

**HelmholtzZentrum münchen**

German Research Center for Environmental Health

# Introduction to Chemoinformatics

Dr. Igor V. Tetko

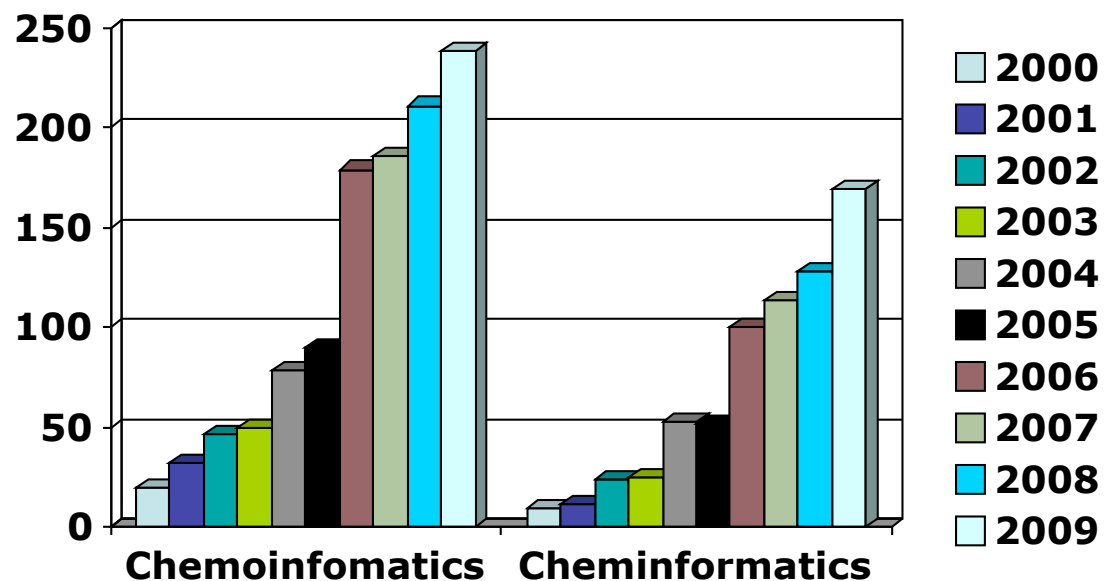
Helmholtz Zentrum München - German Research Center for Environmental Health (GmbH)  
Institute of Bioinformatics & Systems Biology (HMGU)

**Kyiv, 10 August 2009, Summer School**

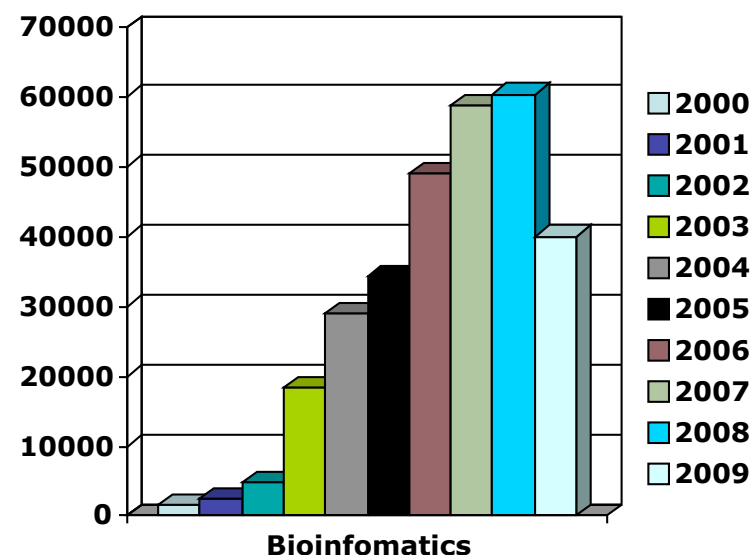
# Layout

- Chemoinformatics??? What does it mean?
- Role of chemoinformatics
  - in drug discovery
  - Chemical Biology, NIH Roadmaps
  - REACH, chemical safety
- Definitions of chemoinformatics
- OCHEM
- Overview of the course

# Cheminformatics & Bioinformatics keywords in SCOPUS



There was no dedicated journal,  
many authors does not explicitly use world  
chemoinformatics



about 10x lower if journals  
and references are excluded

# Dmitrii Ivanovich Mendeleev, 1834-1907

Discoverer of the Periodic Table — An Early “Chemoinformatician”





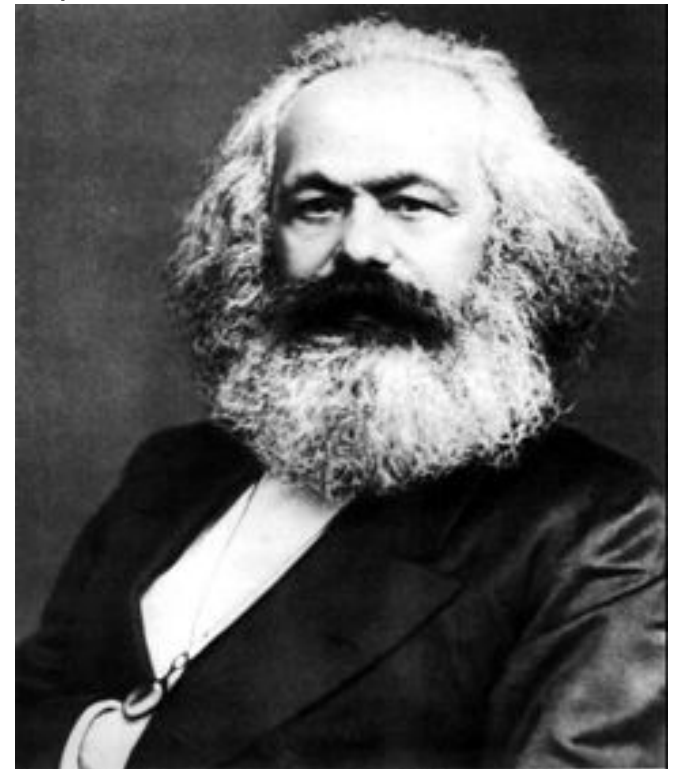
# Why Mendeleev?

Faced with a large amount of data, with many gaps, Mendeleev:

- Sought patterns where none were obvious,
- Made predictions about properties of unknown chemical substances, based on observed properties of known substances,
- Created a great visualization tool!

