

Institute of Structural Biology

OCHEM, openOCHEM and beyond: use of Advanced machine learning (AI) in Chemistry

lgor V. Tetko

Helmholtz Munich and BIGCHEM GmbH

University of Münster, 12 December 2023

HELMHOLTZ MUNICI



Agenda

- Overview of HMGU
- Computational predictions in drug discovery
- OCHEM
- Examples of models
- Consensus modelling
- Multitask learning
- **Representation learning**
- **Reaction predictions**
- Conclusions

Helmholtz Association – Facts and Figures

Germany's largest research organization 19 research centers ٠ Budget: 5 Billion €, more than 42.000 staff . **6 Research Fields** HEALTH AERONAUTICS. EARTH AND MATTER KEY ENVIRONMENT Research Center Jülich SPACE AND TECHNOLOGIES TRANSPORT (FUTURE: INFORMATION) CISPA

6 centers represent the Research Field Health.





BIGCHEM GmbH is a spin-off of the center

Hypothesized Target for African sleeping sickness



Neufeld, C. ... Sattler, M. (2009) EMBO J. 28: 745-754

Peroxins Pex14/Pex5 are responsible for transport of glycosomal enzymes from cytoplasm to glycosomes for glucose metabolism



Dawidowski, M., ... Popowicz, G. M. Science 2017, 355, 1416-1420.

Traditional Process of Drug Discovery



• Profiling and screening in the virtual space helps to identify the most promising candidates

ADMETox filters in Bayer

	Insufficient quality	First approach	Medium model	Good n	nodel	Robust model			
	Endpoint	Model type	Data se	t size	2005	2009	2014	2019	Retraining
	Caco-2 permeation	C (N)	>10 0	00			RF	SVR	Weekly
A bsorption	Caco-2 efflux	C (N)	>10 0	00			RF	SVR	Weekly
	Bioavailability (rat)	С	~200	0				RF	On demand
Distribution	Human serum albumin	N	>30 0	00			PLS	MTNN	On demand
Distribution	Fraction unbound	N	>100	0			PLS	MTNN	On demand
	Microsomal stability (hum)	C (N)	>10 0	00			RF	RF	Weekly
Matabaliana	Microsomal stability (mouse)	C (N)	>10 0	00			RF	RF	Weekly
Metabolism	Microsomal stability (rat)	C (N)	>10 0	00			RF	RF	Weekly
	Hepatocyte stability (rat)	C (N)	>30 0	00			del 2014 2019 Retra RF SVR Weet PLS MTNN On det PLS MTNN On det RF RF Weet RF RF On det RF RF On det RF RF On det SVM SVM On det SVM SVM On det MTNN On det MTNN PLS MTNN On det PLS MTNN On det PLS MTNN On det PLS MTNN On det Indiation Indiation Indiation Indiatin <td< td=""><td>Weekly</td></td<>	Weekly	
	hERG inhibition	С	>10 0	00			RF	SVM	Weekly
	Ames mutagenicity	С	>10 0	00			RF	RF	On demand
Toxicity	CYP inhibition isoforms	С	>10 0	00			RF	RF	On demand
	Phospholipidosis	С	<100	0			SVM	SVM	On demand
	Structure filter tool	Score	n.a		-	-	-	-	On demand
	Solubility (DMSO)	N	>30 ,0	00			DL O	MTNN	On demand
	Solubility (Powder)	N	<10 0	00			PLS	MTNN	On demand
	logD @ pH 7.5	N	>70 0	00			PLS	MTNN	On demand
PhysChem	Membrane affinity	N	<10 0	00			PLS	MTNN	On demand
	рКа	N	>10 0	00			ANN	ANN	On demand
	Oral PhysChem score	Score	n.a		-	-	-	-	On demand
	i.v. PhysChem score	Score	n.a		-	-	-	-	On demand

Göller, A.H. et al Drug Discov. Today 2020, 25 (9), 1702-1709.

OCHEM https://ochem.eu



Home - Database -Models -

Welcome to OCHEM! Your possible actions

Explore OCHEM data

Search chemical and biological data: experimentally measured, published and exposed to public access by our users. You can also upload your data.

Create OSAR models

Build QSAR models for predictions of chemical properties. The models can be based on the experimental data published in our database.

