#### HelmholtzZentrum münchen German Research Center for Environmental Health

## **In Silico** Prediction of ADMET properties with confidence: potential to speed-up drug discovery

Igor V. Tetko

Helmholtz Zentrum München - German Research Center for Environmental Health (GmbH) Institute of Bioinformatics & Systems Biology

> HELMHOLTZ ASSOCIATION

Siofok, 28 September, Conferentia Chemometrica 2009

## Layout of presentation

#### Introduction:

• Why accuracy of prediction is important?

#### Methods:

• What is a Distance to Model? How can we estimate it? What is a property-based space?

Case study 1: Prediction of environmental toxicity
Case study 2: Benchmarking of lipophilicity (logP) predictions
Case study 3: AMES test prediction
Case study 4: CYP450 prediction

#### Conclusions



#### Which common challenges do they face?



#### HelmholtzZentrum münchen





## "One can not embrace the unembraceable." **Possible:** $10^{60}$ - $10^{100}$ molecules theoretically exist Achievable: 1020 - 1024 can be synthesized now by d Ava Mea 10<sup>80</sup> Kozmin Photkowerse Pro les on t 1,00 There is a need for methods which can estimate the accuracy of predictions! HelmholtzZentrum münchen German Research Center for Environmental Health

#### **Representation of Molecules**



HelmholtzZentrum münchen

German Research Center for Environmental Health



## Examples of distances to models (DM) in descriptor space

City-block

- 1) Only two descriptors are used.
- 2) Colors refer to the same values.
- More complex DMs (property-based DMs) also include the target property.<sup>2</sup>

Jaworska et al, *ATLA*, **2005**, 33, 445-459. Tetko et al, *DDT*, **2006**, *11*, 700-7.



80%

60%

50%

40% 30%

20% 10% 1%

HELMHOLTZ

ASSOCIATION

**Probability-density** 

HelmholtzZentrum münchen German Research Center for Environmental Health

#### The descriptor space challenge



We need to know the target property and select correct descriptors!





## **Property-based space illustration**



Do they agree in their votes (**STD**)? Do they have the same pattern of votes (**CORREL**)?

HelmholtzZentrum münchen German Research Center for Environmental Health



#### **Associative Neural Network Property-Based DMs**



HelmholtzZentrum münchen German Research Center for Environmental Health Tetko et al, *DDT*, **2006**, *11*, *700-7*.

![](_page_8_Picture_4.jpeg)

## 1: Estimation of toxicity against *T. pyriformis*

![](_page_9_Picture_1.jpeg)

T. pyriformis

![](_page_9_Picture_3.jpeg)

Prof. T.W. Schultz

The overall goal is to predict and <u>to assess the reliability of predictions</u> toxicity against T. pyriformis for chemicals directly from their structure.

Dataset: 1093 molecules

HelmholtzZentrum münchen German Research Center for Environmental Health

Zhu et al, J. Chem. Inf. Model, 2008, 48, 766-84.

![](_page_9_Picture_9.jpeg)

## CAse studies on the development and application of in-silico techniques for environmental hazard and risk assessment

www.CADASTER.eu

![](_page_10_Picture_2.jpeg)

HelmholtzZentrum münchen C German Research Center for Environmental Health

Challenge (deadline was Sep. 10) co-organized with the European Neural Network Society

![](_page_10_Picture_5.jpeg)

#### Analyzed QSARs (Quantitative Structure Activity Relationship) and distances to models (DM)

country	modeling techniques	descriptors	abbreviation	distances to models (in space)		
				descriptors	property-based	
	ensemble of 192	MolconnZ	kNN-MZ	EUCLID	STD	
	kNN models					
1251	ensemble of 542	Dragon	kNN-DR	EUCLID	STD	
A MARTIN	kNN models					
11 Pr.	SVM	MolconnZ	SVM-MZ			
	SVM	Dragon	SVM-DR			
	SVM	Fragments	SVM-FR			
G	kNN	Fragments	kNN-FR	EUCLID,		
	MLR	Fragments	MLR-FR	TANIMOTO		
	MLR	Molec. properties	MLR-COD			
		(CODESSA-Pro)				
	OLS	Dragon	OLS-DR	LEVERAGE		
C	PLS	Dragon	PLS-DR	LEVERAGE	PLSEU	
	ensemble of 100		ASNN-			
100.0	neural networks	E-state indices	ESTATE		CORREL. STD	
					,	
A 11	consensus model	-	CONS		STD	

German Research Center for Environmental Health

Н

*l*, **2008,** 48(9):1/33-46.