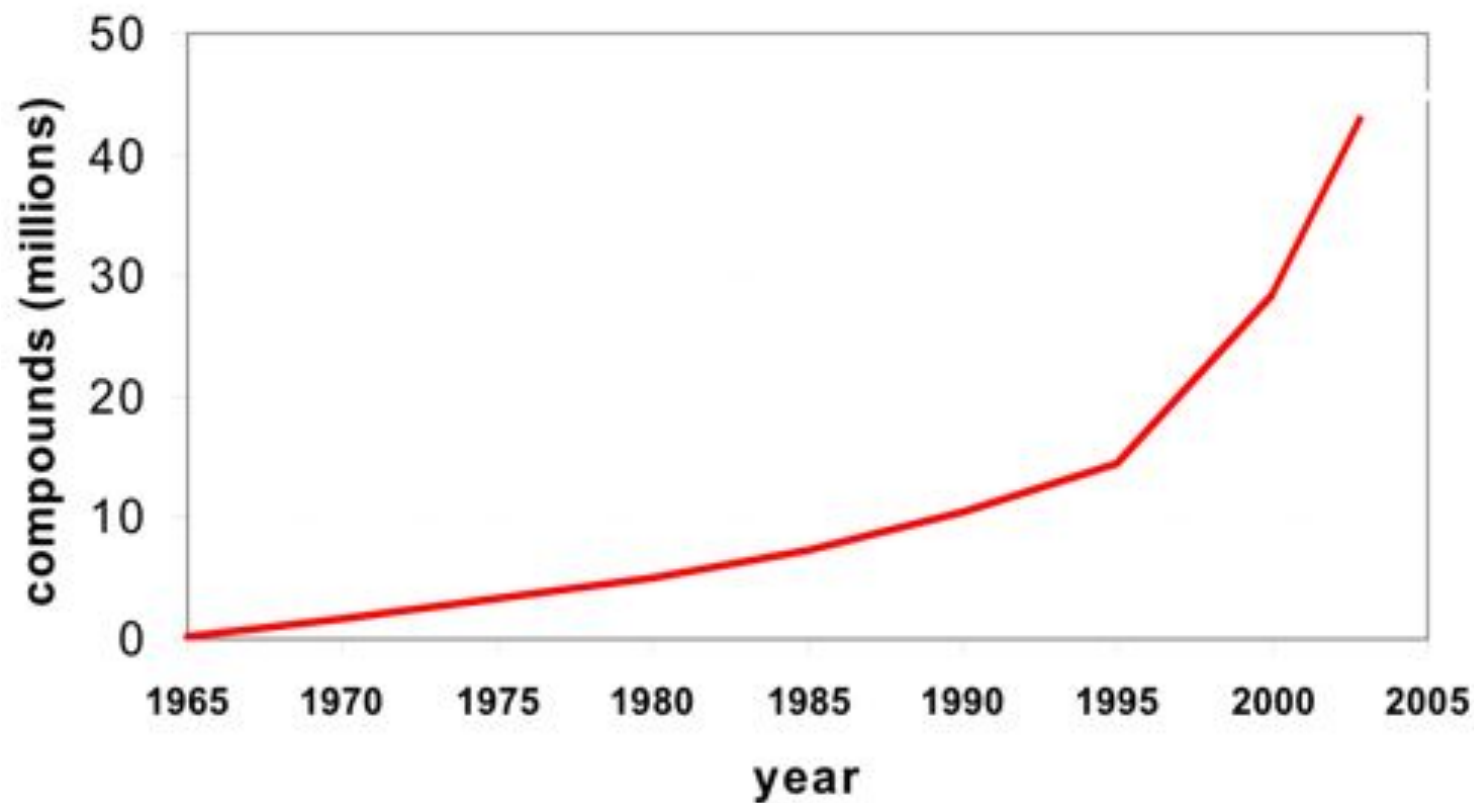
Why he? Why at that time?

“The **Law** of transformation of **Quantity** into **Quality**”  
by Karl Marx



# Quantity aspects of chemoinformatics

compounds published in CAS



# Major challenges of Chemoinformatics

Millions of structures

Thousand of publications

Storage, organization and search of information

QSAR / QSPR studies: prediction of properties and activities

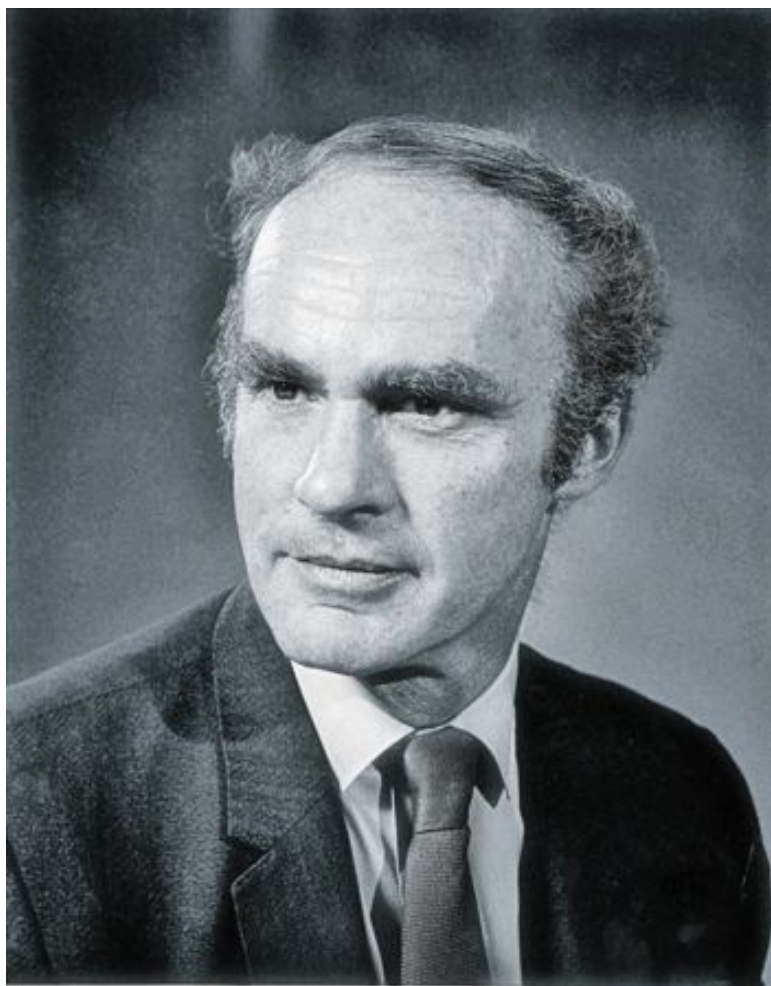
Chemical biology & drug discovery

*In silico* environmental toxicology, REACH

QSAR/QSPR - Quantitative Structure-Activity (Property) Relationship Studies



# Goal of chemistry: Synthesis of Properties



The most fundamental and lasting objective  
of synthesis is  
**not production of new compounds**  
but  
**production of properties**

