

# CHEMOINFORMATICS: PRECURSORS FOR “AI IN MATERIALS”

---

Logan Ward  
Asst. Computational Scientist  
Argonne National Laboratory

20 January 2022

# Historical Perspective

“The application of informatics methods to solve chemical problems” – T. Engel (2003)

**“Chemoinformatics” is “information technology” and chemistry**

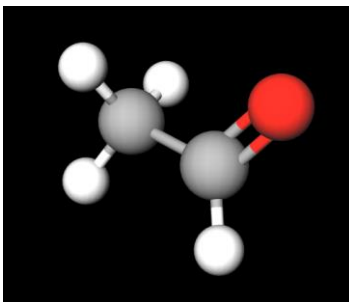
- Chemistry and computing have a 60+ year history
- Ex: Paper on [typing notation for chemicals](#) (1952 Wiswesser)
  - “to sort and list such information with standard tabulating machine or with the new IBM scanning machinery”
  - “The notation in fact has made possible the establishment of the Willson Toxicity Registry, a new activity to catalogue and correlate mammalian toxicity data”

list:	
System	Symbols per Compound
Dyson	27.45
Gruber	29.00 .
Silk	19.36
Wiswesser	12.93

Modern uses of AI techniques in materials follow trends in “chemoinformatics”

# Molecular Descriptors

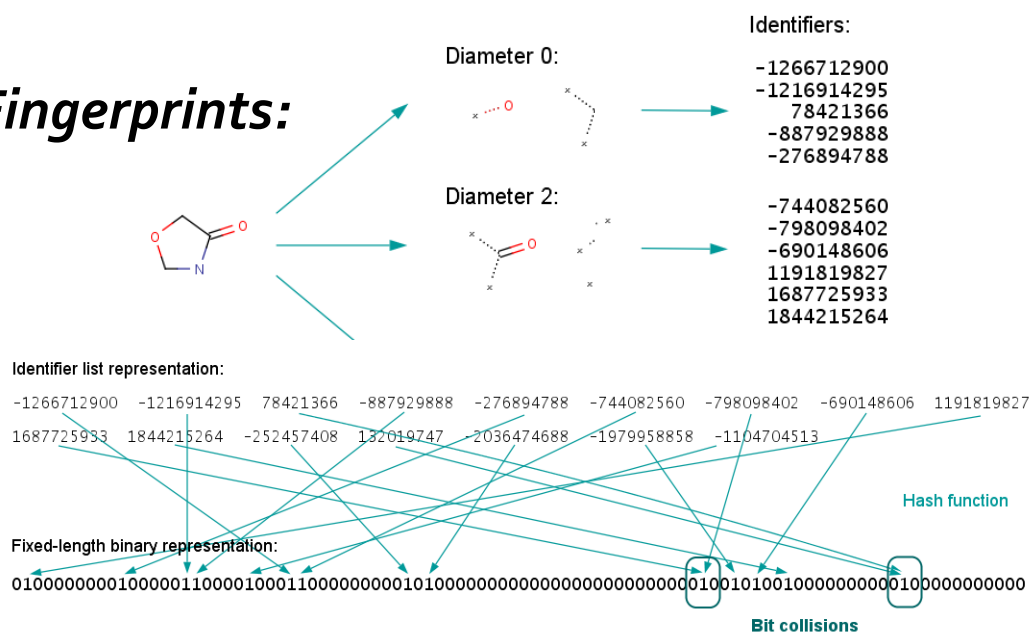
## Discriminative



### Line-notations for structure:

- SMILES (ex: "CC=O")
- InChI (ex: InChI=1S/C2H4O/c1-2-3/h2H,1H3)

### Fingerprints:



## Descriptive

### Extremely Well-Studied



### Handbook of Molecular Descriptors

Author(s): Prof. Dr. Roberto Todeschini, Dr. Viviana Consonni

First published: 22 September 2000

Print ISBN: 9783527299133 | Online ISBN: 9783527613106 | DOI: 10.1002/9783527613106