ASSOCIATION

### **Overview of analyzed distances to models (DMs)**

<b>EUCLID</b> $EU_m = \frac{\sum_{j=1}^k d_j}{k}$ k is number of nearest $EUCLID = E\overline{U}_m$ neighbors, m index of model	<b>TANIMOTO</b> $Tanimoto(a,b) = \frac{\sum_{a,i} x_{a,i} x_{b,i}}{\sum_{a,i} x_{a,i} + \sum_{a,i} x_{b,i} - \sum_{a,i} x_{b,i}}$ $x_{a,i} \text{ and } x_{b,i} \text{ are fragment counts}$				
LEVERAGE	PLSEU (DModX)				
$LEVERAGE = \mathbf{x}^{\mathrm{T}} (\mathbf{X}^{\mathrm{T}} \mathbf{X})^{-1} \mathbf{x}$	Error in approximation (restoration) of the vector of input variables from the latent variables and PLS weights.				
<b>STD</b> $STD = \frac{1}{N-1} \sum (y_i - \overline{y})^2$	CORREL				
	$CORREL(a) = \max_{j} CORREL(a,j) = R^{2}(\mathbf{Y}_{calc}^{a}, \mathbf{Y}_{calc}^{j})$				
$y_i$ is value calculated with model <i>i</i> and $\overline{y}$ is average value	$Y^{a}=(y_{1},,y_{N})$ is vector of predictions of molecule <i>i</i>				

![](_page_12_Picture_3.jpeg)

## STD

#### Property-based space: DM does work!

![](_page_13_Figure_2.jpeg)

HelmholtzZentrum münchen German Research Center for Environmental Health

Tetko et al, J. Chem. Inf. Model, 2008, 48, 1733-46.

![](_page_13_Picture_5.jpeg)

#### **Descriptor space: DM** <u>does not work</u>

![](_page_14_Figure_1.jpeg)

HelmholtzZentrum münchen German Research Center for Environmental Health

Tetko et al, J. Chem. Inf. Model, 2008, 48, 1733-46.

![](_page_14_Picture_4.jpeg)

### Mixture of Gaussian Distributions (MGD)

Idea is to find a MGD, which maximize likelihood (probability)

 $\Pi N(0,\sigma^2(e_i))$ 

of the observed distribution of errors

![](_page_15_Figure_4.jpeg)

![](_page_15_Picture_5.jpeg)

![](_page_15_Picture_6.jpeg)

DM	-	ave	rage rank	highest rank <sup>1</sup>				
	LOO	5-CV	Valid.*	LOO	5-CV	Valid.		
STD-CONS	1	1.8	1.1	12	2	11		
STD-ASNN	2	1.2	2.5		10	1		
STD-kNN-DR	6.6	4.3	4.1					
STD-kNN-MZ	9.2	8.3	5.3					
EUCLID-kNN-DR	7.1	4.9	5.4					
LEVERAGE-PLS	8.4	5	6.3					
EUCLID-kNN-MZ	7.5	7.1	6.4					
TANIMOTO-kNN-FR	7	6.1	6.8					
TANIMOTO-MLR-FR	8.3	8.3	9					
CORREL-ASNN	10.7	10.8	9.4					
LEVERAGE-OLS-DR	12.3	12.6	11.1					
EUCLID-MLR-FR	7	9.3	11.5					
PLSEU-PLS	11.1	11.8	11.5					
EUCLID-kNN-FR	12.1	13.3	12.1					

#### Ranking of Distance to Models (DM)

\*Ordered by performance of the DMs on the validation dataset

HelmholtzZentrum münchen German Research Center for Environmental Health

Tetko et al, J. Chem. Inf. Model, 2008, 48, 1733-46.

![](_page_16_Picture_5.jpeg)

## STD

#### Analysis of DMs for a linear model

![](_page_17_Figure_2.jpeg)

HelmholtzZentrum münchen German Research Center for Environmental Health

Tetko et al, *J Chem Inf Model*, **2008**, 48(9):1733-46.