George S. Hammond  
Norris Award Lecture, 1968

# Where Chemoinformatics is required?

## Pharma companies

- data collection and handling
- model development QSAR/QSPR, *in silico* design
- *In vitro*, *in vivo* data analysis and interpretation

## REACH

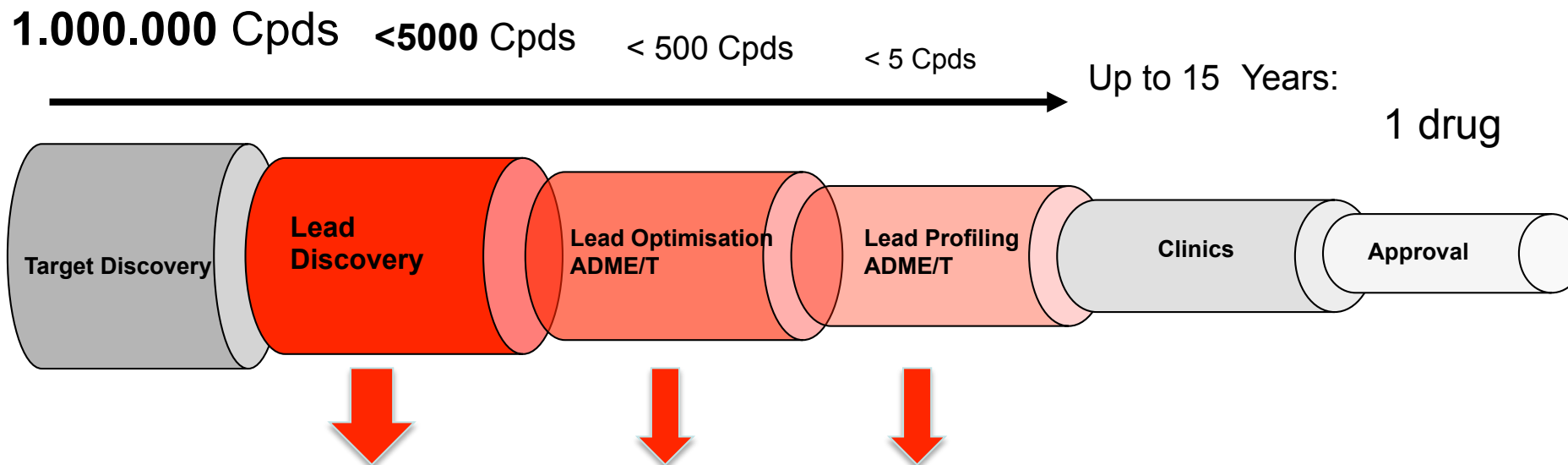
- > 140,000 compounds
- Development of *in silico* methods to predict toxicity of chemical compounds

## Chemical industry

- Design of new properties; chemical synthesis
- Prediction of toxicity of chemicals BEFORE they enter market