Copyright © 2000 WILEY-VCH Verlag GmbH

Book Series: Methods and Principles in Medicinal Chemistry

[Free Access](#)

### Bibliography (Pages: 524-667)

[First Page](#) | [PDF](#) | [References](#) | [Request permissions](#)

### Many types of descriptors:

1. Constitutional (ex: "how many Ns?")
2. Structural (ex: Solvent-Accessible Surface Area)
3. Quantum-chemical (ex: partial charges)

# QSAR: “Quantitative Structure-Activity Relationships”

[CONTRIBUTION FROM THE DEPARTMENT OF CHEMISTRY, POMONA COLLEGE, CLAREMONT, CALIFORNIA]

## $\rho$ - $\sigma$ - $\pi$ Analysis. A Method for the Correlation of Biological Activity and Chemical Structure

BY CORWIN HANSCH AND TOSHIO FUJITA<sup>1</sup>

RECEIVED AUGUST 19, 1963

Using the substituent constant,  $\sigma$ , and a substituent constant,  $\pi$ , defined as  $\pi = \log P_X - \log P_H$  ( $P_H$  is the partition coefficient of a parent compound and  $P_X$  that of a derivative), regression analyses have been made of the effect of substituents on the biological activity of benzoic acids on mosquito larvae, phenols on gram-positive and gram-negative bacteria, phenyl ethyl phosphate insecticides on houseflies, thyroxine derivatives on rodents, diethylaminoethyl benzoates on guinea pigs, and carcinogenic compounds on mice.

Early examples of using “data” + “statistics” date to the 1960s!

General ingredients have not changed:

1. Trusted chemical data resource
2. Informative chemical descriptors
3. Appropriate regression algorithm

$$\log \frac{1}{C} = 0.519\pi + 1.540; \begin{matrix} r^2 & r & s^{12} \\ 0.955 & 0.977 & 0.130 \end{matrix} \quad (13)$$

# QSAR is still a tool you need to know

## SCIENTIFIC REPORTS

OPEN

### Quantitative design rules for protein-resistant surface coatings using machine learning

Received: 10 September 2018

Accepted: 22 November 2018

Published online: 22 January 2019

Tu C. Le<sup>1</sup>, Matthew Penna<sup>1,2</sup>, David A. Winkler<sup>3,4,5,6</sup> & Irene Yarovsky<sup>1,2</sup>

Preventing biological contamination (biofouling) is key to successful development of novel surface and nanoparticle-based technologies in the manufacturing industry and biomedicine. Protein adsorption is a crucial mediator of the interactions at the bio – nano – materials interface but is not well understood.

Science has created better descriptors, better supervised learning algorithms, but the basic idea is still the same.

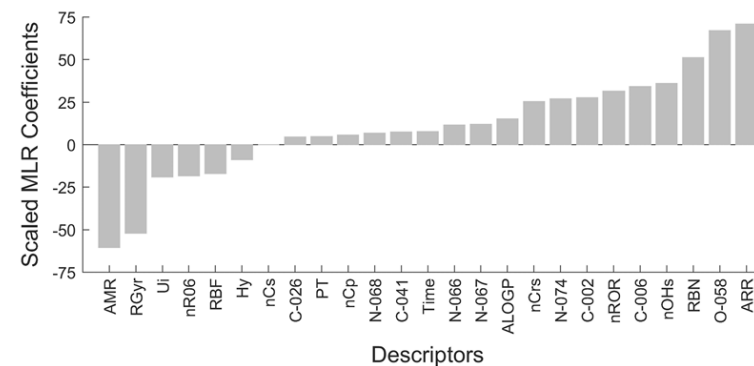


Figure 5. Scaled MLR coefficients of the most relevant descriptors selected from the pool 67 descriptors.

