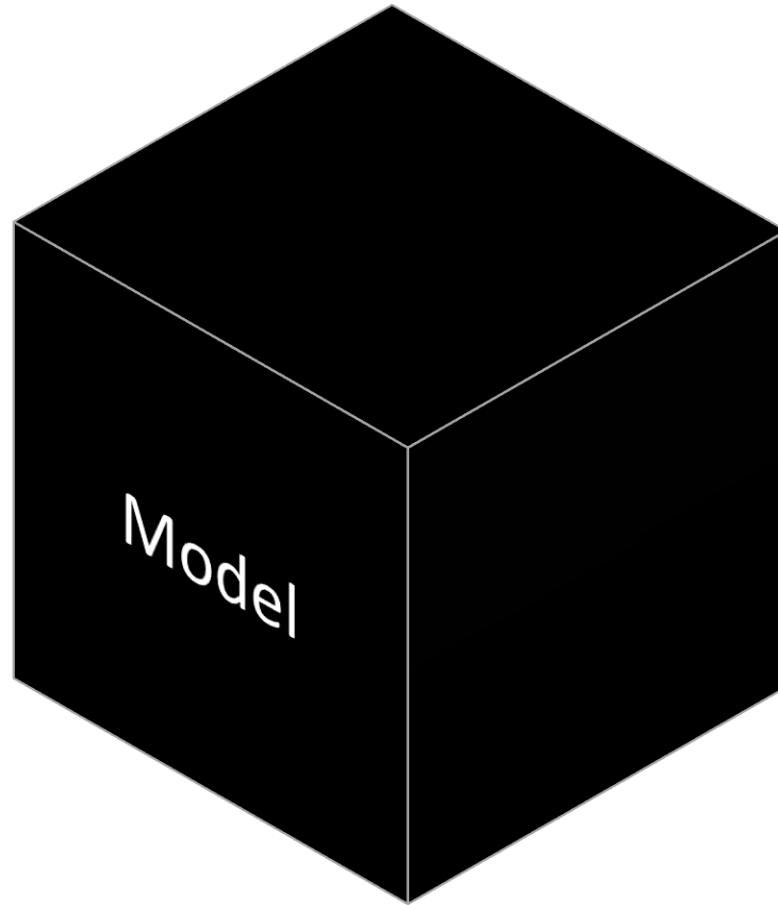


Explainable artificial intelligence: evolution, achievements and perspectives

Pavel Polishchuk

Institute of Molecular and Translational Medicine
Faculty of Medicine and Dentistry
Palacky University

pavlo.polishchuk@upol.cz



plant growth inhibition activity of
phenoxyacetic acids

$$1/C = 4.08\pi - 2.14\pi^2 + 2.78\sigma + 3.38$$

Hansch equation

rate of penetration of membranes
in the plant cell

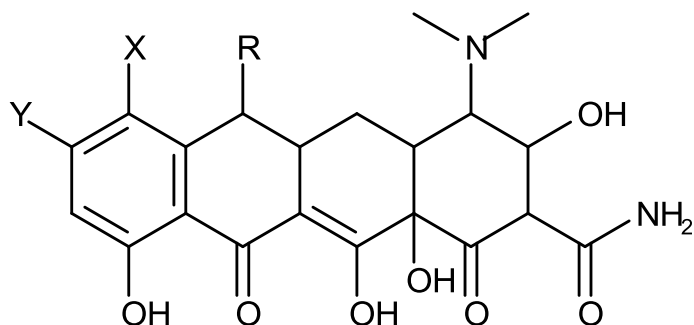
electronic factors

$$\pi = \log P_X - \log P_H$$

σ - Hammett constant

Hansch, C.; Fujita, T., ρ - σ - π Analysis. A Method for the
Correlation of Biological Activity and Chemical Structure.
Journal of American Chemical Society **1964**, 86, 1616-1626.

Free-Wilson models

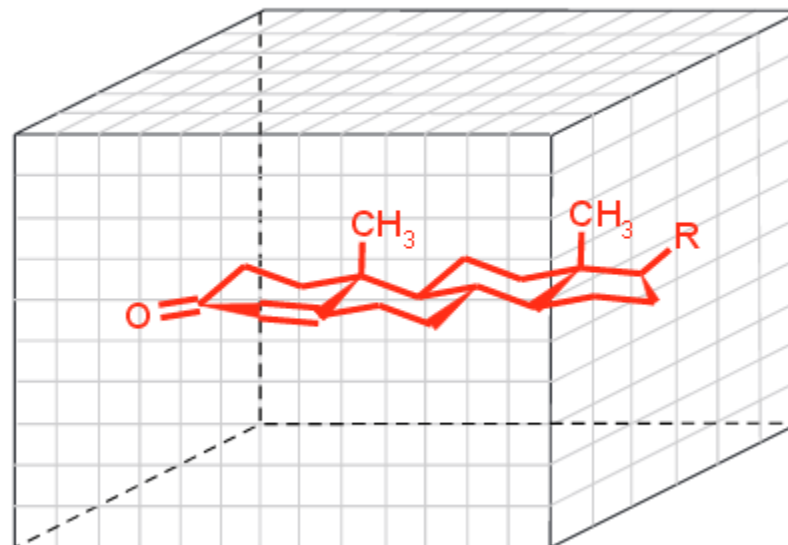
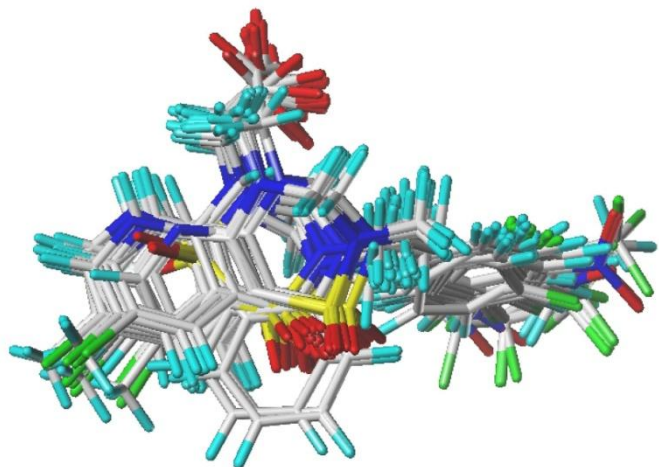


Inhibition activity of compounds
against *Staphylococcus aureus*

R is H or CH₃;
X is Br, Cl, NO₂ and
Y is NO₂, NH₂, NHC(=O)CH₃

$$\text{Act} = 75R_H - 112R_{\text{CH}_3} + 84X_{\text{Cl}} - 16X_{\text{Br}} - 26X_{\text{NO}_2} + 123Y_{\text{NH}_2} + 18Y_{\text{NHC(=O)CH}_3} - 218Y_{\text{NO}_2}$$

CoMFA: Comparative molecular field analysis



probe: electrostatic (H^+) steric (Csp^3)

Mol	Electrostatic field descriptors						Steric field descriptors							
	(0,0,0)	(0,0,1)	(G,G,G)	(0,0,0)	(0,0,1)	(G,G,G)
1														
2														
...														
N														

PLS model

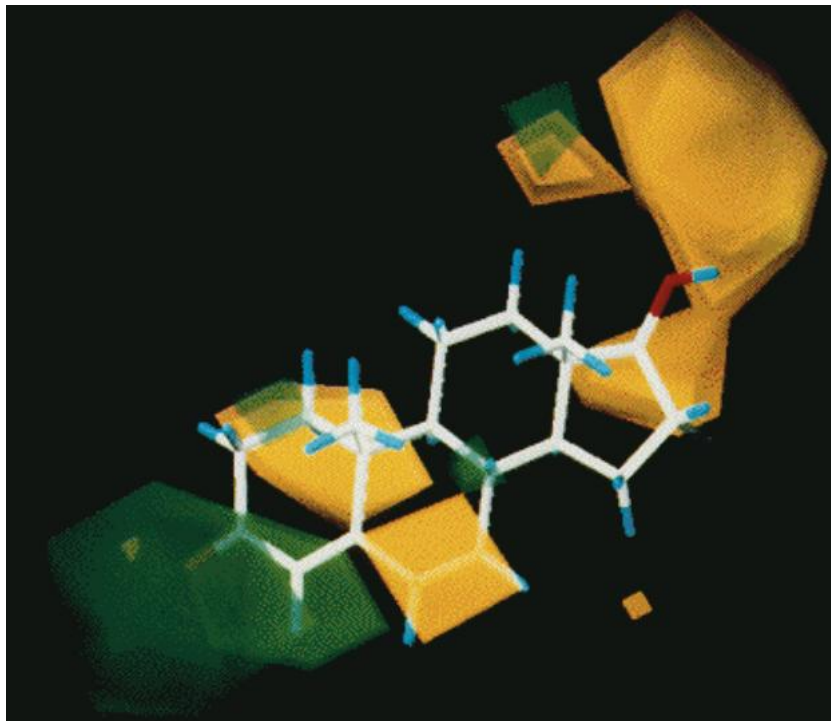
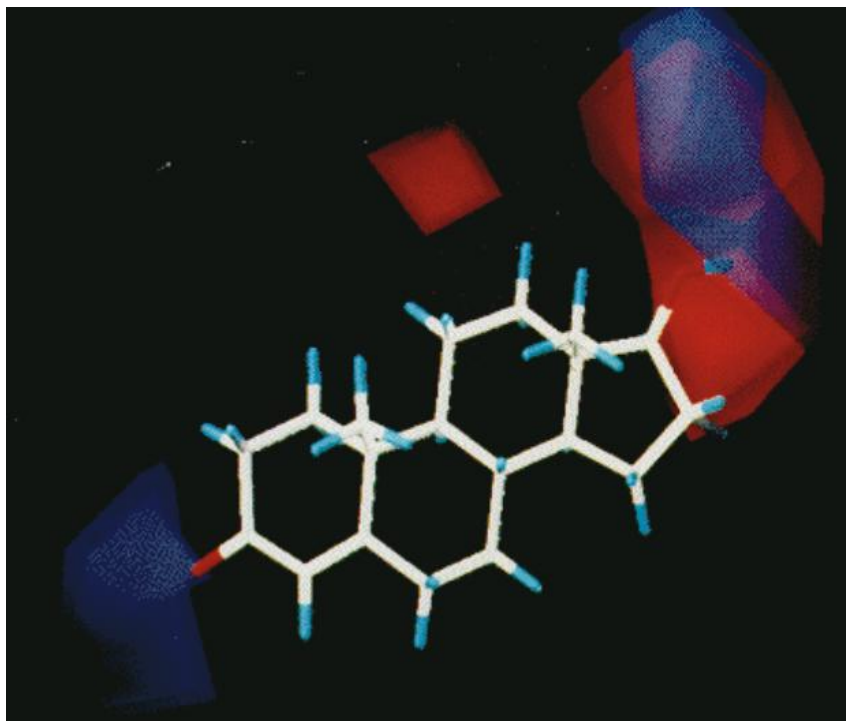


Y
2.4
3.7
...
8.1

CoMFA: Comparative molecular field analysis

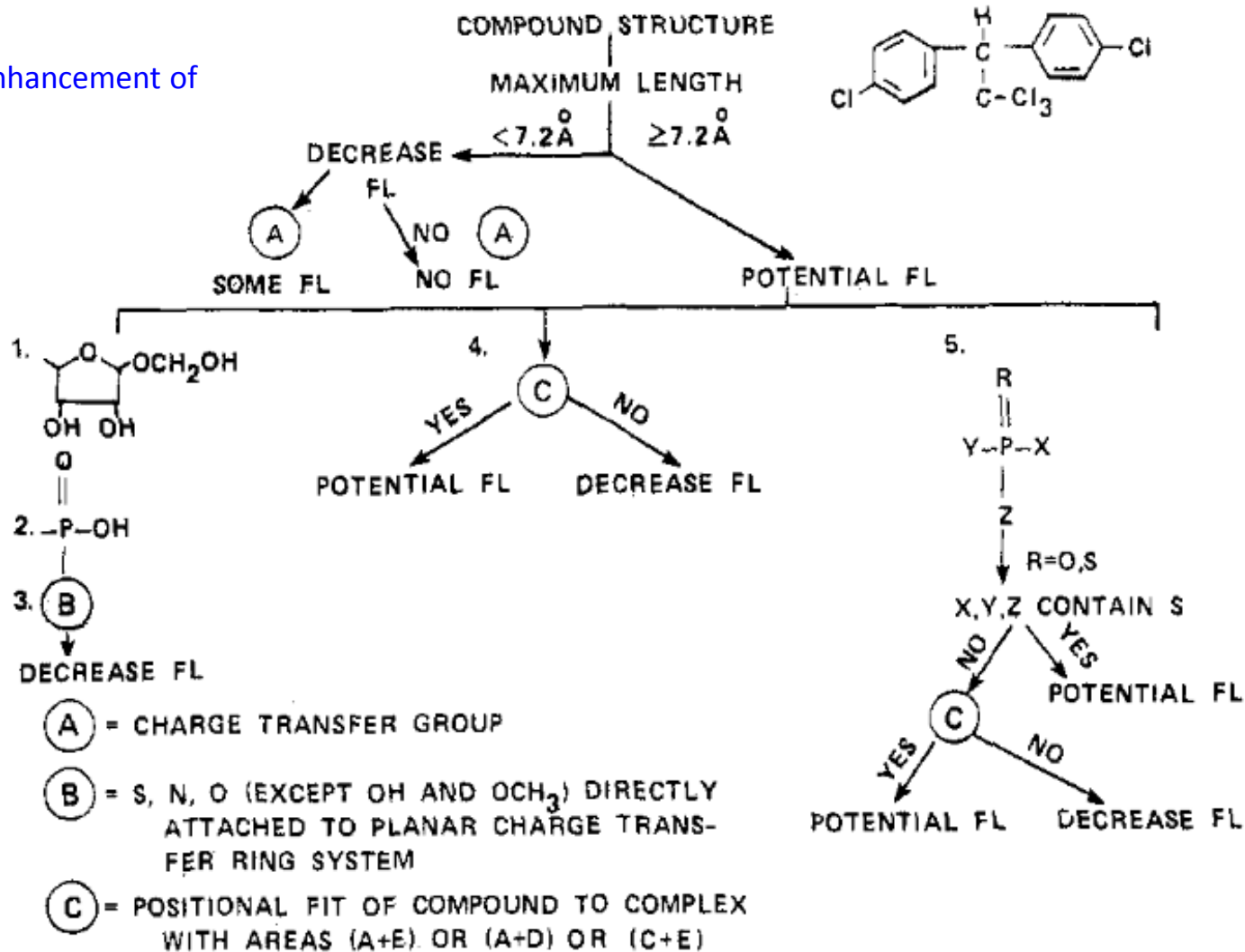
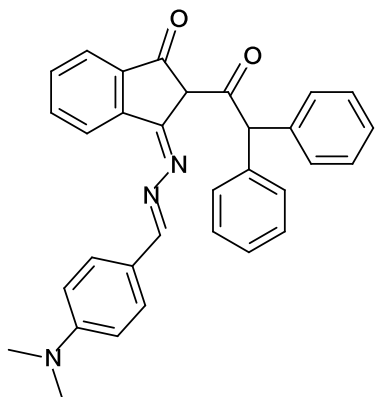
$$Y = \sum_{i=1}^n b_i x_i + c$$