# Pharma R&D: Cost and Productivity issues

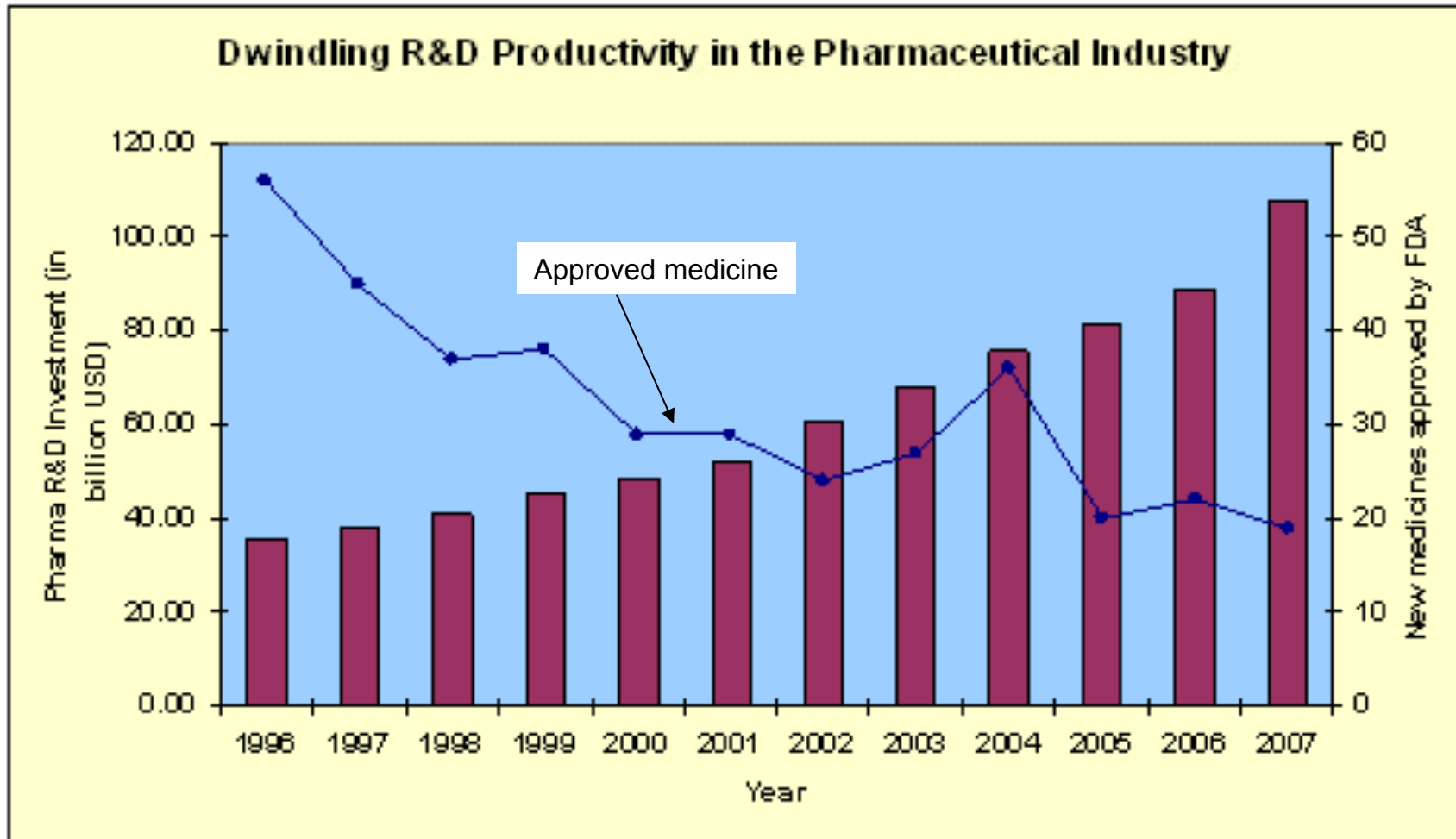
## Compound numbers



In vitro and in vivo property determination:  
Millions of screens for solubility, stability, absorption, metabolism, transport,  
reactive products, drug interactions, etc etc

Preclinics Costs: > \$300m PER COMPOUND to reach approval

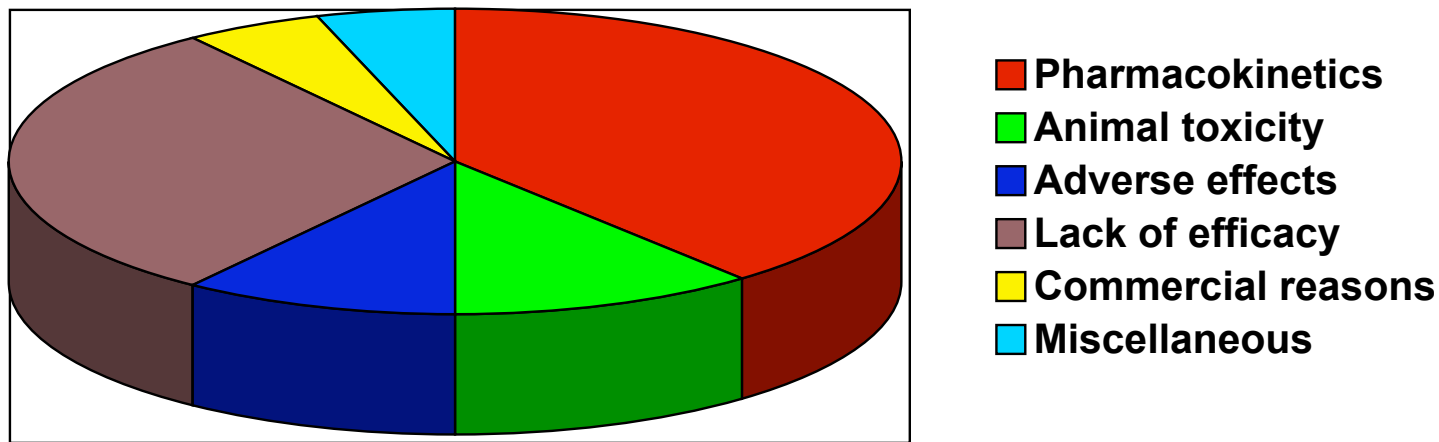
# Declining R&D productivity in the pharmaceutical industry



See <http://www.frost.com/prod/servlet/market-insight-top.pag?docid=128394740>

Source : PhRMA 2007, FDA

# Pharma R&D Cost and productivity: Reasons for compound failure



TOP four reasons are connected to compound **A**bsorption, **D**istribution, **M**etabolism and **E**xcretion, all of which may contribute to lack of efficacy and **T**oxicity: **ADME/T** issues

**Chemoinformatics!**

# NIH Roadmaps



Chemical Biology



## **Dr. Elias Zerhouni**

Former NIH director (2002-2008)  
\$29.5 billions in 2008 for  
27 Institutes

Ukraine GDP is ca \$180 billions  
Belarus GDP is ca \$60 billions

# NIH Roadmaps

New pathways to **discovery**

- understand **complex biological systems**
- understand their **connections**
- build **better "toolbox"** for researchers

Research **teams** of the future

Re-engineering the **clinical research** team

*Molecular Libraries Screening Center Network (MLSCN)*

- *10 centers were created*
- *250 assays to screen >200,000 molecules*

**Chemoinformatics!**

- **PubChem** database to handle data

# REACH

Registration, Evaluation, Authorisation and  
Restriction of Chemical substances



European Chemicals Agency (ECHA) in Helsinki





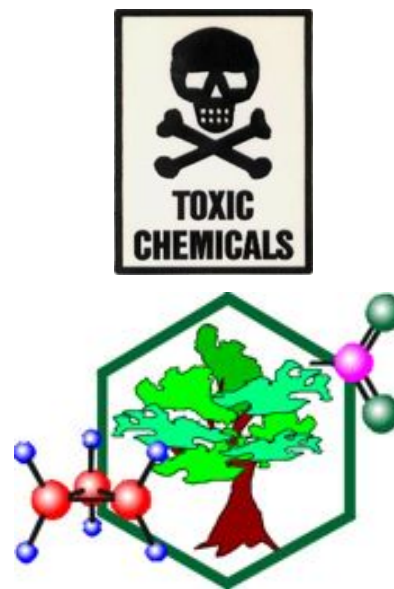
# REACH and Environmental Chemoinformatics

> 140,000 chemicals to be registered ... is a lot!