Run predictions

Apply one of the available models to predict property you are interested in for your set of compounds.

Screen compounds with ToxAlerts

Screen your compound libraries against structural alerts for such endpoints as mutagenicity, skin sensitization, aqueous toxicity, etc.

Tutorials

Check our video tutorials to know more about the OCHEM features.

Our acknowledgements

Feedback and help

User's manual Check an online user's manual



v.4.2.282

👮 log in create account

A+ a- Privacy statement



OCHEM Database schema

Properties

Modeling iterative workflow

Representation of chemical structures

Definition of molecular descriptors

The molecular descriptor is the final result of a logic and mathematical procedure which transforms chemical information encoded within a symbolic representation of a molecule into a useful number, or the result of some standardized experiment.

Roberto Todeschini & Viviana Consonni

Molecular Descriptors for Chemoinformatics

WILEY-VCH

Roberto Todeschini, Viviana Consonni

Second, Revised and Enlarged Edition Volume 1: Alphabetical Listing Volume 1: Monte of the second second

Examples of descriptors

✓ alvaDesc v.2.0.4 (5666/3D)

[select all] [select none] [select 3D] [unselect 3D]

- Constitutional descriptors (50)
- ✓ Topological indices (79)
- Connectivity indices (37)
- 2D matrix-based descriptors (608)
- Burden eigenvalues (96)
- ETA indices (40)
- Geometrical descriptors (3D, 38)
- ✓ 3D autocorrelations (3D, 80)
- ✓ 3D-MoRSE descriptors (3D, 224)
- GETAWAY descriptors (3D, 273)
- Functional group counts (3D, 154)
- ✓ Atom-type E-state indices (346)
- 2D Atom Pairs (1596)
- Charge descriptors (3D, 15)
- ✓ Drug-like indices (30)
- WHALES (3D, 33)
- Chirality (70)

- Ring descriptors (35)
- Walk and path counts (46)
- Information indices (51)
- ✓ 2D autocorrelations (213)
- P_VSA-like descriptors (69)
- Edge adjacency indices (324)
- ✓ 3D matrix-based descriptors (3D, 132)
- RDF descriptors (3D, 210)
- ✓ WHIM descriptors (3D, 114)
- Randic molecular profiles (3D, 41)
- Atom-centred fragments (115)
- Pharmacophore descriptors (165)
- ✓ 3D Atom Pairs (3D, 36)
- Molecular properties (3D, 27)
- CATS 3D (3D, 300)
- **MDE** (19)