![](_page_17_Picture_5.jpeg)

## **2: Benchmarking of logP calculators**

Existing Dogma:

- Prediction of physico-chemical properties, in particular log P, is simple
- There is no need to measure them
- We have enough number of good computational methods

## Is this true?

![](_page_18_Picture_7.jpeg)

#### **Data & background models**

18 methods (major commercial providers and public software)

*in house* data: 95809 molecules from Prizer 889 molecules from Nycomed

#### Arithmetic Average Model (AAM):

mean log*P* was used as a prediction (one value for all molecules)

Rank III: models with errors (RMSE)  $\geq$  AAM, i.e. <u>non-predictive</u> Rank I: models with RMSE identical or close to the best method Rank II: remaining models

![](_page_19_Picture_7.jpeg)

#### Benchmarking of logP methods for in-house data of Pfizer & Nycomed Pfizer set (N = 95 809)

#### Large number of methods could not perform better than the **AAM** model !

Catastrophe !?

	Pfizer set ( <i>N</i> = 95 809)					Nycomed set (N = 882)					
Method	RMSF	rank	% in error range			RMSE,	RMSF	rank	% in error range		
		- unix	<0.5	0.5- 1	· >1	zwitterions excluded <sup>2</sup>	, and L	- unit	<0.5	0.5- 1	>1
Consensus logP	0.95	I	48	29	24	0.94	0.58	I	61	32	7
ALOGPS	1.02	Т	41	30	29	1.01	0.68	Т	51	34	15
S+logP	1.02	Т	44	29	27	1.00	0.69	Т	58	27	15
NC+NHET	1.04	II	38	30	32	1.04	0.88	Ш	42	32	26
MLOGP(S+)	1.05	II	40	29	31	1.05	1.17	Ш	32	26	41
XLOGP3	1.07	II	43	28	29	1.06	0.65	Т	55	34	12
MiLogP	1.10	II	41	28	30	1.09	0.67	Т	60	26	14
AB/LogP	1.12	II	39	29	33	1.11	0.88	Ш	45	28	27
ALOGP	1.12	II	39	29	32	1.12	0.72	II	52	33	15
ALOGP98	1.12	II	40	28	32	1.10	0.73	II	52	31	17
OsirisP	1.13	II	39	28	33	1.12	0.85	II	43	33	24
AAM	1.16	Ш	33	<b>29</b>	38	1.16	0.94	Ш	<b>42</b>	31	27
CLOGP	1.23	Ш	37	28	35	1.21	1.01	Ш	46	28	22
ACD/logP	1.28	Ш	35	27	38	1.28	0.87	Ш	46	34	21
CSlogP	1.29	Ш	37	27	36	1.28	1.06	Ш	38	29	33
COSMOFrag	1.30	Ш	32	27	40	1.30	1.06	Ш	29	31	40
QikProp	1.32	Ш	31	26	43	1.32	1.17	Ш	27	24	49
KowWIN	1.32	Ш	33	26	41	1.31	1.20	Ш	29	27	44
QLogP	1.33	Ш	34	27	39	1.32	0.80	Ш	50	33	17
XLOGP2	1.80	Ш	15	17	68	1.80	0.94	Ш	39	31	29
MLOGP(Dragon)	2.03	Ш	34	24	42	2.03	0.90	Ш	45	30	25

HelmholtzZentrum münchen German Research Center for Environmental Health

Mannhold et al, J. Pharm. Sci., 2009, 98, 861-893.

![](_page_20_Picture_6.jpeg)

http://vcclab.org

#### Virtual Computational Chemistry Laboratory

#### ALOGPS 2.1

•LogP: 75 variables, 12908 molecules, RMSE=0.35, MAE=0.26

•LogS: 33 variables, 1291 molecules, RMSE=0.49, MAE=0.35

Tetko et al, *J. Comput. Aided Mol. Des.* **2005**, 19, 453-63.

Tetko & Tanchuk, *J. Chem. Info. Comput. Sci.*, **2002**, 42, 1136-45.

HelmholtzZentrum münchen German Research Center for Environmental He

![](_page_21_Figure_8.jpeg)

#### **CORREL** ALOGPS self-learns new data to cover new scaffolds

N=95809 (*in house* Pfizer data)

ALOGPS Blind prediction

ALOGPS LIBRARY

![](_page_22_Figure_4.jpeg)

HelmholtzZentrum münchen German Research Center for Environmental Health

Tetko et al, QSAR Comb. Sci., 2009, 28, 845-9.