It is expensive to measure all of them (\$200,000 per compound), a lot of animal testing => **€24 billions over next 10 years**

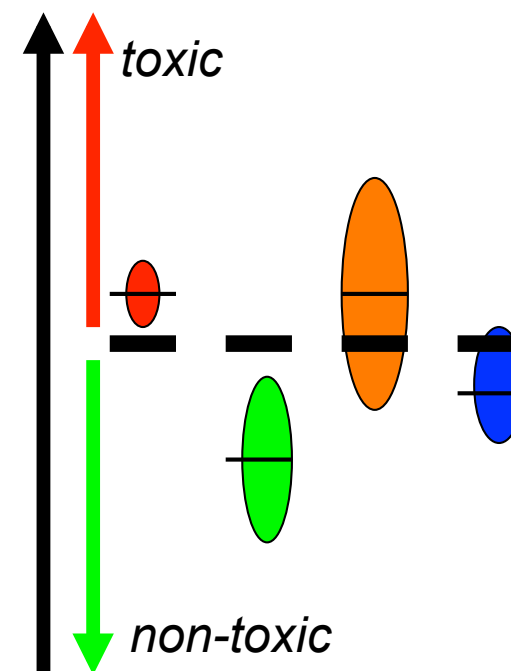
QSAR models can be used to prioritize compounds

- Compound is predicted to be toxic
  - Biological testing will be done to prove/disprove the models
- Compound is predicted to be not toxic
  - tests can be avoided, saving money, animals
  - but ... only if we are confident in the predictions



# Requirements of biological testing following QSAR model prediction

model prediction	prediction confidence	
	high	low
toxic, $EC_{50} > LIMIT$	<b>strong need</b>	<b>moderate need</b> (depends on other properties)
non-toxic, $EC_{50} < LIMIT$	<b>no need</b>	<b>low need</b> (depends on other properties)



Acceptance of decisions is more accurate if confidence intervals (prediction errors) are known and are taken into analysis: concept of applicability domain.

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

[www.CADASTER.eu](http://www.CADASTER.eu)

[Home](#) [REACH](#) [CHALLENGE](#) [Partners](#) [WP](#) [Meetings](#) [News](#) [Forum](#) [Login](#)



CAsE studies on the Development and Application of in-Silico Techniques for Environmental hazard and Risk assessment

Search this site:

### Related topics

- Open Positions
- People
- Publications
- Related Projects
- Links
- Contact

### Latest news

- TRISK is now open for application
- Challenge on [www.CADASTER.eu](http://www.CADASTER.eu)
- We are online !!!

## Home



### About CADASTER

Implementation of REACH requires demonstration of the safe manufacture and use of chemicals. REACH aims to achieve a proper balance between societal, economic and environmental objectives, and attempts to efficiently use the scarce and scattered information available on the majority of substances. Thereupon REACH aims to reduce animal testing by optimized use of in silico and in vitro information on related compounds.

The REACH regulation advocates the use of non-animal testing methods, but guidance is needed on how these methods should be used. The procedures include alternative methods such as chemical and biological read-across, in vitro results, in vivo information on analogues, (Q)SARs, and exposure-based waiving. The concept of Intelligent Testing Strategies for regulatory endpoints has been outlined to facilitate the assessments. Intensive efforts are needed to translate the concept into a workable, consensually acceptable, and scientifically sound strategy.

**CADASTER** aims at providing the practical guidance to integrated risk assessment by carrying out a full hazard and risk assessment for chemicals belonging to four compound classes. A Decision Support System (DSS) will be developed that will be updated on a regular basis in order to accommodate and integrate the alternative methods mentioned above.

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

 **HELMHOLTZ**  
ASSOCIATION

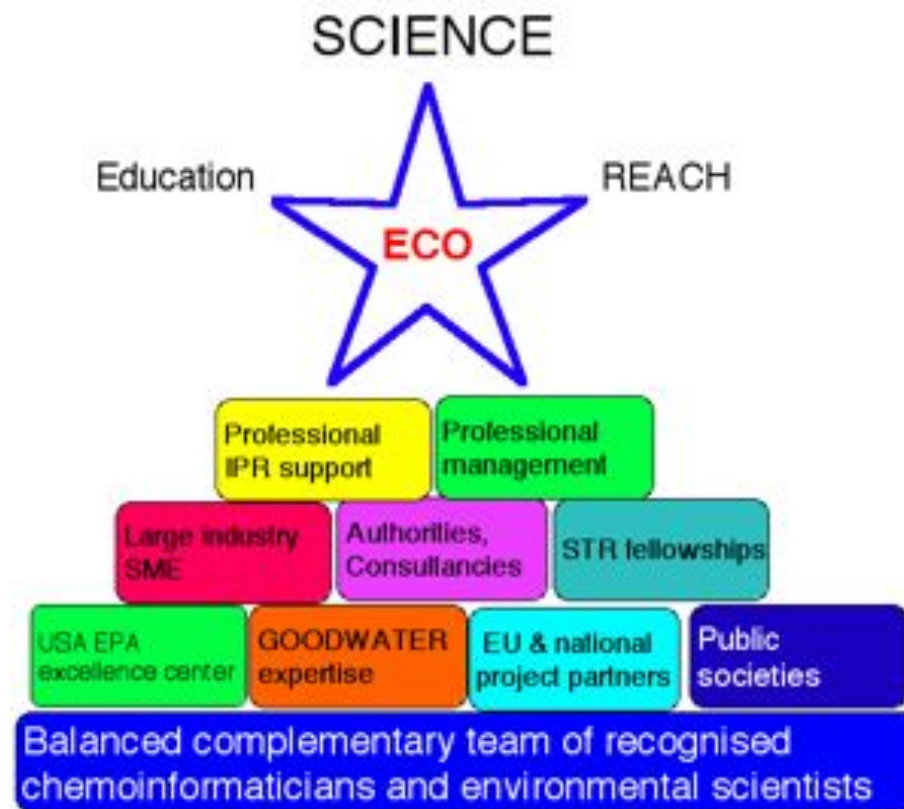
# Environmental ChemOinformatics (ECO) Marie Curie Initial Training Center

**Main goal:** training of new personnel with respect to REACH

**Team:** 7 partners from 5 EU countries

**Offered training:**

- 11 PhD students
- 1 postdoc
- 37 short-term fellowships (3-12 months)
- **in total:** ca 50 years of positions