QSPR/QSAR modelling in OCHEM

Select the molecular descriptors 🌒		Create a model 🕤 Select the training and validation sets, the machine learning method and the validation protocol
Recommended descriptor types (2D)	Predictions by OCHEM's featured models 🕕	
✓ OEState	Ames levenberg	
 Bonds Indices Counts only 	Toxicity against T. Pyriformis ALogPS 3.0 CYP142 Estate+ALogPS	Select the training and validation sets:
 ALogPS (2) Mold2 (777) CDDD JPlogP SIRMS ISIDA fragments The in Hashed Atom Pair fingerprint (MAP4) GSFragment (1138) QNPR Multilevel Neighborhoods of Atoms (MNA) Structural alerts (ToxAlerts and Functinal Groups) Recommended descriptor types (3D) alvaDesc v.2.0.4 (5666/3D) Dragon v. 7 (5270/3D) CDK 2.7.1 descriptors (256/3D) Chemaxon descriptors (499/3D) RDKit descriptors (32) MORDRED descriptors (1826/3D) MORDRED descriptors (1826/3D) KrakenX descriptors (1826/3D) MPAC2016 descriptors (18251/3D) MERA descriptors (529/3D) MERA descriptors (529/3D) MERA descriptors (529/3D) MERAY descriptors (42/3D) MERAY descriptors (529/3D) 	 CYP2C9 Estate+ALogPS CYP2C19 Estate+ALogPS CYP2C19 Estate+ALogPS CYP2AE Estate+ALogPS CYP2AE Estate+ALogPS Pyrolysis point prediction (best Estate) Melting Point prediction (best Estate) Water solubility model based on logP and Melti ALOGPS 2.1 logP ALOGPS 2.1 logS Outputs of other OCHEM models Obsolete/Additional descriptor types CDK 2.0 descriptors (256/3D) CDK 1.4.11 descriptors (256/3D) E-state Dragon v. 5.4 (1644/3D) Dragon v. 5.6 (3224/3D) Dragon v. 6 (4885/3D) MOPAC 7.1 descriptors (25/3D) 	Training set (required): peptidesregr [details] Add a validation set The model will predict this property: LogD using unit: Log unit Skip model configuration and use the predefined settings Choose the learning method: Suggested modeling method: ASNN: Associative Neural Networks doi:10.1007/978-1-60327-101-1_10 (New) Attentive FP doi: 10.1021/acs.jmedchem.9b00959 ChemProp MPNN for property prediction (GPU) doi:10.1007/978-3-030-30493-5_79 ChemProp MPNN for property prediction (GPU) doi:10.1007/978-3-030-30493-5_79 Transformer-CNF model Consensus model (based on models developed for the same set) DEEPCHEM: several methods from DeepChem (GPU) arXiv:1703.00564 (New) DIMENET - Directional Message Passing Neural Network arXiv:2003.03123 Deep Learning Consensus Architecture (DLCA) doi:10.1021/acs.jcim.9b00526 DNN: Deep Neural Network (GPU) doi:10.1021/acs.jcim.8b00685 EAGCNG - Edge Attention based Multi-relational Graph Convolutional Networks (GPU) arXiv:1802.04944 FSMLR: Fast Stagewise Multiple Linear Regression doi:10.1134/S0012500807120026 GNN- Graph Semorphism Network (GPU) arXiv:1910.13124
Spectrophores (144/3D)		 KNN: k - Nearest Neighbors KPLS - Kernel Partial Least Squares doi:10.1109/IJCNN.2006.246832 LibS/W. arid exacts bacteria patientical doi:10.1146/1961190.1961190
Special descriptors (scaffolds, fingerprints):		O LSSVMG: Least Squares Support Vector Machine (GPU) doi:10.11023/A:1018628609742
Chemaxon Scaffolds Silicos-It Scaffolds ECFP Fingerprints MolPrint Fingerprints		 MLR: Multiple Linear Regression PLS: Partial Least Squares doi:10.1016/S0169-7439(01)00155-1 RFR: Random Forest regression and classification doi:10.1023/A:1010933404324 Transformer-CNN - Transformer Convolutional Neural Network (GPU) doi:10.1186/s13321-020-00423-w Transformer-CNNi - faster Transformer-CNN (GPU) doi:10.1186/s12321-020-00423-w WERA- IAB: Werks C4 5 deriving register on put classification are with banging doi:10.11145/1656274.1656278
Conditions of experiments		WEKA-RF: Random Forest, only classification doi:10.1023/4:1010933404324
□ pH □ Ionisable		O XGBoost: Scalable and Flexible Gradient Boosting doi:10.1145/2939672.2939785

Model validation

Validation method: N-Fold cross-validation -

Number of folds: 5 Stratified cross-validation (classification only) Treat each record as a new molecule 0

You can create a model from template: import an XML model template or use another model as a template

Examples of recent studies using OCHEM

What Features of Ligands Are Relevant to the Opening of **Cryptic Pockets in Drug Targets?**

Zhonghua Xia ¹⁽⁰⁾, Pavel Karpov ¹, Grzegorz Popowicz ¹⁽⁰⁾, Michael Sattler ^{1,2} and Igor V. Tetko ^{1,3,*}

Article

More Is Not Always Better: Local Models Provide Accurate **Predictions of Spectral Properties of Porphyrins**

Aleksey I. Rusanov¹, Olga A. Dmitrieva¹, Nugzar Zh. Mamardashvili¹ and Igor V. Tetko^{1,2,3,*}

Highly Accurate Filters to Flag Frequent Hitters in AlphaScreen Assays by Suggesting their Mechanism

Dipan Ghosh,^[a] Uwe Koch,^[a] Kamyar Hadian,^[b] Michael Sattler,^[c, d] and Igor V. Tetko^{*[d, e, f]}