# (Shallow) Machine Learning and QSPR

**Table 2c.** Melting-Point QSPR for Nitrocyanamide Salts

Eq	N	r <sup>2</sup>	q <sup>2</sup>	F	s <sup>2</sup>	Term	Coefficient	t-test	Descriptor
4	7	0.960	0.894	120	3.71	0	$-1.12 \times 10^2 (\pm 1.03 \times 10^1)$	-10.9	Intercept
						1	-7.11 ( $\pm 0.649$ )	-10.9	HOMO <sub>subst</sub> : Highest occupied molecular orbital located on the substituent.

Source: [Trohalaki, Pachter. QSAR & Comb. Sci. \(2005\)](#)

Sophisticated methods still useful,  
but don't discount linear regression!

## Why?

- Small dataset sizes
- Informative descriptors
- Interpretability
- Limited computational costs



# Deep Learning is not Cure-All

Table 3: Summary of performances(test subset): conventional methods versus graph-based methods. Graph-based models outperform conventional methods on 11/17 datasets.

Category	Dataset	Metric	Best performances - conventional methods	Best performances - graph-based methods
Quantum Mechanics	QM7	MAE	KRR(CM): 10.22	<b>DTNN: 8.75</b>
	QM7b	MAE	KRR(CM): 1.05	<b>DTNN: 1.77*</b>
			Multitask: 0.0150	<b>MPNN: 0.0143</b>
			Multitask(CM): 4.35	<b>DTNN: 2.35</b>
			XGBoost: 0.99	<b>MPNN: 0.58</b>
			XGBoost: 1.74	<b>MPNN: 1.15</b>
			XGBoost: 0.799	<b>GC: 0.655</b>
Biophysics		AUC-PRC	Logreg: 0.129	<b>GC: 0.136</b>
	MUV	AUC-PRC	<b>Multitask: 0.184</b>	Weave: 0.109
	HIV	AUC-ROC	<b>KernelSVM: 0.792</b>	GC: 0.763
	BACE	AUC-ROC	<b>RF: 0.867</b>	Weave: 0.806
Physiology	PDBbind(full)	RMSE	<b>RF(grid): 1.25</b>	GC: 1.44
	BBBP	AUC-ROC	<b>KernelSVM: 0.729</b>	GC: 0.690
	Tox21	AUC-ROC	KernelSVM: 0.822	<b>GC: 0.829</b>
	ToxCast	AUC-ROC	Multitask: 0.702	<b>Weave: 0.742</b>
	SIDER	AUC-ROC	<b>RF: 0.684</b>	GC: 0.638
	ClinTox	AUC-ROC	Bypass: 0.827	<b>Weave: 0.832</b>

\* As discussed in section 4.4, DTNN outperforms KRR(CM) on 14/16 tasks in QM7b while the mean-MAE is skewed due to different magnitudes of labels.

For benchmark problems of learning from molecular data, conventional ML better for **6/17** cases

# An Example: Computational Toxicology

<https://ntp.niehs.nih.gov/whatwestudy/niceatm/comptox/ct-opera/opera.html>

Home » What We Study » NICEATM: Alternative Methods » Computational Toxicology » OPERA

## Computational Toxicology

Adverse Outcome Pathways

Computational Models of Chemical Activity

ICE: Integrated Chemical Environment

In Vitro to In Vivo Extrapolation

Integrated Approaches to Testing and Assessment

OPERA

## OPERA

### Open Structure-activity/property Relationship App

Quantitative structure-activity/property relationship (QSAR/QSPR) models provide predictions of chemical activity that can augment non-animal approaches for predicting toxicity. However, the performance of QSAR models highly depends on the quality of the data and modeling methodologies used.

To provide robust QSAR/QSPR models for chemical properties of environmental interest that can be used for regulatory purposes, EPA NCCT created the Open Structure-activity/property Relationship App (OPERA) (Mansouri et al. 2018 [↗](#)). OPERA is a free and open-source/open-data suite of QSAR models providing predictions for physicochemical properties, environmental fate parameters, and toxicity endpoints. All OPERA models were built on curated data and QSAR-ready chemical structures standardized using an open-source workflow (Mansouri et al. 2016 [↗](#)).