# Definitions of Chemoinformatics

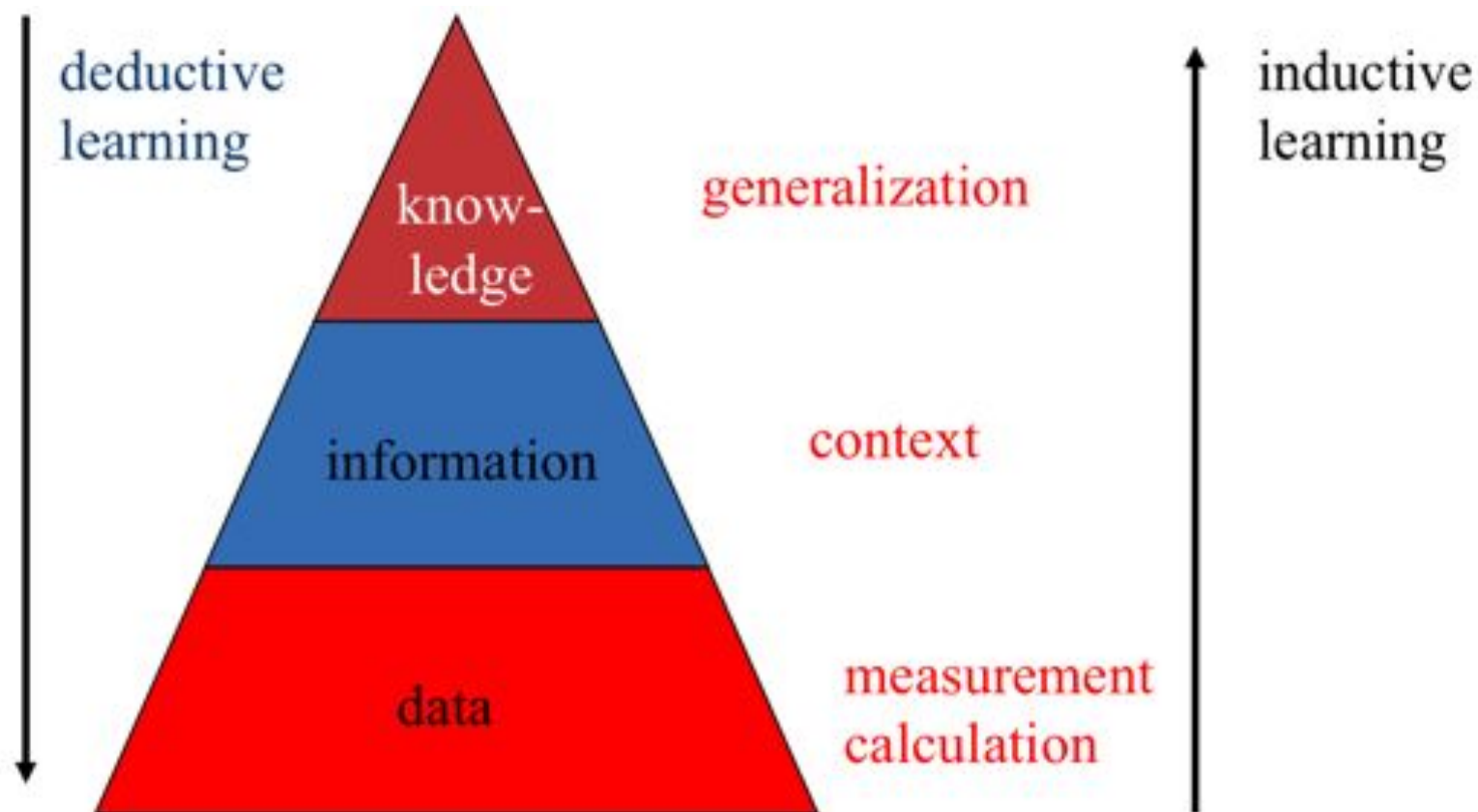
*Chemoinformatics* is a generic term that encompasses the design, creation, organization, management, retrieval, analysis, dissemination, visualization, and use of chemical information. *G. Paris, 1988*

*Chemoinformatics* is the mixing of those information resources to transform data in to information and information in to knowledge for the intended purpose of making better decisions faster in the area of drug lead identification and optimization. *F.K. Brown, 1998*

*Chemoinformatics* is the application of informatics methods to solve chemical problems. *J. Gasteiger, 2004*

*Chemoinformatics* is a field dealing with molecular objects (graphs, vectors) in multidimensional chemical space. *A. Varnek & A. Tropsha, 2007.*

# Chemoinformatics: data to knowledge



# What is required to create a good model?

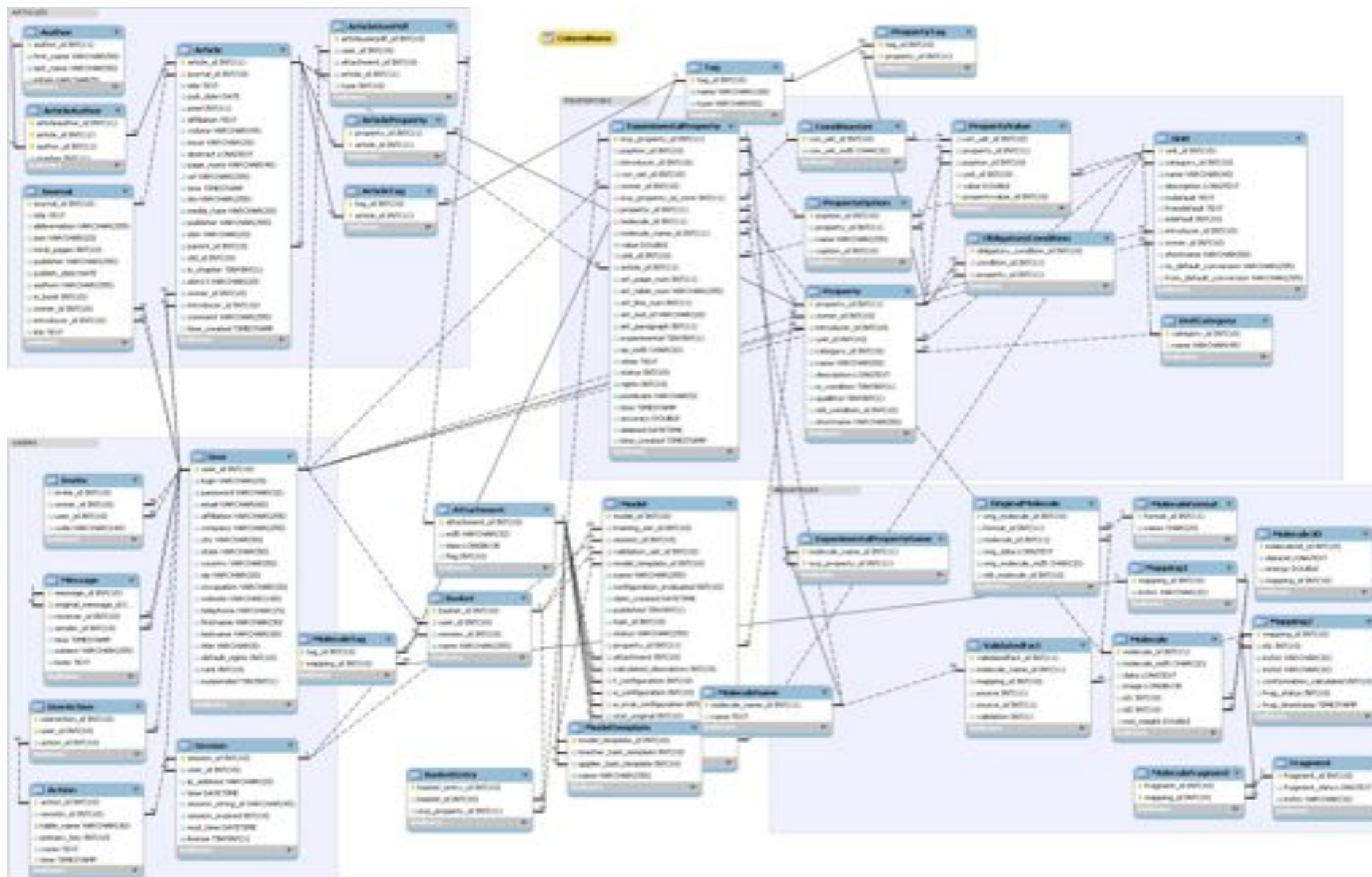
- Data -- the most essential part!
- Appropriate representation of molecules (choice of descriptors is crucial!)
- Good statistical methods