b – contribution of steric or electrostatic field in a particular cell

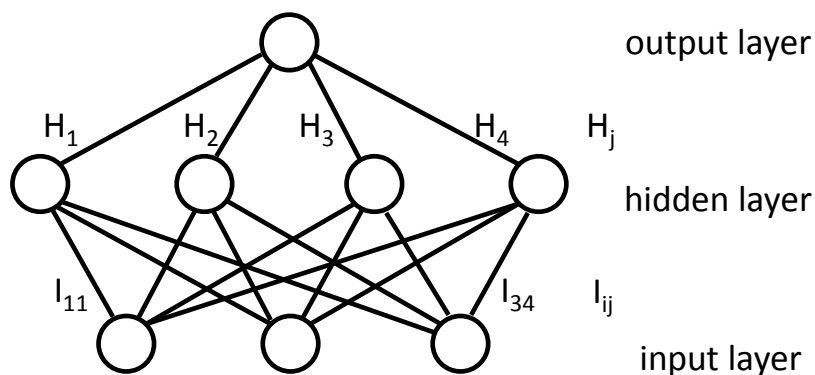


Decision tree

solid-phase fluorescence enhancement of



Neural networks: relative importance



$$P_{ij} = |I_{ij}| \times |H_j|$$

$$Q_{ij} = \frac{P_{ij}}{\sum_i P_{ij}}$$

$$S_i = \sum_j Q_{ij}$$

$$\text{relative importance}_i = \frac{S_i}{\sum_i S_i}$$

Garson, G. D. Interpreting neural-network connection weights. *AI Expert* **1991**, *6*, 46-51

adsorbability of 55 organic compounds on activated carbon fibers

$$\log K = 3.33 - 1.55 {}^3\chi^v + 0.58 {}^5\chi^v + 3.52 {}^6\chi^v - 1.42 {}^3\chi_c + 2.29 {}^4\chi_{pc}^v$$

$$n = 49, R^2_{adj} = 0.648, SE = 0.199$$

	${}^3\chi^v$	${}^5\chi^v$	${}^6\chi^v$	${}^3\chi_c$	${}^4\chi_{pc}^v$
relative importance, %	20.3	17.3	34.4	11.6	16.4
influence on logK	↓	↑	↑	↓	↑

${}^3\chi^v$ is related to bulky and branched molecules

${}^5\chi^v$ is related to heteroatomic contents

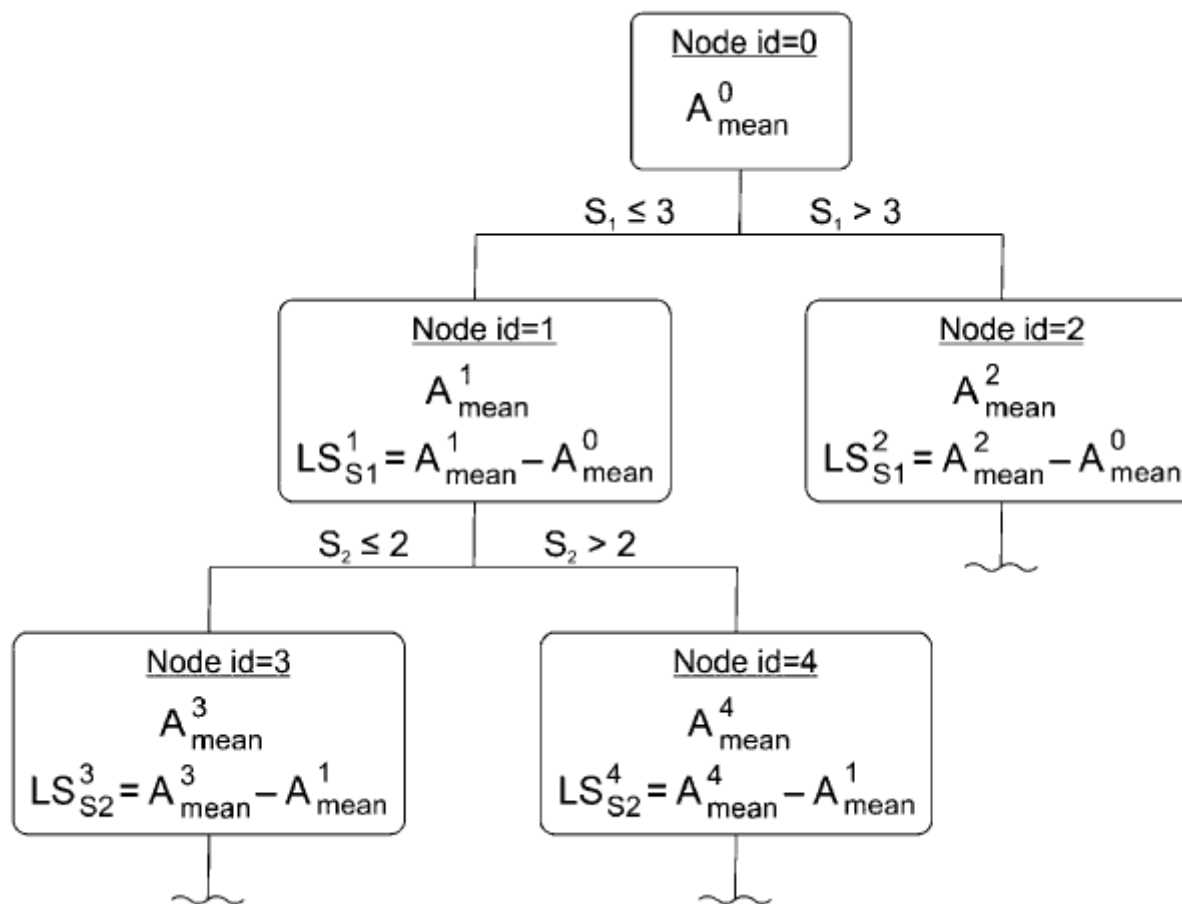
${}^6\chi^v$ is related to highly branched compounds, like atrazin

${}^3\chi_c$ is related to highly substituted compounds comprising *tert*-butyl groups or having more than three substituents

${}^4\chi_{pc}^v$ is related to compounds with more than four substituents

Brasquet, C.; Bourges, B.; Le Cloirec, P. Quantitative Structure–Property Relationship (QSPR) for the Adsorption of Organic Compounds onto Activated Carbon Cloth: Comparison between Multiple Linear Regression and Neural Network. *Environ. Sci. Technol.* **1999**, *33*, 4226-4231.

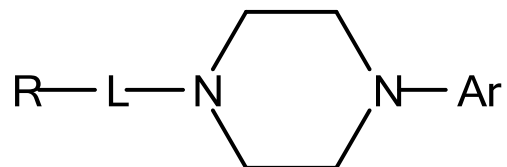
Random Forest: descriptor contributions



347 agonists of 5-HT_{1A} receptor

PLS, 72 descriptors, $R^2_{5CV} = 0.64$

RF, 2500 descriptors, $R^2_{OOB} = 0.70$

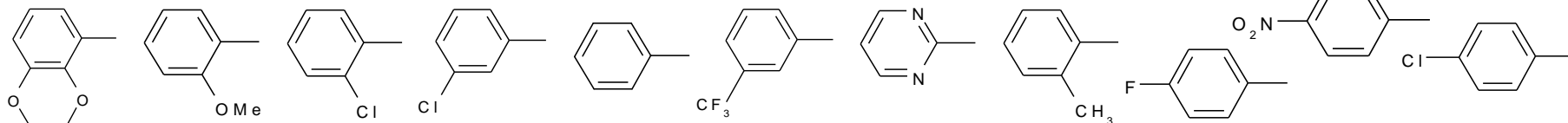


Ar - substituted (hetero)aryls

L - polymethylene chain

R - various (poly)cyclic residues

Ar



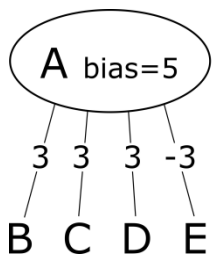
PLS	0.84	0.18	0.03	-0.04	-0.06	-0.09	-0.11	-0.66	-0.73	-0.94	-0.96
RF	0.27	0.24	0.04	0.07	-0.02	0.11	0.04	-0.04	-0.55	-0.66	-0.66

L	-(CH ₂) ₆ -	-(CH ₂) ₅ -	-(CH ₂) ₄ -	-(CH ₂) ₃ -	-(CH ₂) ₂ -	-CH ₂ -
PLS	0.8	0.71	0.81	0.08	-0.04	0.06
RF	0.14	0.19	0.14	-0.01	-0.03	0.05

Rule extraction approaches

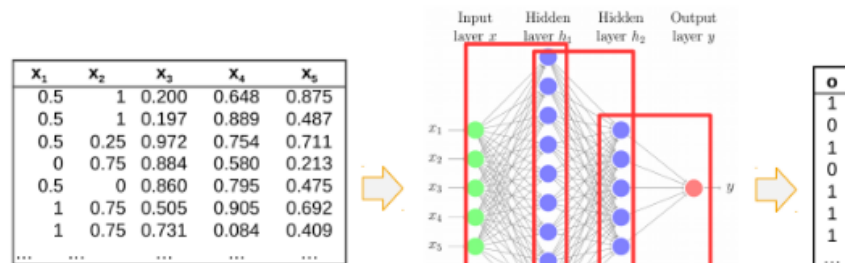
Decompositional (use knowledge about internal structure of a model, e.g. NN)

DeepRED



if B, C and not(E), then A
 if B, D and not(E), then A
 if C, D and not(E), then A
 if B, C, D, then A

Fu, L., Rule learning by searching on adapted nets. In *Proceedings of the ninth National conference on Artificial intelligence - Volume 2*, AAAI Press: Anaheim, California, **1991**; pp 590-595.



```

IF x1>0.5 AND x2>0.6 THEN h11<=0.4
IF x1>0.5 AND x2<=0.6 THEN h11>0.4
IF x1<=0.5 ...
...

IF h12>0.4 AND h110<=0.1 THEN h23<=0.5
IF h12>0.4 AND h110>0.1 THEN h24>0.3
IF h12<=0.4 AND h11<=0.4 THEN h21>0.6
IF h12<=0.4 AND h11 >0.1 THEN h21<=0.6

IF h21>0.6 AND h24>0.3 THEN o=0
IF h21>0.6 AND h24<=0.3 THEN o=1
IF h21<=0.6 THEN o=1
    
```