CATMoS: Collaborative Acute Toxicity Modeling Suite

Kamel Mansouri,^{1,41} Agnes L. Karmaus,¹ Jeremy Fitzpatrick,² Grace Patlewicz,³ Prachi Pradeep,^{3,4} Domenico Alberga,⁵ Nathalie Alepee,⁶ Timothy E.H. Allen, ⁷ Dave Allen,¹ Vinicius M. Alves,^{8,9} Carolina H. Andrade,⁹ Tyler R. Auernhammer,¹⁰ Davide Ballabio,¹¹ Shannon Bell,¹ Emilio Benfenati,¹² Sudin Bhattacharya,¹³ Joyce V. Bastos,⁹ Stephen Boyd,¹⁴ J.B. Brown,¹⁵ Stephen J. Capuzzi,⁸ Yaroslav Chushak,^{16,17} Heather Ciallella,¹⁸ Alex M. Clark,¹⁹ Viviana Consonni,¹¹ Pankaj R. Daga,²⁰ Sean Ekins,¹⁹ Sherif Farag,⁸ Maxim Fedorov,²¹ Denis Fourches,^{22,23} Domenico Gadaleta,¹² Feng Gao,¹⁴ Jeffery M. Gearhart,^{16,17} Garett Goh,²⁴ Jonathan M. Goodman,⁷ Francesca Grisoni,¹¹ Christopher M. Grulke,³ Thomas Hartung.²⁵ Matthew Hirn.²⁶ Pavel Karpov.²⁷ Alexandru Korotcov.²⁸ Giovanna J. Lavado.¹² Michael Lawless.²⁰ Xinhao Li,²² Thomas Luechtefeld,²⁵ Filippo Lunghini,²⁹ Giuseppe F. Mangiatordi,⁵ Gilles Marcou,²⁹ Dan Marsh,²⁵ Todd Martin,³⁰ Andrea Mauri,³¹ Eugene N. Muratov,^{8,9} Glenn J. Myatt,³² Dac-Trung Nguyen,³³ Orazio Nicolotti,⁵ Reine Note,⁶ Paritosh Pande,²⁴ Amanda K. Parks,¹⁰ Tyler Peryea,³³ Ahsan H. Polash,¹⁵ Robert Rallo,²⁴ Alessandra Roncaglioni,¹² Craig Rowlands,²⁵ Patricia Ruiz,³⁴ Daniel P. Russo,¹⁸ Ahmed Sayed,³⁵ Risa Sayre,³⁴ Timothy Sheils,³³ Charles Siegel,²⁴ Arthur C. Silva,⁹ Anton Simeonov,³³ Sergey Sosnin,²¹ Noel Southall,³³ Judy Strickland,¹ Yun Tang,³⁶ Brian Teppen,¹⁴ Igor V. Tetko.^{27,37} Dennis Thomas.²⁴ Valery Tkachenko.²⁸ Roberto Todeschini,¹¹ Cosimo Toma.¹² Ignacio Tripodi.³⁸ Daniela Trisciuzzi,⁵ Alexander Tropsha,⁸ Alexandre Varnek,²⁹ Kristijan Vukovic,¹² Zhongyu Wang,³⁹ Liguo Wang,³⁹ Katrina M. Waters,²⁴ Andrew J. Wedlake,⁷ Sanjeeva J. Wijevesakere,¹⁰ Dan Wilson,¹⁰ Zijun Xiao,³⁹ Hongbin Yang,³⁶ Gergely Zahoranszky-Kohalmi,³³ Alexey V. Zakharov,³³ Fagen F. Zhang,¹⁰ Zhen Zhang,⁴⁰ Tongan Zhao,³³ Hao Zhu,¹⁸ Kimberley M. Zorn,¹⁹ Warren Casey,⁴¹ and Nicole C. Kleinstreuer⁴¹ 18

¹Integrated Laboratory Systems, LLC, Morrisville, North Carolina, USA

²ScitoVation, Research Triangle Park, North Carolina, USA

³Center for Computational Toxicology and Exposure, U.S. Environmental Protection Agency, Research Triangle Park, North Carolina, USA ⁴Oak Ridge Institute for Science and Education (ORISE) Research Participation Program, U.S. Environmental Protection Agency, Research Triangle Park, North Carolina, USA