# OCHEM - On-line CHEmical database and Modeling environment









### Compounds properties browser

Search for numerical compounds properties linked to scientific articles

Area of your interest:  
no tags selected [change]

#### SOURCE

Article/Source [select]

Page Table

#### PROPERTY

Activity/Property [select]

No unit filter

#### CONDITIONS

#### MOLECULE

Name / QID / InchiKey

[search by fragment]

#### MISCELLANEOUS

Current set [7]

Show all

Records by introducers:

All users

- ☐ Original records
- ☐ Not validated
- ☐ Error records
- ☐ Error in chies
- ☐ Mismatching names
  - ☐ Include stereochem.
- ☐ Empty molecules
- ☐ Duplicates
- ☐ No stereochemistry

Sort by:

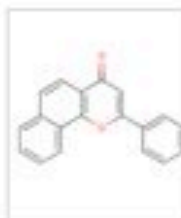
Creation time Asc

REFRESH RESET

Basket Records

1 - 5 of 53861

5 Items on page 1 of 10773 > >>

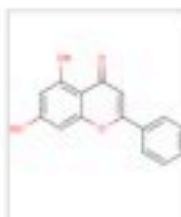


● IC50 CYP450 Inhibition = 0.5  $\mu$ M

Campbell, DR  
Flavonoid inhibition of aromatase enzyme activity  
in human p...  
P: 383  
J. Steroid Biochem. Mol. Biol. 1993; 46 (3) 381-8  
7,8-BF

CYP450 Type = CYP19  
CYP450 etalon reaction = aromatization  
Inhibition mechanism = competitive

14:27, 4 Aug 09  
vkovallshyn 82

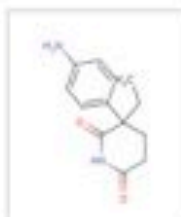


● IC50 CYP450 Inhibition = 4.6  $\mu$ M

Campbell, DR  
Flavonoid inhibition of aromatase enzyme activity  
in human p...  
P: 383  
J. Steroid Biochem. Mol. Biol. 1993; 46 (3) 381-8  
Chrysin

CYP450 Type = CYP19  
CYP450 etalon reaction = aromatization  
Inhibition mechanism = competitive

14:26, 4 Aug 09  
vkovallshyn 82

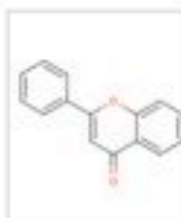


● IC50 CYP450 Inhibition = 7.4  $\mu$ M

Campbell, DR  
Flavonoid inhibition of aromatase enzyme activity  
in human p...  
P: 383  
J. Steroid Biochem. Mol. Biol. 1993; 46 (3) 381-8  
Aminoglutethimide

CYP450 Type = CYP19  
CYP450 etalon reaction = aromatization  
Inhibition mechanism = competitive

14:25, 4 Aug 09  
vkovallshyn 82



● IC50 CYP450 Inhibition = 68.0  $\mu$ M

Campbell, DR  
Flavonoid inhibition of aromatase enzyme activity  
in human p...  
P: 383  
J. Steroid Biochem. Mol. Biol. 1993; 46 (3) 381-8  
Flavone

CYP450 Type = CYP19  
CYP450 etalon reaction = aromatization  
Inhibition mechanism = competitive

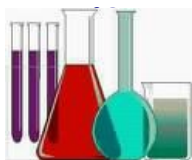
14:20, 4 Aug 09  
vkovallshyn 82

# Database schema

## Simplified overview

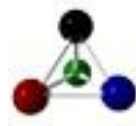
### Property

$\log P = 0.5$   
Melting Point = 100



### Filtering

Toxicology, Biology,  
Partition coefficient.



### Condition

Temperature,  
pH, species,  
tissue, method



### Tags

Toxicology, Biology,  
Partition coefficient.



### Data Point

LD50 oral = 10.0 g/kg  
Species = rat  
Schmeiser, HH  
Evaluation of health risks caused by musk ketone...  
P: 295 T: 1 L: 1  
Int J Hyg Environ Health 2001; 203 (4) 293-299  
Fragrances CADASTER molecules  
81-14-1 ; FRA-043 ; Musk ketone  
09:02, 23 Mar 09  
maraluni

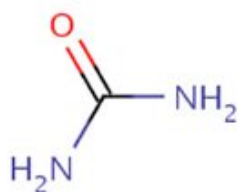
Introducer  
Bill G., Sergey B.

Date of modification  
Informationssystem



### Structure

Benzene. Urea. ...



### Manipulation

Editing  
Organization  
Working sets<



### Article

Garberg, P  
"In vitro models for ..."