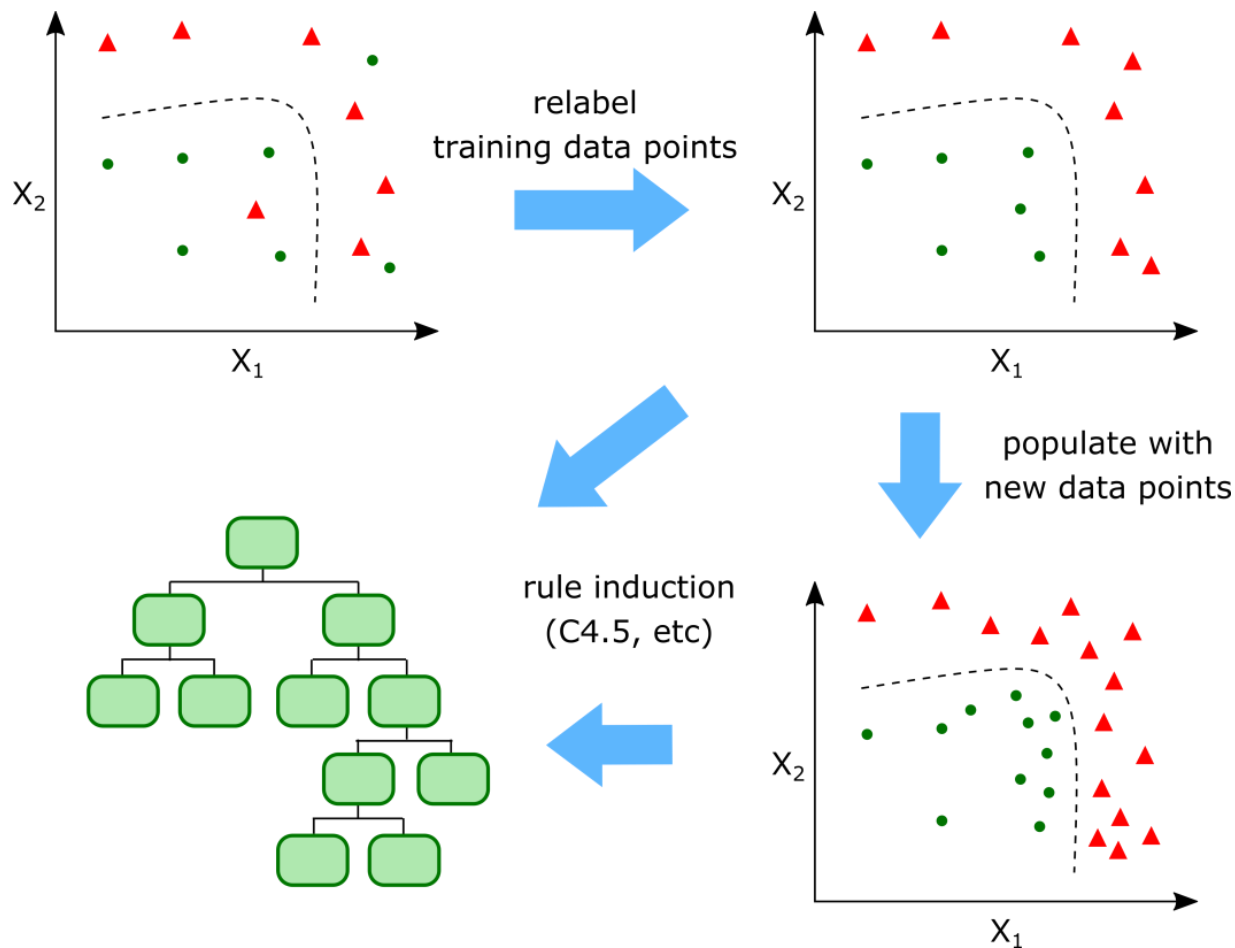
```

IF x1<0.5 AND x2>0.75 THEN o=1
IF x1>0.9 THEN o=1
IF x1>0.5 AND x1<0.9 AND x3>0.2 THEN o=1
IF x2>0.2 AND x3<0.5 AND x5<0.5 THEN o=1
IF x2>0.4 AND x3<0.7 THEN o=1
IF x2<0.2 THEN o=1
IF x4>0.8 THEN o=1
IF x3<0.7 AND x3>0.2 AND x4<0.3 THEN o=1
    
```

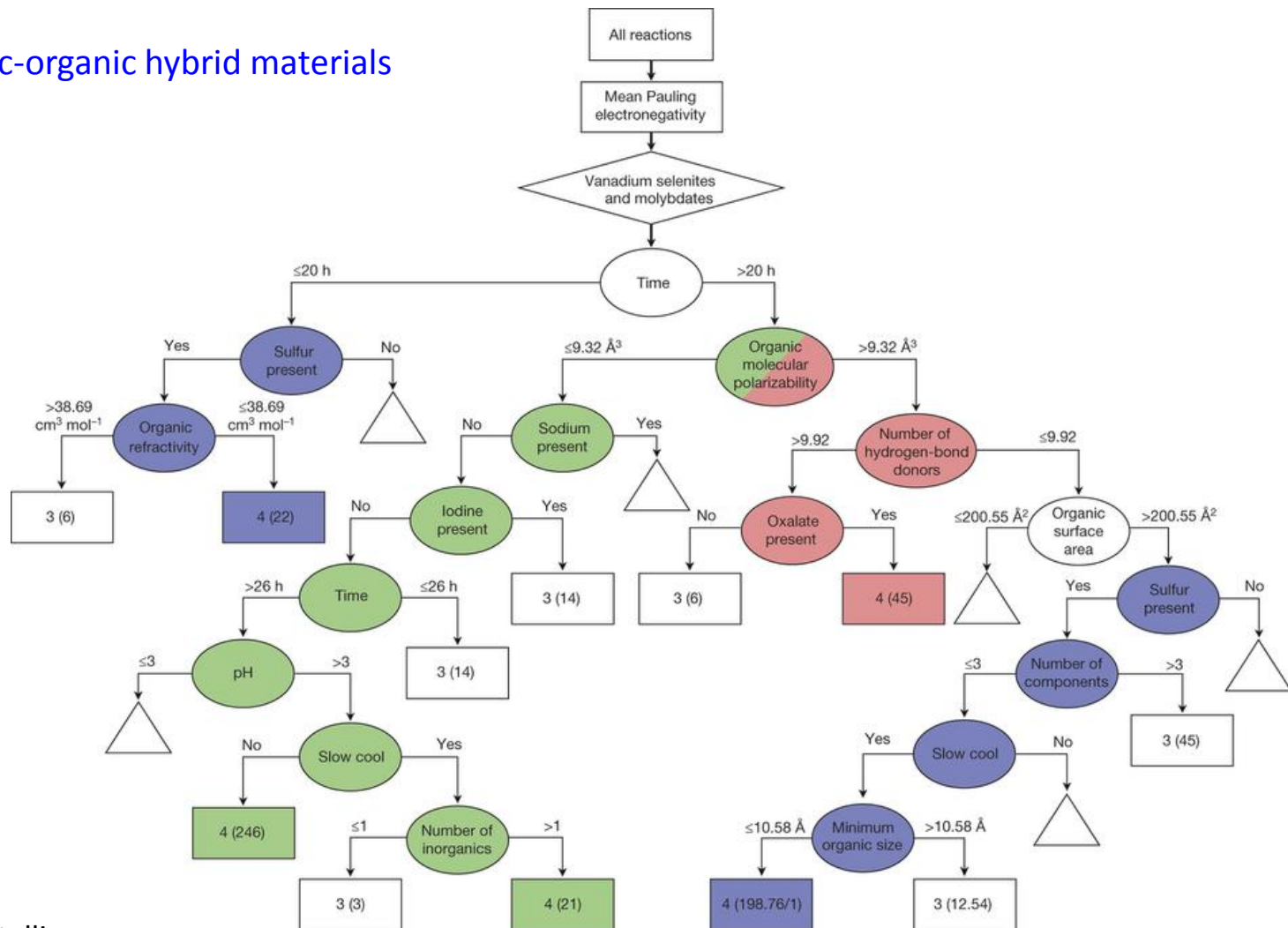
Zilke, J. R.; Loza Mencía, E.; Janssen, F. DeepRED – Rule Extraction from Deep Neural Networks. In *Discovery Science: 19th International Conference, DS 2016, Bari, Italy, October 19–21, 2016, Proceedings*, Calders, T.; Ceci, M.; Malerba, D., Eds. Springer International Publishing: Cham, **2016**; pp 457-473.

Rule extraction approaches

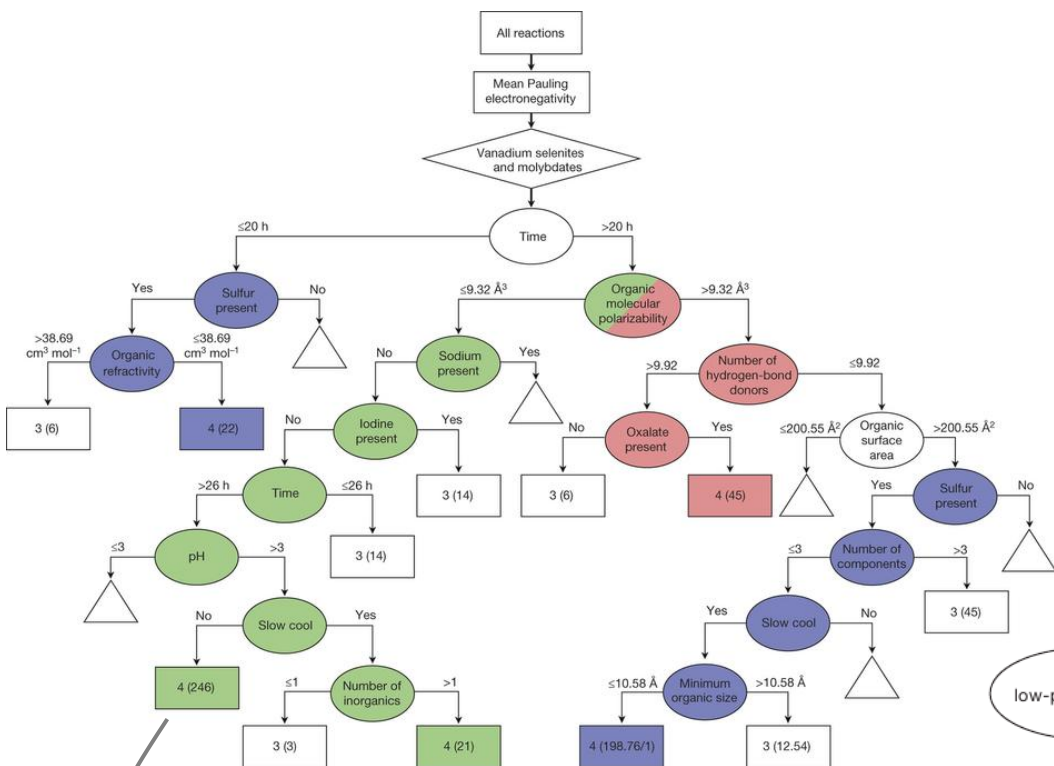
Pedagogical / surrogate modeling (treat a model as a black box)



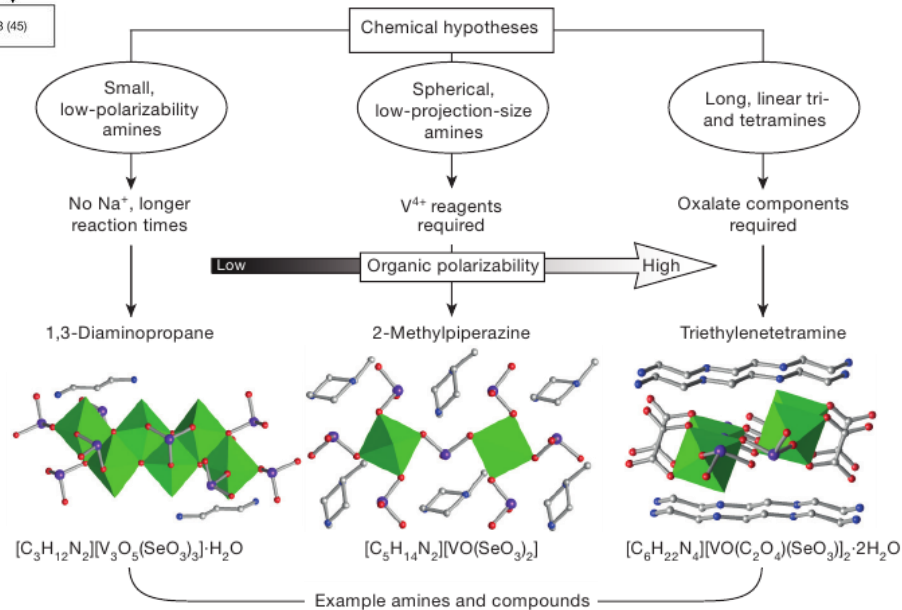
inorganic-organic hybrid materials



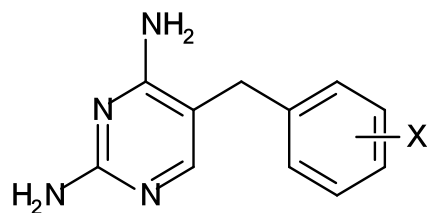
output:
 3 – polycrystalline
 4 – single-crystal



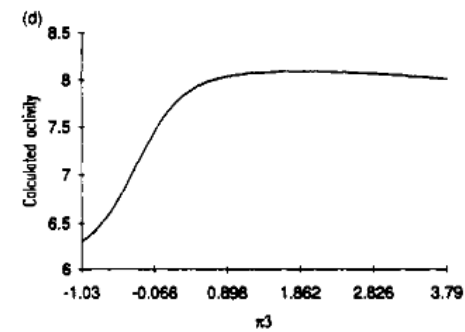
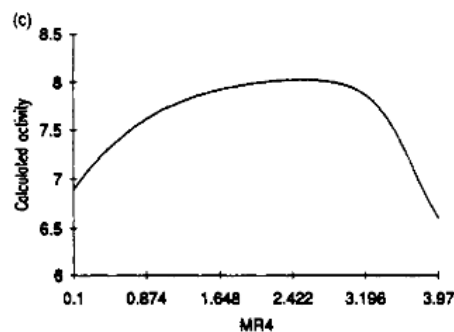
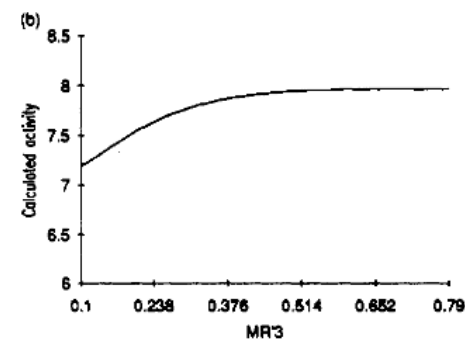
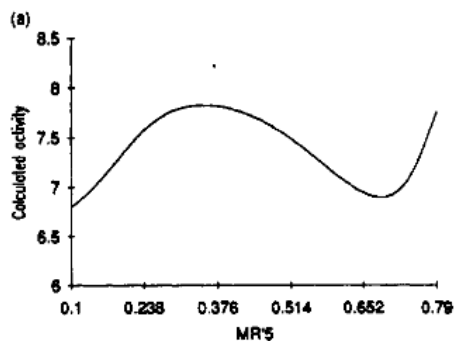
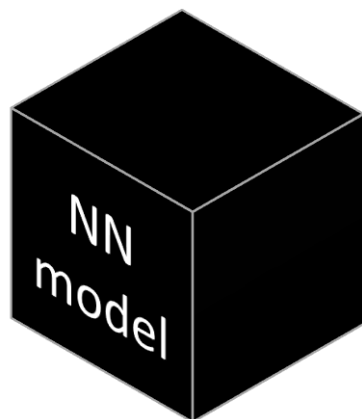
output:
3 – polycrystalline
4 – single-crystal