OPERA is an ongoing collaboration between [NICEATM](#) and EPA. Recent additions to OPERA include predictions for:

SHARE THIS:       1  
<https://ntp.niehs.nih.gov/go/opera> 



**QSAR is not going away any time soon.**

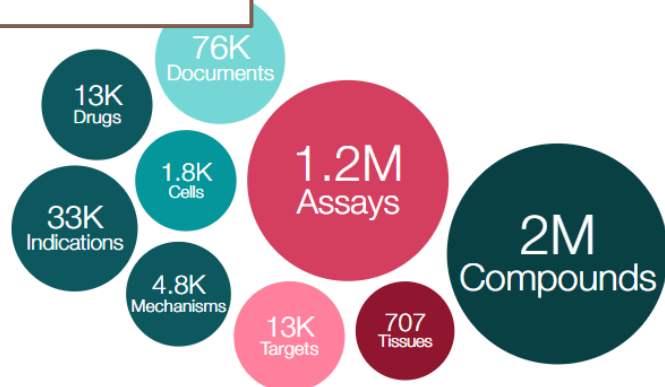


# PRACTICAL CONSIDERATIONS: WHAT TOOLS ARE OUT THERE?

---

# Chemical Databases Abound

ChEMBL



<https://www.ebi.ac.uk/chembl/>



<http://www.chemspider.com/>

Chemicalize Calculation Search Drawing Batch

Enter a molecule name, registry number, SMILES, or InChI (e.g. niacin)

**Structures**

**Results**

Input	Isopiperitenone
Molar mass	150.221 g/mol
Exact mass	150.104465071 Da
Formula	C <sub>10</sub> H <sub>14</sub> O
Composition	C (79.96%), H (9.39%), O (10.65%)
Lipinski's rule of five	✓

<https://chemicalize.com/>

TOOL

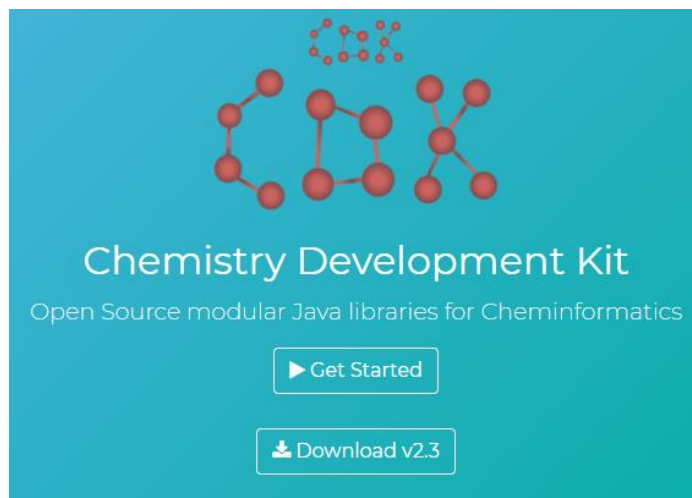
## UL Cheminformatics Tool Kit

With a curated database of 70 million structures and 80,908 chemicals with 833,844 labeled hazard endpoints, our digital solution utilizes an advanced algorithm, machine learning, and analysis of millions of chemical combinations to predict chemical hazards.

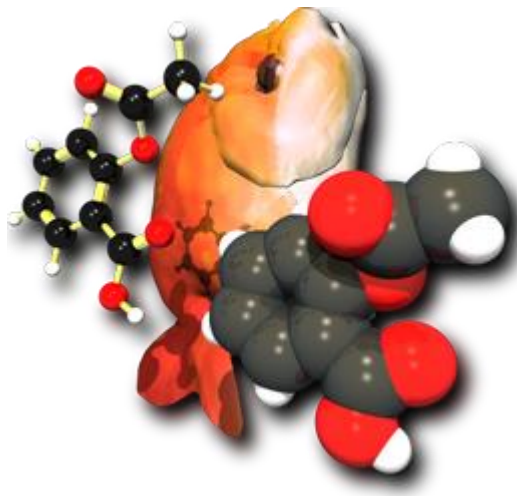
[Learn more](#)