# Descriptors

## Model editor

You have to configure the model before you can build it.

### Select descriptors blocks

Please select which types of descriptors do you want to use:

☐ E-state

☒ Dragon

- ☒ constitutional descriptors
- ☒ walk and path counts
- ☒ information indices
- ☒ edge adjacency indices
- ☒ topological charge indices
- ☒ Randic molecular profiles
- ☒ RDF descriptors
- ☒ WHIM descriptors
- ☒ functional group counts
- ☒ charge descriptors

- ☒ topological descriptors
- ☒ connectivity indices
- ☒ 2D autocorrelations
- ☒ BCUT descriptors
- ☒ eigenvalue-based indices
- ☒ geometrical descriptors
- ☒ 3D-MoRSE descriptors
- ☒ GETAWAY descriptors
- ☒ atom-centred fragments
- ☒ molecular properties

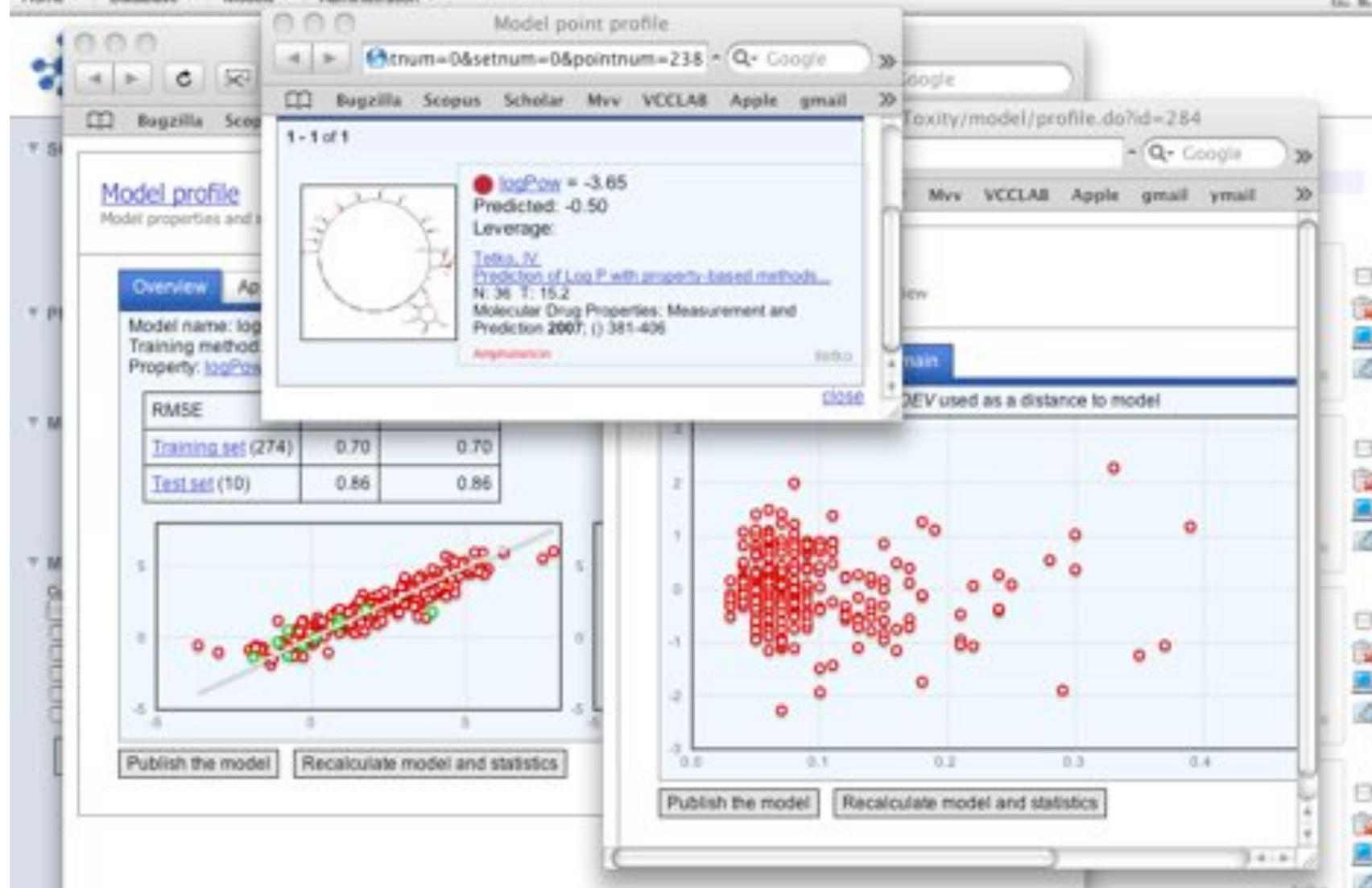
☒ Fragmentors

Fragmentors' length from  to

Type of fragments

<<Back

Next>>



# Offered course

Monday (10.08, 10:00-13:00)

- Introduction to Chemoinformatics
- From 2D to 4D descriptors (Prof. V.E. Kuzmin, OSU, Ukraine)
- Developing chemoinformatics models: One can't embrace unembraceable (Dr I.V. Tetko, HMGU, Germany)



Tuesday (11.08, 10:00-13:00)

- Role of 3D information in chemoinformatics: problem of chirality (Prof. V.E. Kuzmin)
- Introduction to predictive toxicology (Prof. W. Peijnenberg, 2 lectures)



Wednesday (12.08, 10:00-13:00)

- Introduction to OCHEM: database (Mr S. Novotarskyi, HMGU)
- Practical session (Mr S. Novotarskyi & I. Sushko, HMGU)



Thursday (13.08, 10:00-13:00)

- Introduction to OCHEM: models (Mr I. Sushko, HMGU, Germany)
- Overview of machine learning methods available at OCHEM (Dr. I. Tetko)
- Practical session (Mr S. Novotarskyi & I. Sushko, HMGU, Germany)





# HMGU Team



# Acknowledgement

Prof. W. Peijnenburg  
Prof. V. Kuz'min

Mr I. Sushko  
Mr S. Novotarskyi

## Team members

Mr A.K. Pandey  
Mr R. Körner  
Mr S. Brandmayer  
Dr. M. Rupp

BMBF GO-Bio grant 0313883, CADASTER FP7 EU project,  
Germany-Ukraine bilateral grant

Thank you for your attention!