Sensitivity analysis



DHFR inhibitors



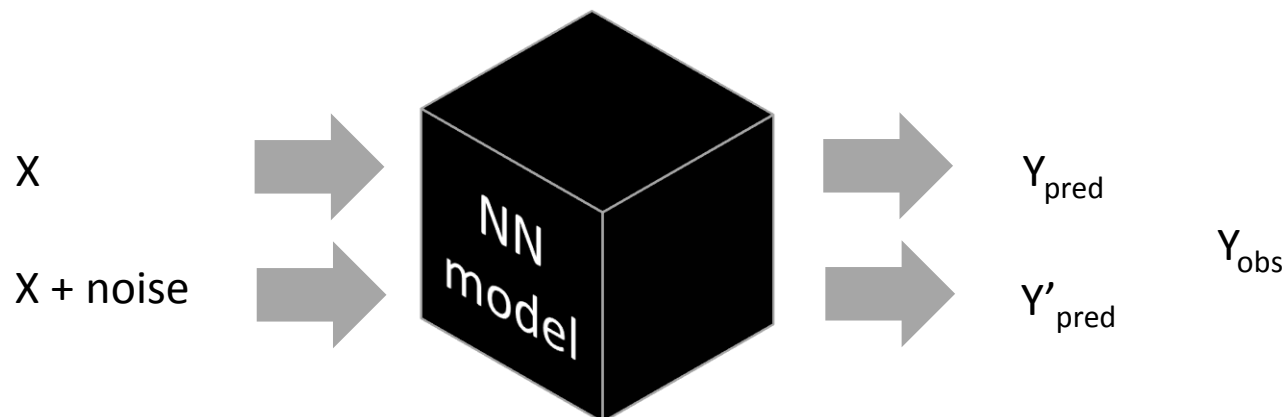
$$\log(1/K) = 0.95MR_5 + 0.89MR_3 + 0.80MR_4 - 0.21MR_4^2 + 1.58\pi_3 - 1.77\log(\beta \times 10^{\pi_3} + 1) + 6.65$$

$$RMSE = 0.093$$

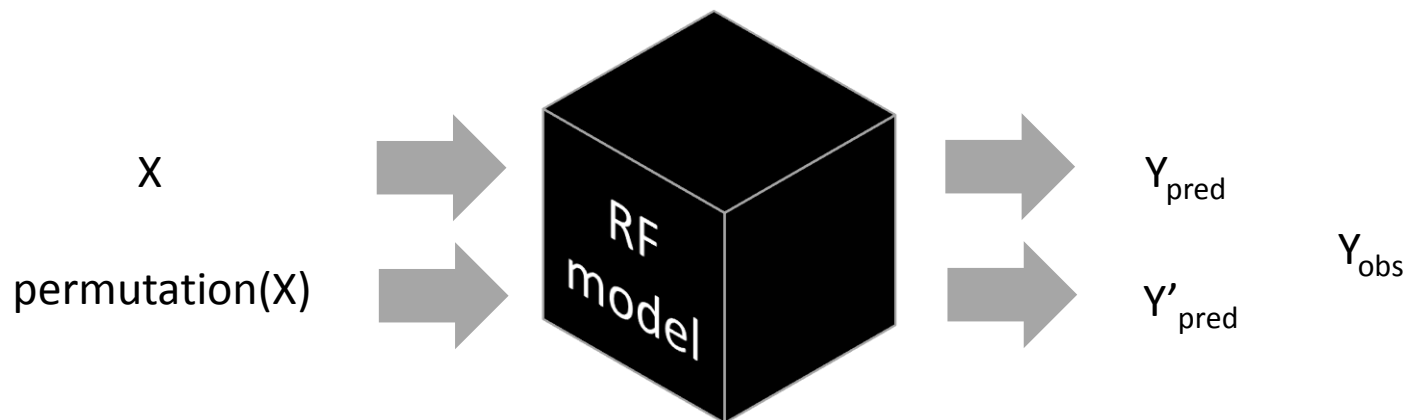
$$\log(1/K) = 11.79MR_5^3 - 15.74MR_5^2 + 6.55MR_5 + 0.89MR_3 + 0.80MR_4 - 0.21MR_4^2 + 1.58\pi_3 - 1.77\log(\beta \times 10^{\pi_3} + 1) + 6.24$$

$$RMSE = 0.074$$

Variable importance

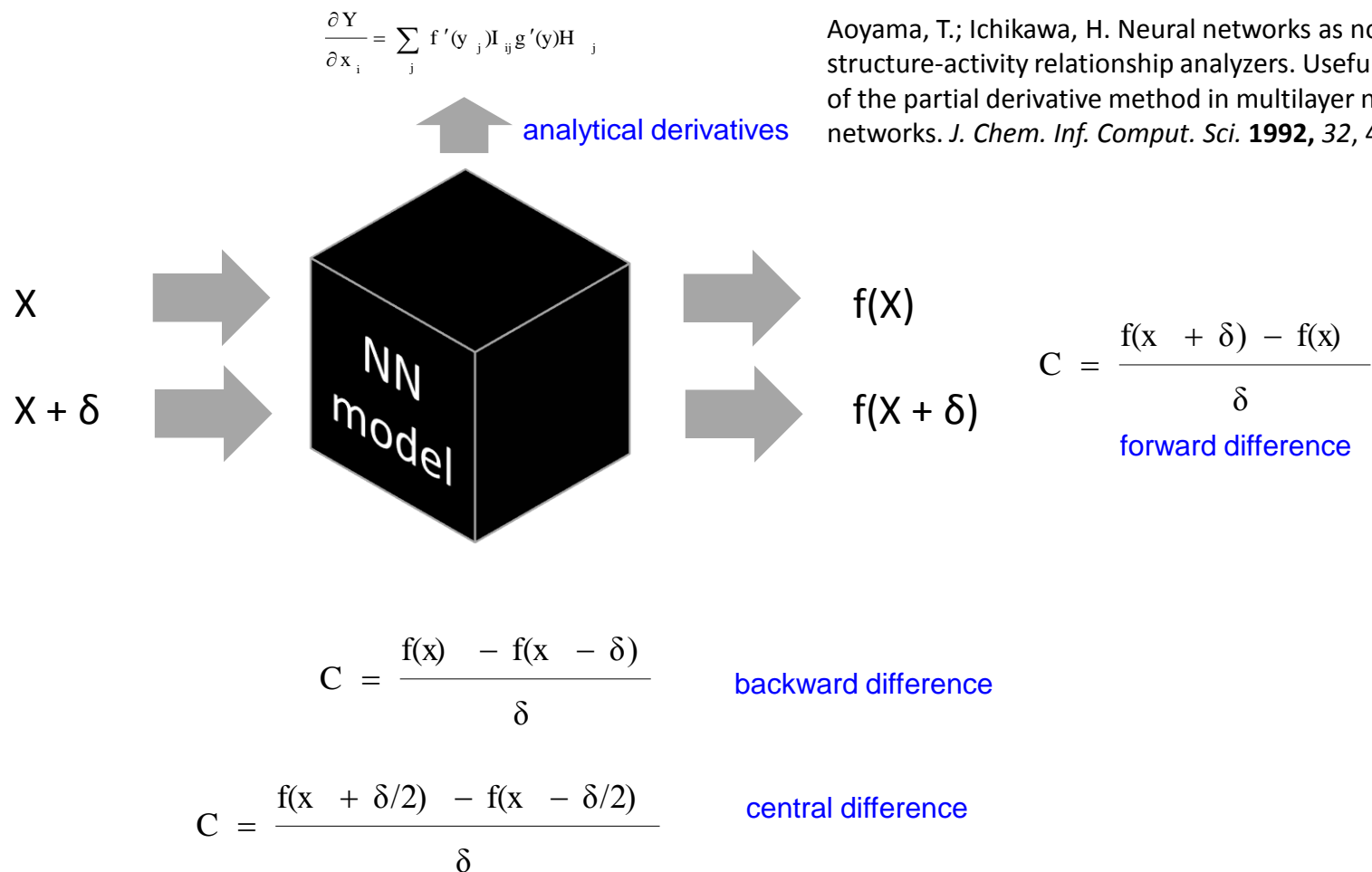


Györgyi, G. Inference of a rule by a neural network with thermal noise. *Phys. Rev. Lett.* **1990**, *64*, 2957-2960.



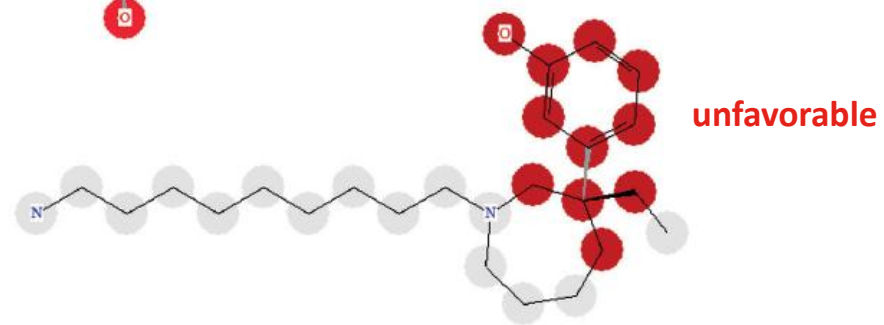
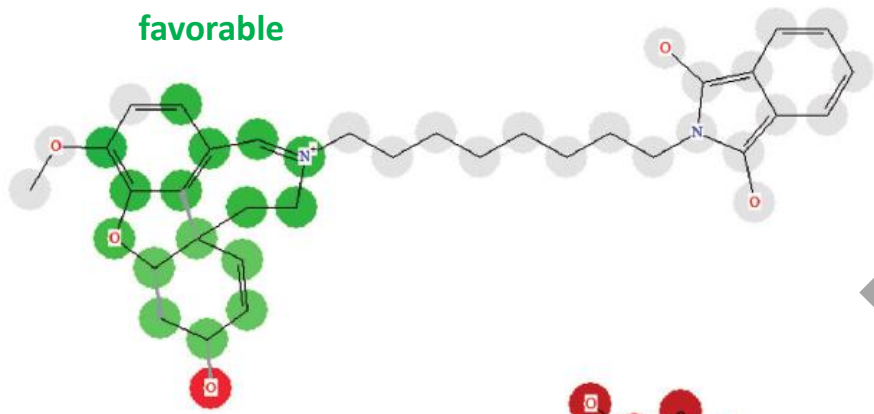
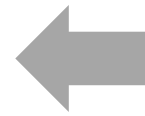
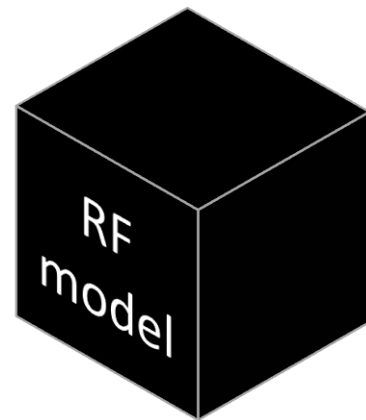
Breiman, L., Random Forests. *Machine Learning* **2001**, *45*, 5-32.

Partial derivatives



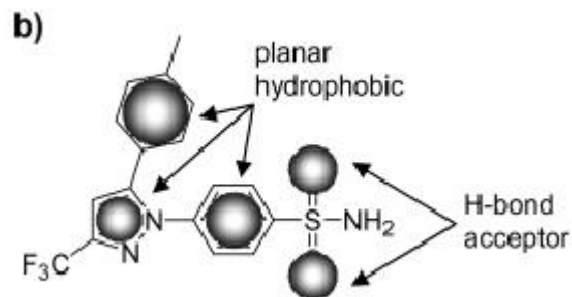
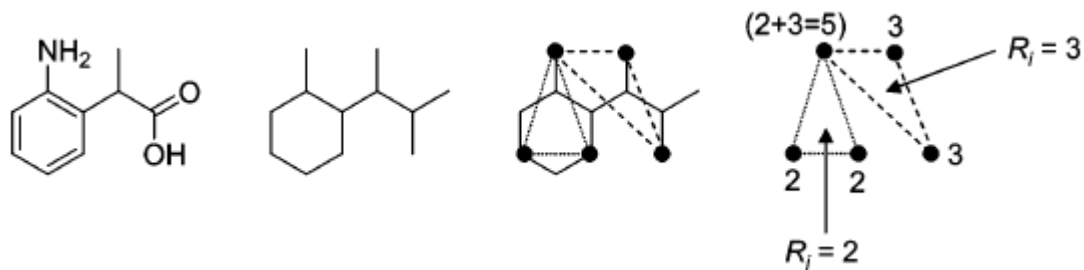
Aoyama, T.; Ichikawa, H. Neural networks as nonlinear structure-activity relationship analyzers. Useful functions of the partial derivative method in multilayer neural networks. *J. Chem. Inf. Comput. Sci.* **1992**, 32, 492-500.