<https://www.ul.com/resources/apps/ul-cheminformatics-tool-kit>

# Toolkits for Building QSAR Models



<https://cdk.github.io/>



[http://openbabel.org/wiki/Main\\_Page](http://openbabel.org/wiki/Main_Page)

README.md

## RDKit

Azure Pipelines succeeded docs passing DOI 10.5281/zenodo.4318717

RDKit is a collection of cheminformatics and machine-learning software written in C++ and Python.

- [BSD license](#) - a business friendly license for open source
- Core data structures and algorithms in C++
- [Python 3.x wrapper](#) generated using Boost.Python
- Java and C# wrappers generated with SWIG
- 2D and 3D molecular operations

<https://rdkit.org>



<https://chm.kode-solutions.net/>

# Final Note: Chemoinformatics is a Mature Field

There is still a lot materials informatics can learn from

- 60+ years of history
- [Textbooks](#), [more textbooks](#), codes and review papers
- Decades of “lessons learned”



## Handbook of Molecular Descriptors

Author(s): Prof. Dr. Roberto Todeschini, Dr. Viviana Consonni

First published: 22 September 2000

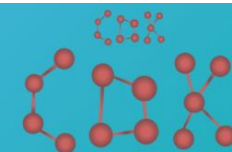
Print ISBN: 9783527299133 | Online ISBN: 9783527613106 | DOI: 10.1002/9783527

Copyright © 2000 WILEY-VCH Verlag GmbH

Book Series: Methods and Principles in Medicinal Chemistry

Table 3. Common QSPR modelling pitfalls and methods of avoiding them

pitfall	recommendation to minimize or avoid
use of uninformative descriptors	use descriptors that are related to the molecular structure where possible, use virtual screening methods when complex descriptors are necessary, and develop new materials descriptors
overfitting, and grossly underdetermined systems	reduce size of descriptor pool before building models, <sup>24</sup> monitor number of fitted parameters (descriptor weights or neural network weights) to ensure they are substantially less than the number of experiments, and check that training and test set statistics are similar
descriptor selection and chance correlations	use Topliss criteria <sup>22,23</sup> to estimate probability of chance correlations and descriptor scrambling; avoid methods where repeated sampling of a larger pool of descriptors is done to obtain the optimum subset of descriptors; and use sparse, context-dependent feature-selection methods <sup>17,18</sup>
modeling complex, nonlinear structure–property relationships	avoid overly complex nonlinear models, compare nonlinear model statistics with linear models, and use regularizing methods that attempt to optimize model complexity <sup>16,17</sup>
validating QSPR models	synthesize new materials that models predict to be superior and test if feasible, use independent test sets to assess model predictivity otherwise, and employ cross-validation methods with caution <sup>27,28</sup>
domain of applicability of models	calculate the range of all descriptors used to develop the model, <sup>29–32</sup> avoid extrapolations using descriptor space distant from that used in model, and use probabilistic modeling methods (e.g., Bayesian regularization <sup>16</sup> ) that allow estimation of likely prediction error
incorrect handling of outliers	avoid removing outliers wherever possible, check whether outlier lies well within domain before removing it, remove outliers sparingly and describe why they were omitted, and retest properties for outliers to eliminate measurement or transcription errors



**Chemistry Development Kit**

Open Source modular Java libraries for Cheminformatics

[▶ Get Started](#)

[Download v2.3](#)

Good Refs: [T. Le et al. Chem. Rev. \(2012\), 2889](#)

# Practical Exercises

---

First for the course!

[WardLT/applied-ai-for-materials:molecular-property-prediction/chemoinformatics](#)

Learn how to:

- Manipulate chemical data with RDKit and Pandas
- Train conventional machine learning models with *descriptors* and *fingerprints*