![](_page_22_Picture_7.jpeg)

![](_page_23_Picture_0.jpeg)

# Local correction of a model based on nearest neighbors

![](_page_23_Picture_2.jpeg)

HelmholtzZentrum münchen German Research Center for Environmental Health

![](_page_23_Picture_4.jpeg)

![](_page_24_Picture_0.jpeg)

## **Estimation of the model accuracy by the distance to nearest neighbors**

![](_page_24_Picture_2.jpeg)

HelmholtzZentrum münchen German Research Center for Environmental Health

![](_page_24_Picture_4.jpeg)

#### ALOGPS distinguishes reliable vs. non-reliable predictions in property-based space (CORREL)

![](_page_25_Figure_1.jpeg)

HelmholtzZentrum münchen German Research Center for Environmental Health

Tetko et al, *Chemistry & Biodiversity*, **2009**, in press.

![](_page_25_Picture_4.jpeg)

![](_page_26_Picture_0.jpeg)

#### **ALOGPS dramatically improves accuracy**

![](_page_26_Figure_2.jpeg)

Only reliable predictions (and we know "who is who"!) have much higher accuracy.

![](_page_26_Picture_5.jpeg)

### **3: Prediction of Ames Mutagenicity set**

http://ml.cs.tu-berlin.de/toxbenchmark Toxicity against *Salmonella typhimurium* 

Data set: 4361 molecules

67% with mutagenic effect (**background model**)

Large international collaboration effort of 13 labs from USA, Canada, EU, Russia, Ukraine & China

![](_page_27_Picture_5.jpeg)

Prof. Bruce N. Ames Inventor of the test (1975)

HelmholtzZentrum münchen German Research Center for Environmental Health

<sup>1</sup>Schwaighofer et al, *JCIM*, **2008**, 48, 785-96.

![](_page_27_Picture_9.jpeg)

#### STD Associative Neural Network analysis of Ames set

![](_page_28_Figure_1.jpeg)

Only reliable predictions (15% of all data points) are 22%/5% = 4 times more accurate!

HelmholtzZentrum münchen German Research Center for Environmental Health

in preparation

![](_page_28_Picture_5.jpeg)

## 4: Prediction of CYP450 1A2 inhibitors

http://pubchem.ncbi.nlm.nih.gov/assay/assay.cgi?aid=410

Bioassay AID 410

One of the test performed within NIH Roadmap

4177 active molecules3680 inactive molecules

53% were inhibitors of CYP (**background accuracy**)

![](_page_29_Picture_6.jpeg)

**Dr. Elias Zerhouni** Former NIH director (2002-2008)

HelmholtzZentrum münchen German Research Center for Environmental Health

![](_page_29_Picture_9.jpeg)

#### STD Associative Neural Network analysis of CYP450 set

![](_page_30_Figure_1.jpeg)

The most reliable predictions (30% of all molecules) are 21%/5% = 4 times more accurate!

HelmholtzZentrum münchen German Research Center for Environmental Health

in preparation

![](_page_30_Picture_5.jpeg)

### **ADMETox and in silico challenges**

![](_page_31_Figure_1.jpeg)

Developed methodology allows navigation in space of molecules with a confidence and:

- $\checkmark$  to develop targeted (local) models to cover specific series.
- $\checkmark$  to reliably estimate which compounds can/can't be reliably predicted.
- $\checkmark$  to provide experimental design and to minimize costs of new measurements.
- □ This is our expertise and "know-how" that we are applying to new data.

![](_page_31_Picture_8.jpeg)

## Acknowledgements

![](_page_32_Picture_1.jpeg)

#### Funding

GO-Bio BMBF <u>http://qspr.eu</u> Germany-Ukraine grant UKR 08/006 FP7 Marie Curie ITN ECO <u>http://eco-itn.eu</u> FP7 CADASTER <u>http://www.cadaster.eu</u>

#### My group

Iurii Sushko Sergii Novotarskyi Anil K. Pandey Robert Körner Stefan Brandmaier Matthias Rupp Vijyant Srivastava Wolfram Teetz

Collaborators: Dr. G. Poda (Pfizer) Dr. C. Ostermann (Nycomed) Dr. C. Höfer (DMPKore) Prof. A. Tropsha (NC, USA) Prof. T. Oprea (New Mexico, USA) Prof. A. Varnek (Strasbourg, France) Prof. R. Mannhold (Düsseldorf, Germany) Prof. R. Todeschini (Milano, Italy) + many other colleagues

![](_page_32_Picture_7.jpeg)

## HELMHOLTZ

HelmholtzZentrum münchen

German Research Center for Environmental Health