AChE inhibitors  fragmental descriptors

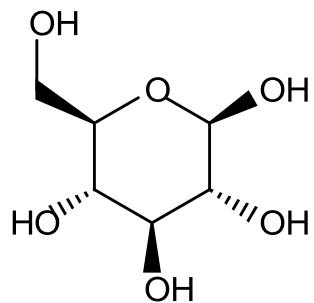


COX2 inhibitors

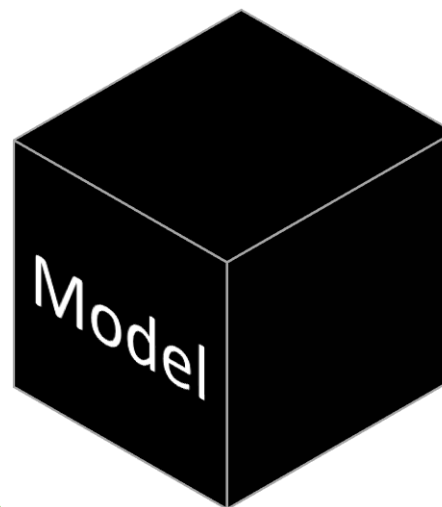
$$R_i = f(\mathbf{x}(F_i = 1)) - f(\mathbf{x}(F_i = 0))$$



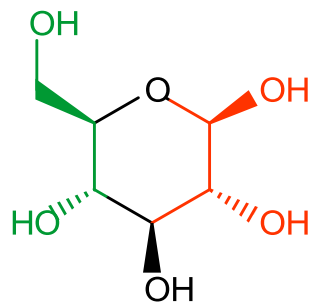
Learning



D ₁	D ₂	D ₃	...	D _N
1	0	9	...	1
4	0	1	...	1
0	2	3	...	3
...
4	0	0	...	1



Output



(optional)

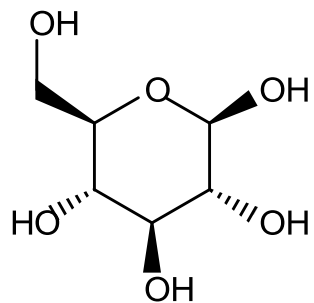
Descriptor contributions



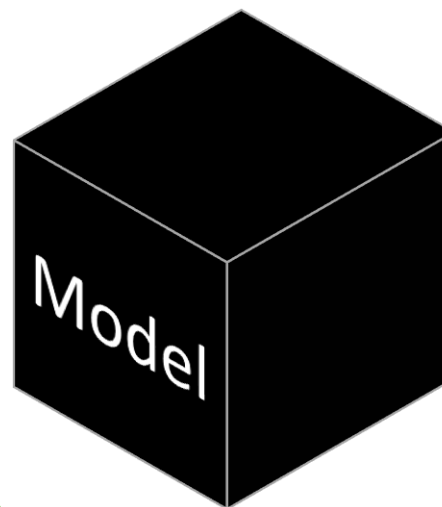
Interpretation

“model → descriptors → structure” paradigm

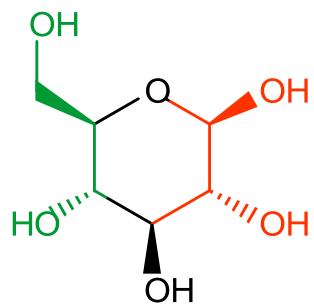
Learning



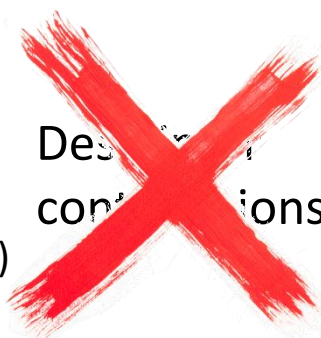
D ₁	D ₂	D ₃	...	D _N
1	0	9	...	1
4	0	1	...	1
0	2	3	...	3
...
4	0	0	...	1



Output



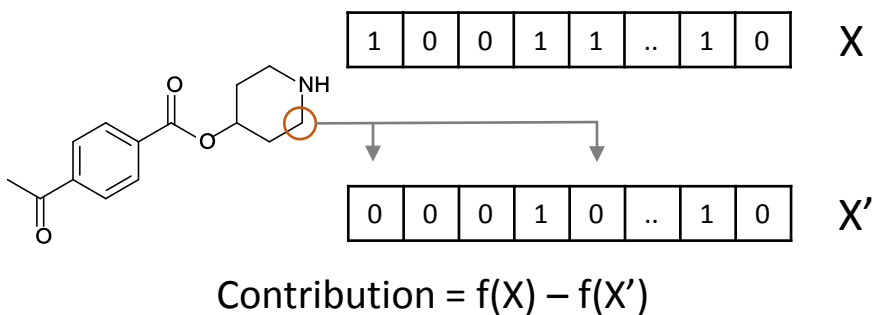
(optional)



Interpretation

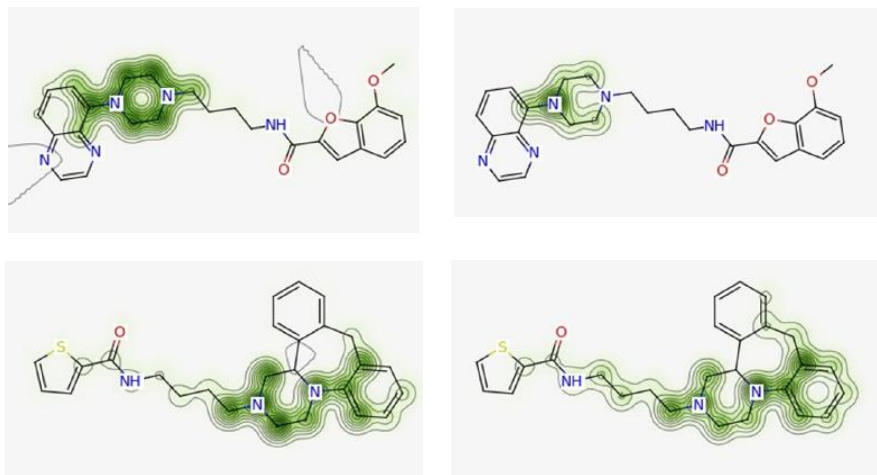
“model → structure” paradigm

Similarity maps



RF

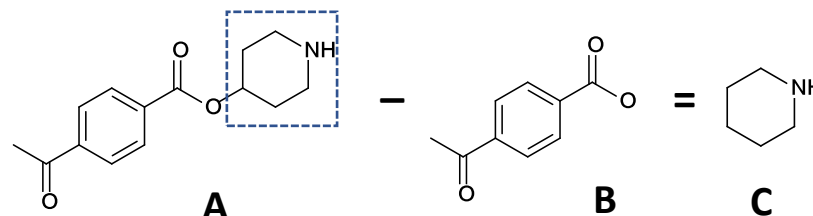
NB



dopamine D3 ligands

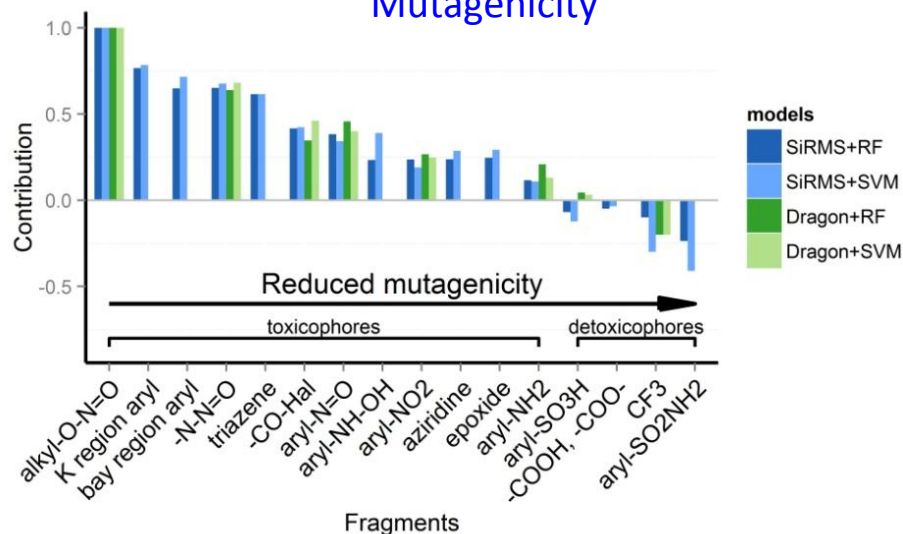
Riniker, S.; Landrum, G. Similarity maps - a visualization strategy for molecular fingerprints and machine-learning methods. *J. Cheminf.* **2013**, *5*, 43.

Universal interpretation



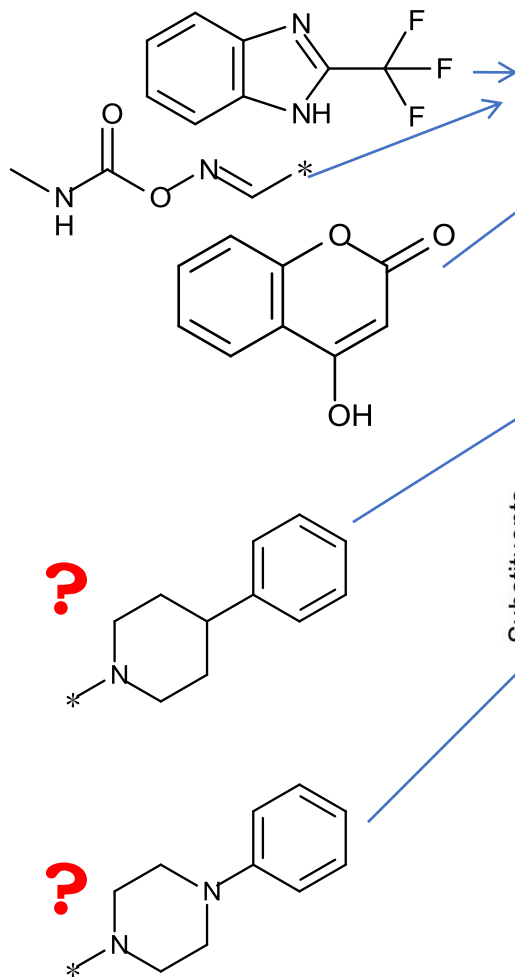
Activity _{pred} (A)	Activity _{pred} (B)	Contribution(C)
$f(A)$	$f(B)$	$W(C) = f(A) - f(B)$

Mutagenicity



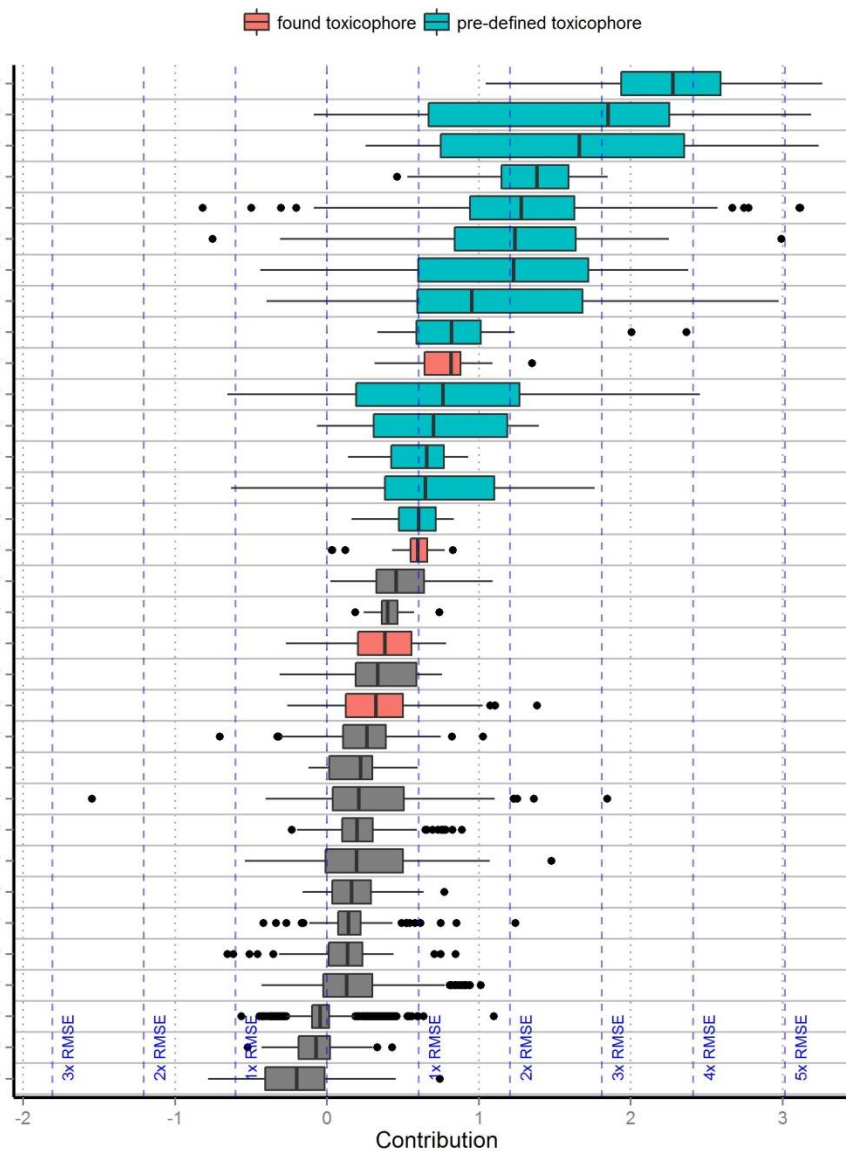
Polishchuk, P. G.; Kuz'min, V. E.; Artemenko, A. G.; Muratov, E. N. Universal Approach for Structural Interpretation of QSAR/QSPR Models. *Molecular Informatics* **2013**, *32*, 843-853
<https://github.com/DrrDom/spci>

Acute oral toxicity on rats



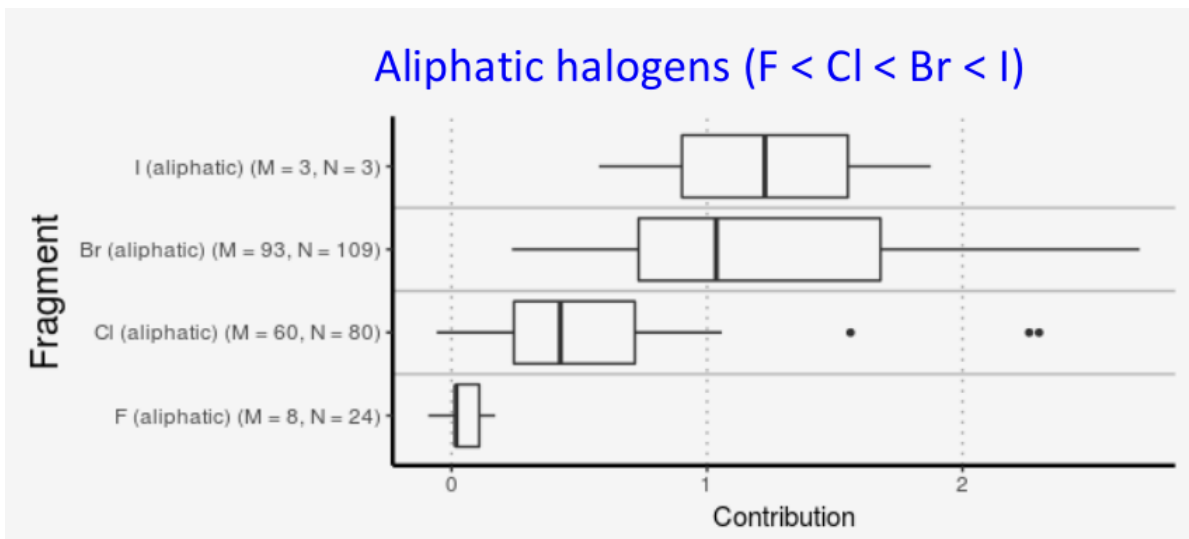
Substituents

- 2-trifluoromethylbenzimidazole (M = 128, N = 128)***
- O-(methylaminocarbonyl)oxime (M = 91, N = 101)***
- 4-hydroxycoumarin (M = 14, N = 16)***
- dinitrophenoxy (M = 22, N = 24)***
- phosphorodithioate (M = 172, N = 180)***
- phosphorothionate (M = 142, N = 148)***
- phosphoryl (M = 127, N = 133)***
- hexachlorononbornene (M = 22, N = 23)***
- 1,3-indandione (M = 16, N = 16)***
- 4-phenylpiperidine (M = 15, N = 16)***
- carbamate (M = 333, N = 358)***
- NHNH2 (not hydrazide) (M = 11, N = 12)**
- 2-fluoroacetyl (M = 17, N = 18)***
- thiourea (M = 45, N = 51)***
- 2-(2,4-dichlorophenoxy)acetyl (M = 14, N = 14)***
- phenylpiperazine (M = 32, N = 32)***
- nitrosamine (M = 119, N = 126)***
- pyrrole (M = 14, N = 14)***
- piperazine (M = 100, N = 100)***
- aziridine (M = 17, N = 39)***
- piperidine (M = 114, N = 128)***
- furan (M = 69, N = 75)***
- I (aliphatic) (M = 10, N = 14)**
- nitrile (M = 260, N = 309)***
- N(CH3)2 (aliphatic) (M = 156, N = 176)***
- pyridine (M = 177, N = 192)***
- Br (aliphatic) (M = 98, N = 155)***
- Br (aromatic) (M = 122, N = 181)***
- cyclic carbamate (M = 42, N = 42)*
- NO2 (aromatic) (M = 353, N = 441)***
- OH (aliphatic) (M = 844, N = 1328)***
- COOH (aromatic) (M = 106, N = 115)***
- SO3H (M = 51, N = 72)***



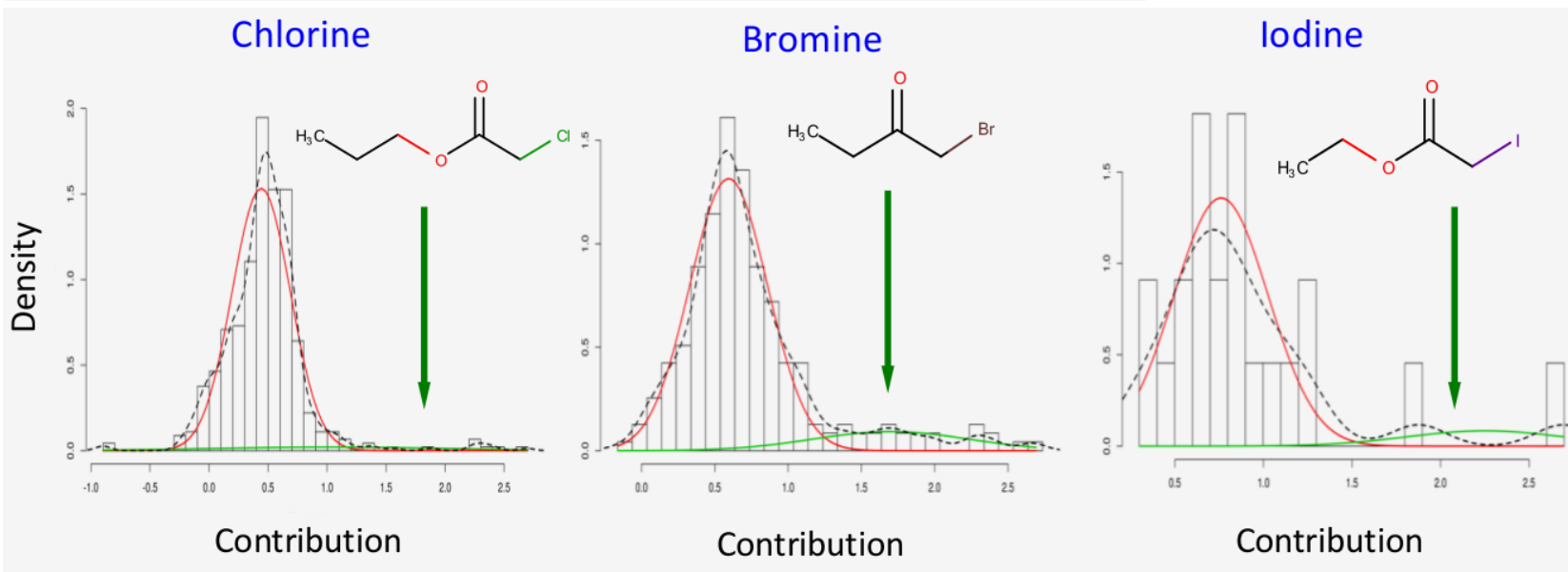
Toxicity towards *Tetrahymena Pyriformis*

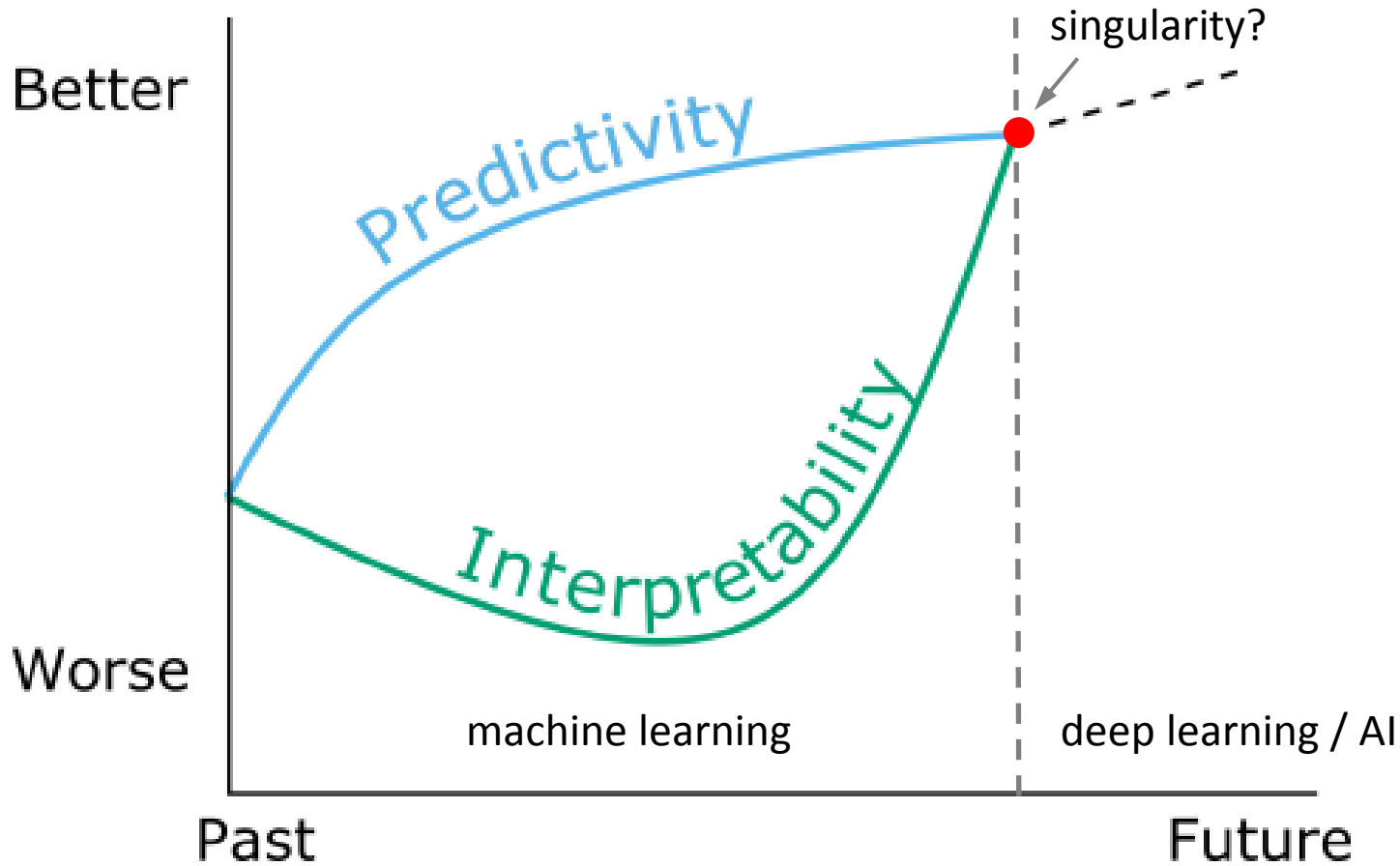
Aliphatic halogens (F < Cl < Br < I)

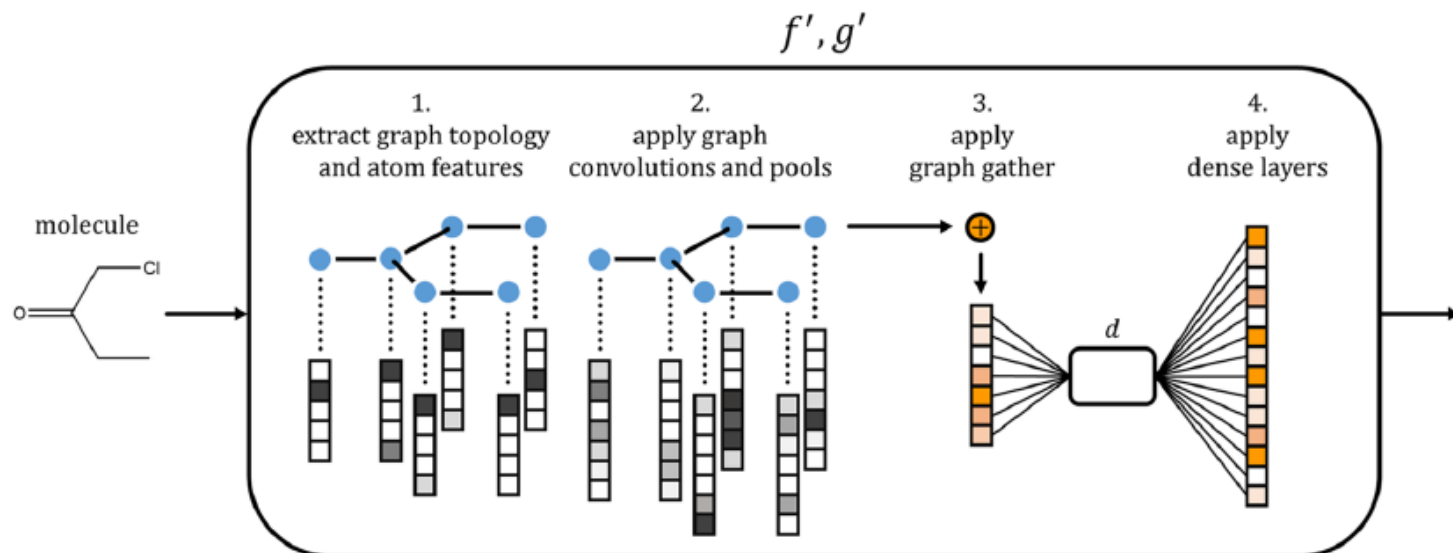
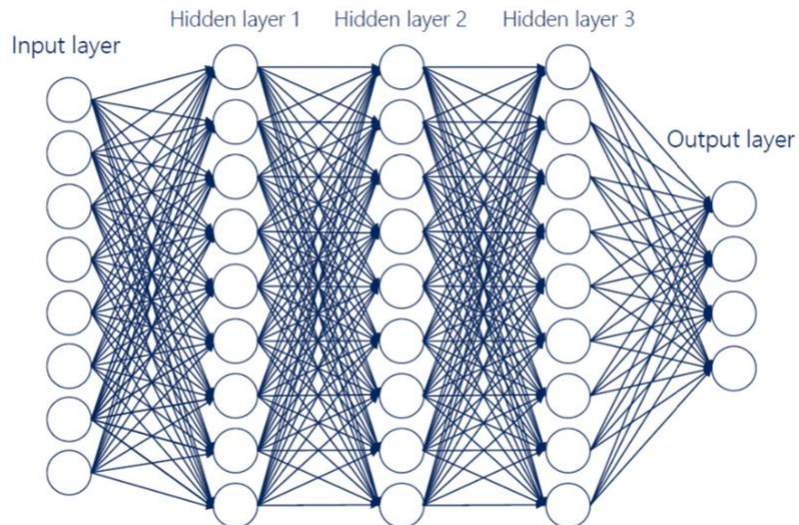


Gaussian Mixture Modeling

SMARTSminer









Self-explaining

attention-based approaches

Gradient-based

GradCAM

GradInput

Integrated Gradients

Perturbation-based

SmoothGrad

Class activation map (CAM)



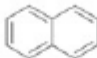



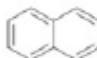
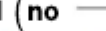






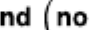




















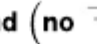




Layer-wise propagation


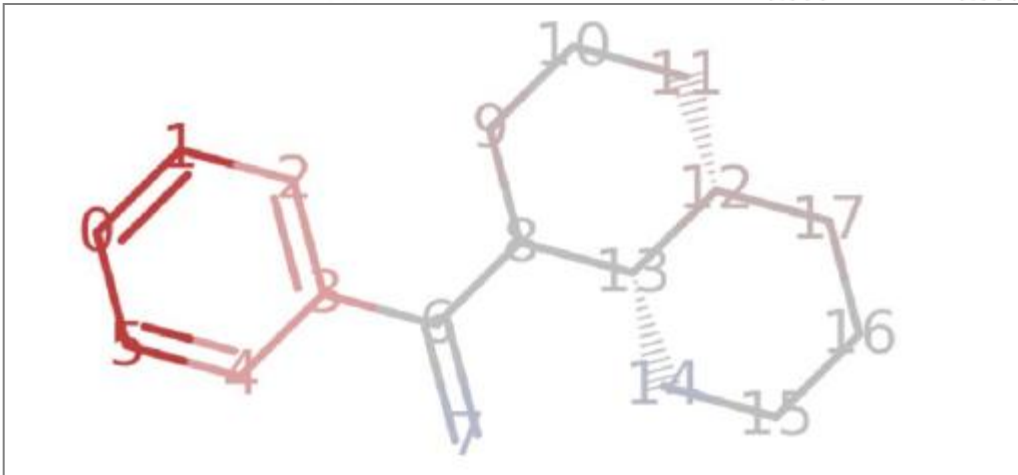
GNNExplainer



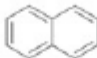



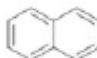














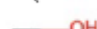













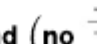




Surrogate modeling

Local Interpretable Model-Agnostic Explanation (LIME)

Shapley additive explanation (SHAP)

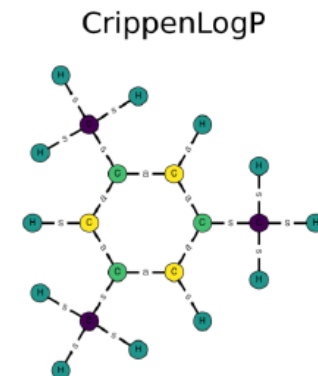
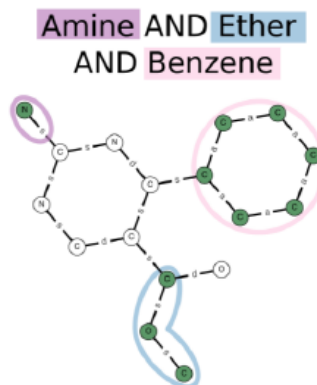
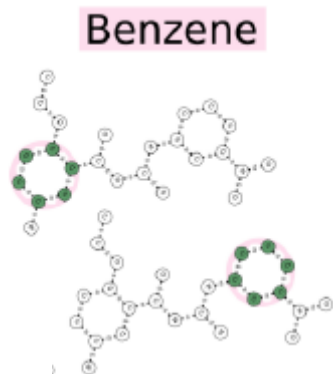
Logic number	Synthetic binding logic	GC Zinc AUC	GC Attribution AUC
0.		1.000	0.980
1.		0.995	0.980
2.		1.000	1.000
3.	no  —NH ₂	1.000	0.970
4.	 or (no )	0.992	0.910
5.	 and (no  —NH ₂)	0.999	0.890
6.	 —F and  =O	1.000	0.770
7.	 and  =O	1.000	0.790
8.	 —F and  —OH and (no )	1.000	0.930
9.	 —NH ₂ and  and 	0.995	0.700
10.	( —NH ₂ or no ) and (no )	0.999	0.860
11.	 —OH and (no  —F) and (no )	1.000	0.880
12.	 —F and  and ( —NH ₂ or no  —OH)	0.999	0.670
13.	( and no  =O) or ( and no )	1.000	0.700
14.	( or no  —OH) and  =O and (no )	1.000	0.750
15.	 and (no ) and  —NH ₂ and  =O	0.996	0.760

Logic number	Synthetic binding logic	GC Zinc AUC	GC Attribution AUC																														
0.		1.000	0.980																														
1.		0.995	0.980																														
2.																																	
3.																																	
4.																																	
5.																																	
6.																																	
7.																																	
8.																																	
9.																																	
10.																																	
11.		<table border="1"> <thead> <tr> <th>Atom index</th> <th>Attribution scores ranked in decreasing order</th> <th>Involved in ground truth binding logic</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>0.29</td> <td>1</td> </tr> <tr> <td>5</td> <td>0.29</td> <td>1</td> </tr> <tr> <td>0</td> <td>0.28</td> <td>1</td> </tr> <tr> <td>2</td> <td>0.09</td> <td>1</td> </tr> <tr> <td>4</td> <td>0.09</td> <td>1</td> </tr> <tr> <td>3</td> <td>0.07</td> <td>1</td> </tr> <tr> <td>9</td> <td>0.03</td> <td>0</td> </tr> <tr> <td>11</td> <td>0.02</td> <td>0</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> </tr> </tbody> </table>	Atom index	Attribution scores ranked in decreasing order	Involved in ground truth binding logic	1	0.29	1	5	0.29	1	0	0.28	1	2	0.09	1	4	0.09	1	3	0.07	1	9	0.03	0	11	0.02	0	
Atom index	Attribution scores ranked in decreasing order	Involved in ground truth binding logic																															
1	0.29	1																															
5	0.29	1																															
0	0.28	1																															
2	0.09	1																															
4	0.09	1																															
3	0.07	1																															
9	0.03	0																															
11	0.02	0																															
...																															
12.																																	
13.																																	
14.																																	
15.																																	

Logic number	Synthetic binding logic	GC Zinc AUC	GC Attribution AUC
0.		1.000	0.980
1.		0.995	0.980
2.		1.000	1.000
3.	no  NH ₂	1.000	0.970
4.	 or (no )	0.992	0.910
5.	 and (no  NH ₂)	0.999	0.890
6.	 F and 	1.000	0.770
7.	 and 	1.000	0.790
8.	 F and  OH and (no )	1.000	0.930
9.	 NH ₂ and  and 	0.995	0.700
10.	( NH ₂ or no ) and (no )	0.999	0.860
11.	 OH and (no  F) and (no )	1.000	0.880
12.	 F and  and ( NH ₂ or no  OH)	0.999	0.670
13.	( and no ) or ( and no )	1.000	0.700
14.	( or no  OH) and  and (no )	1.000	0.750
15.	 and (no ) and  NH ₂ and 	0.996	0.760

binary classification

regression



	Benzene				Amine AND Ether AND Benzene				CrippenLogP			
	GCN	MPNN	GraphNets	GAT	GCN	MPNN	GraphNets	GAT	GCN	MPNN	GraphNets	GAT
Random Baseline	0.61	0.61	0.61	0.61	0.5	0.5	0.5	0.5	0.13	0.13	0.13	0.13
GradInput	0.72	0.54	0.54	0.56	0.52	0.53	0.55	0.41	0.12	0.09	0.13	0.1
SmoothGrad(GI)	0.71	0.54	0.54	0.53	0.51	0.55	0.59	0.38	0.15	0.11	0.15	0.11
GradCAM-last	0.74	0.72	0.66	0.66	0.54	0.74	0.55	0.46	0.04	0.33	0.24	0.07
GradCAM-all	0.75	0.68	0.84	0.62	0.54	0.62	0.7	0.44	0.05	0.27	0.27	0.09
IG	0.97	0.89	0.94	0.95	0.69	0.59	0.72	0.54	0.31	0.24	0.24	0.27
CAM	0.98	0.96	0.76	0.99	0.75	0.76	0.6	0.65	0.2	0.37	0.28	0.23
Attention Weights	--	--	--	0.51	--	--	--	0.51	--	--	--	-0.06

AUC
R_{Kendall}

Benchmark suite for interpretation approaches

regression data sets:

N data set:

N = +1

others = 0

N-O data set:

N = +1

O = -1

others = 0

N+O data set:

N = +0.5

O = +0.5

n(N) = n(O) in molecules

amide data set:

amide group = +1

classification data sets:

amide data set:

presence of an amide group

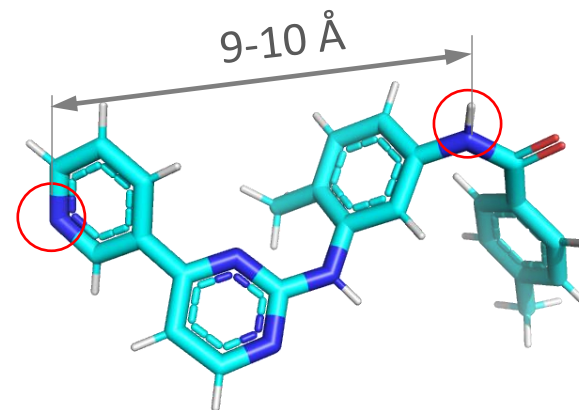
pharmacophore data set:

H-bond donor ... 9-10Å ... H-bond acceptor
(for at least one conformer)

training: 7000

test: 3000

Control possible correlations in data sets
to avoid hidden biases



<https://github.com/ci-lab-cz/ibenchmark>

Benchmark suite for interpretation approaches

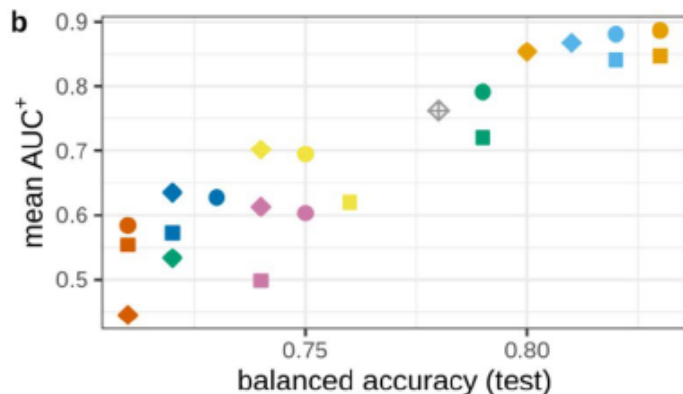
metrics:

AUC – how well a model/approach rank all atoms

top-n – how well a model/approach rank ground truth atoms on top

RMSE = how precisely a model/approach calculate atom contribution

pharmacophore data set

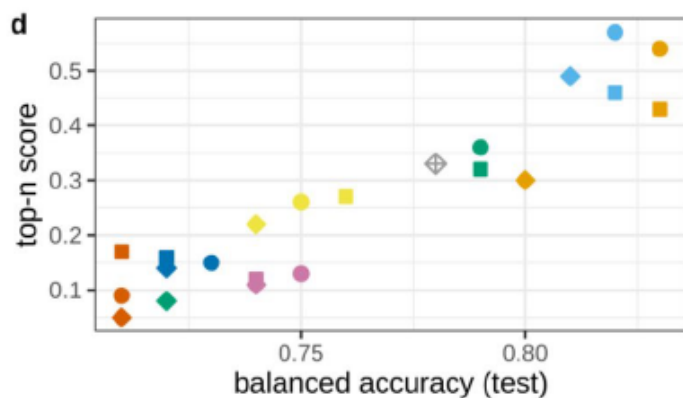


Model

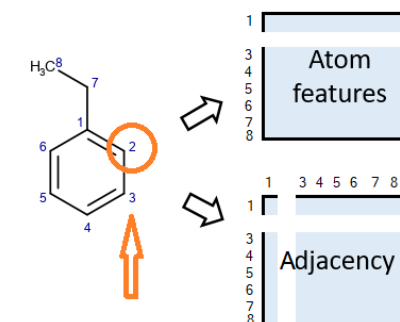
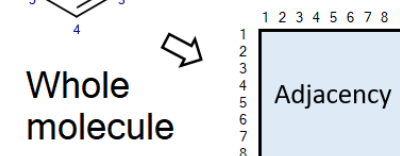
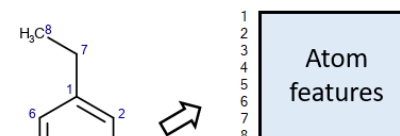
- gbm
- rf
- ◇ svm
- ◇ GC

fp

- AP
- bAP
- MG2
- bMG2
- RDK
- bRDK
- TT
- GC




GCN interpretation



Atom being removed

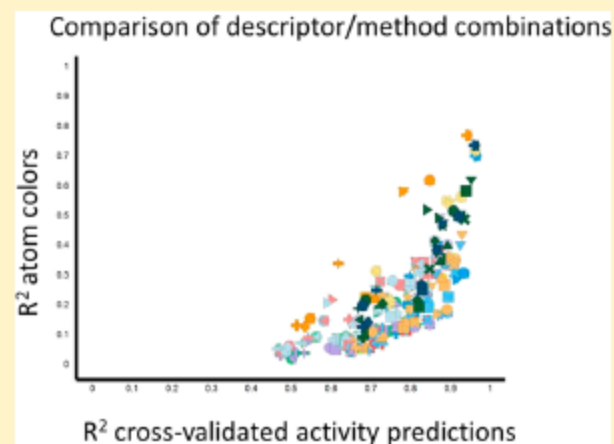
Interpretation of QSAR Models by Coloring Atoms According to Changes in Predicted Activity: How Robust Is It?

 Robert P. Sheridan* 

Modeling and Informatics, Merck & Co. Inc., Kenilworth, New Jersey 07065, United States

Supporting Information

ABSTRACT: Most chemists would agree that the ability to interpret a quantitative structure–activity relationship (QSAR) model is as important as the ability of the model to make accurate predictions. One type of interpretation is coloration of atoms in molecules according to the contribution of each atom to the predicted activity, as in “heat maps”. The ability to determine which parts of a molecule increase the activity in question and which decrease it should be useful to chemists who want to modify the molecule. For that type of application, we would hope the coloration to not be particularly sensitive to the details of model building. In this Article, we examine a number of aspects of coloration against 20 combinations of descriptors and QSAR methods. We demonstrate that atom-level coloration is much less robust to descriptor/method combinations than cross-validated predictions. Even in ideal cases where the contribution of individual atoms is known, we cannot always recover the important atoms for some descriptor/method combinations. Thus, model interpretation by atom coloration may not be as simple as it first appeared.



1. Do we need new interpretation approaches?

Yes, but they should be properly validated and compared with state-of-the-art approaches.

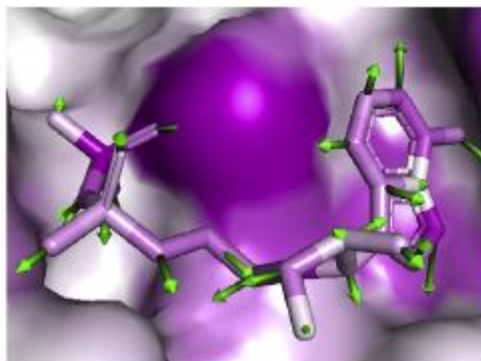
2. What should be the main focus of further research?

Development of approaches able to retrieve new types of knowledge from models (we can already retrieve atom contributions by different methods, let's search for something new and useful)

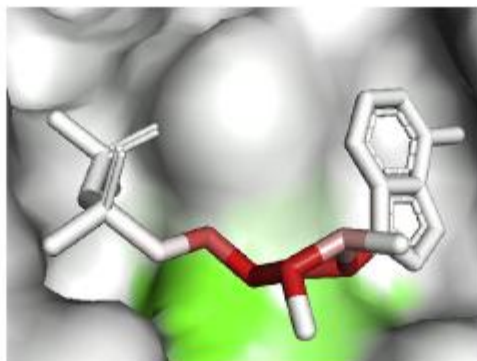
There is a need to introduce interpretation into routine decision making pipelines to supports decisions of medicinal chemists and other researchers

Gnina interpretation

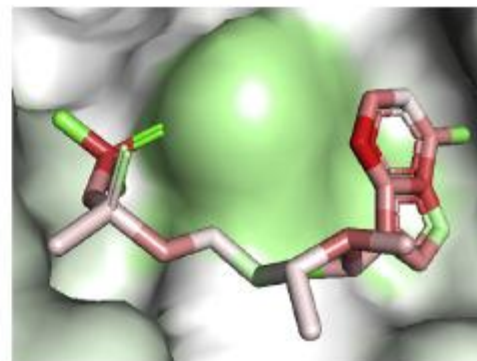
Affinity Prediction Score = 2.698



(a) Gradient

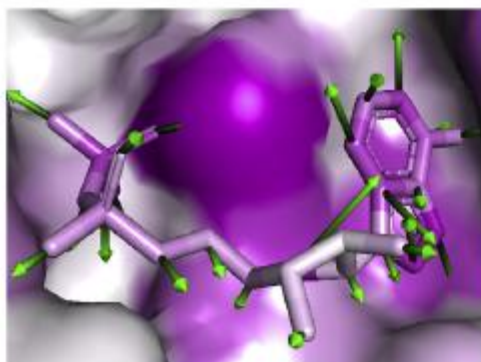


(b) CLRP

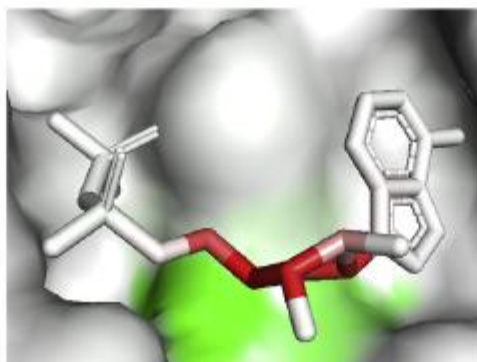


(c) Masking

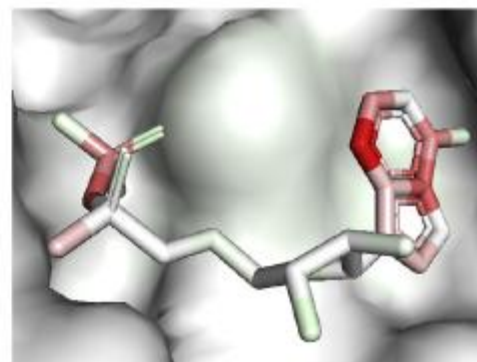
Pose Score = 0.255



(d) Gradient

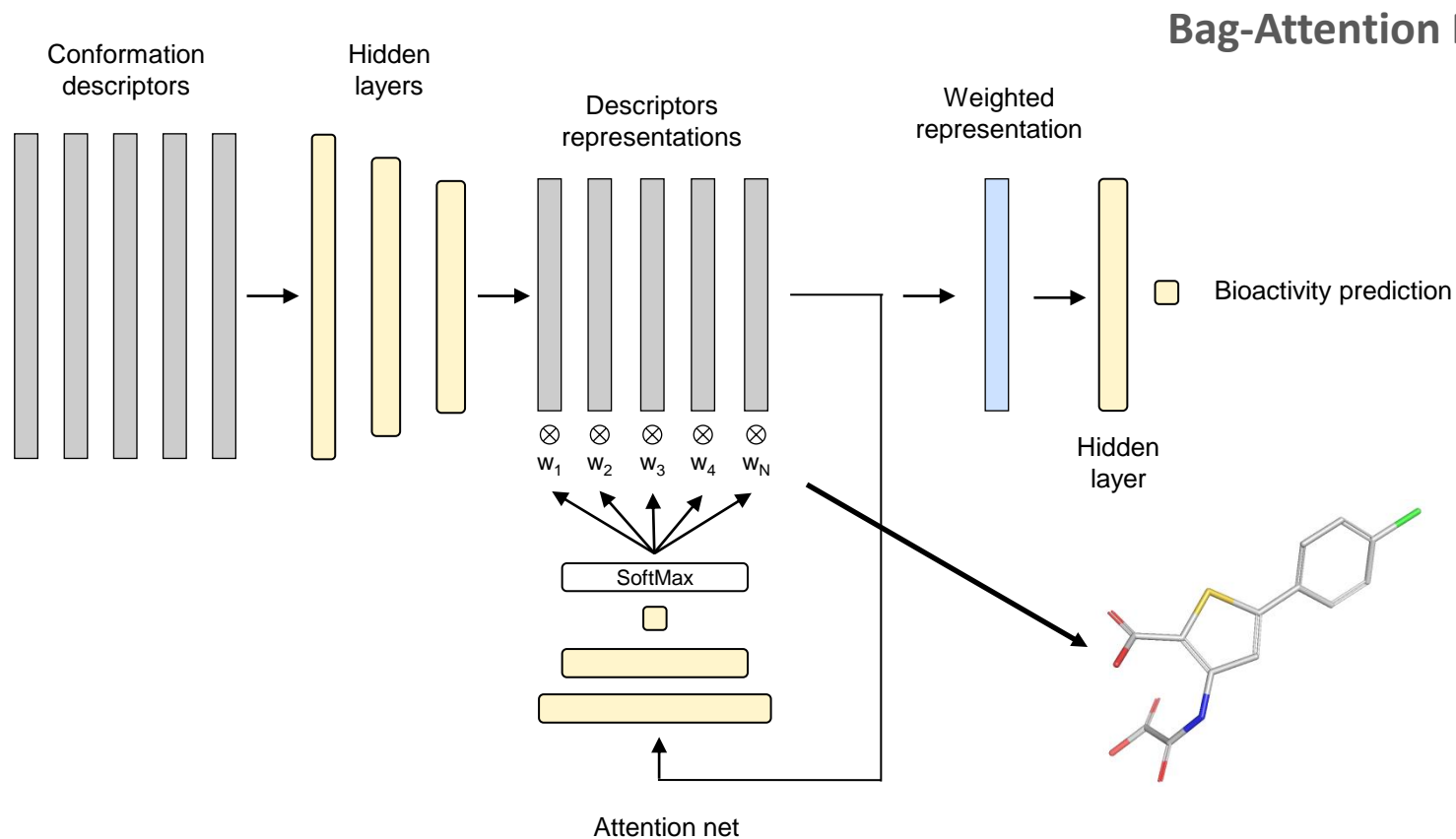


(e) CLRP

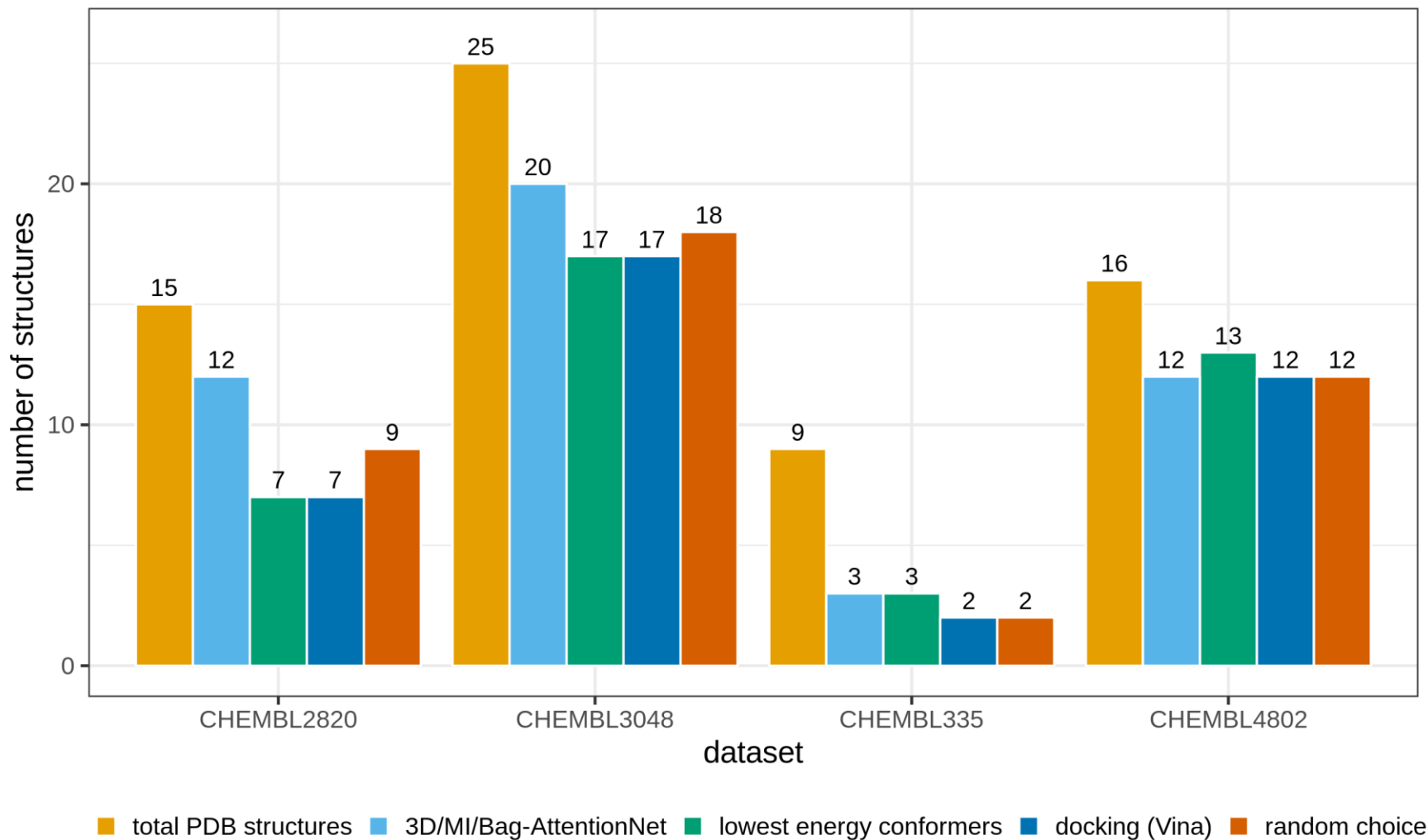


(f) Masking

Identification of bioactive conformers



Selected compounds from test set had average RMSD of generated conformers > 2Å relative to PDB structure



Take-home messages

- There are a lot of hidden treasures in models which can improve our understanding of complex phenomena and augment our knowledge and our goal is to retrieve them
- Do not use anecdotal evidence for evaluation of interpretation performance, do systematic evaluation (use available benchmarks or develop your own)
- The more predictive a model the better interpretation performance, however, even for well predictive models interpretation may be rather low
- Not all interpretation approaches are applicable to chemical problems
- Any predictive model is interpretable (“model \rightarrow structure” paradigm)



pharmaceuticals

an Open Access Journal by MDPI

IMPACT
FACTOR
5.863

Covered in:
PubMed

Pharmacophore Modeling and Applications in Drug Discovery: Challenges and Recent Advances

Guest Editors

Dr. Pavel Polishchuk, Dr. Thomas Seidel

Deadline

28 February 2022

Special Issue

mdpi.com/si/91899

Invitation to submit