⁵Dipartimento di Farmacia-Scienze del Farmaco, Università degli Studi di Bari "Aldo Moro", Bari, Italy

⁶I 'Oréal Research & Innovation Aulnay-sous-Bois France

New low-molecular bioregulators as effective agents against multi-drug resistant Acinetobacter baumannii clinical isolate

Received: 18 October 2019	Revised: 27 January 2020	Accepted: 3 March 2020
DOI: 10.1111/cbdd.13678		

RESEARCH ARTICLE

WILEY

In silico and in vitro studies of a number PILs as new antibacterials against MDR clinical isolate Acinetobacter haumannii

Maria M. Trush¹ | Vasyl Kovalishyn¹ | Diana Hodyna¹ | Olexandr V. Golovchenko¹ Svitlana Chumachenko¹ | Igor V. Tetko^{2,3} | Volodymyr S. Brovarets¹ | Larysa Metelytsia¹

Figure 1. Inhibition zone diameters of the six studied PILs (content on a disk, 1.25 µmoles) of an MDR clinical isolate of A. baumannii on agar plates.

Article

Theoretical and Experimental Studies of Phosphonium Ionic Liquids as Potential Antibacterials of MDR Acinetobacter baumannii

Larysa O. Metelytsia¹, Diana M. Hodyna¹, Ivan V. Semenyuta¹, Vasyl V. Kovalishyn¹, Sergiy P. Rogalsky¹, Kateryna Yu Derevianko¹, Volodymyr S. Brovarets¹ and Igor V. Tetko^{2,3,*}

Overview Applicability domain

Model name: M4 Consensus AcinBaum Class - 334799 , published in In silico and in vitro studies of a number PILs as new antibacterials against MDR clinical isolate Acinetobacter baumannii Public ID is 783

Predicted property: AcinBaum Class modeled in CLASS Training method: Consensus

Data Set	#	Accuracy	Balanced Accuracy	MCC	AUC
• Training set: A_Baumanii_Set1 (training)	210 records	83% ± 2.0	82% ± 3.0	0.66 ± 0.05	0.9 ± 0.02
• Test set: A_Baumanii_Set1 (test) [x]	53 records	83% ± 5.0	83% ± 5.0	0.7 ± 0.1	0.92 ± 0.04

Show ROC curves

$Real{\downarrow}/Predicted{\rightarrow}$	inactive	active	Hit rate	$Real{\downarrow}/Predicted{\rightarrow}$	inactive
inactive	69	21	0.77	inactive	22
active	14	106	0.88	active	5
Precision	0.83	0.83		Precision	0.81
Training (Original)				Tes	st (Original)

Hit rate

0.85

0.81

active 4

22

0.85

Contents lists available at ScienceDirect

Journal of Molecular Liquids

journal homepage: www.elsevier.com/locate/molliq

Check fo

updates

20

Beware of proper validation of models for ionic Liquids!

D.M. Makarov^{a,*}, Yu.A. Fadeeva^a, L.E. Shmukler^a, I.V. Tetko^{a,b,c}

^a G. A. Krestov Institute of Solution Chemistry of the Russian Academy of Sciences, Ivanovo, Russia

^b Institute of Structural Biology, Helmholtz Zentrum München-Research Center for Environmental Health (GmbH), Ingolstädter Landstraße 1, 85764 Neuherberg, Germany ^c BIGCHEM GmbH, Valerystr. 49, 85716 Unterschleißheim, Germany

dataset)

Title:	Beware of proper validation of models for lon	ic Liquids!
Authors:	Makarov, D.M.; Fadeeva, Yu.A.; Shmukler, L.	E.; Tetko, I.V.;
Journal reference:	Journal of Molecular Liquids, 2021; 344 (); 11772	https://ochem.eu/article/135195
Internal identifier:	A135195	
Data and mode	ls	
This article is refere Melting Point for le	enced from 249 experimental records This article is onic Liquids (TRANSNNI)	connected to 3 predictive model(s) : trained using the dataset MP ILs FULL set (
Melting Point for le descriptors)	onic Liquids (RFR based on KrakenX	trained using the dataset MP ILs FULL set (Zview dataset profile or export the dataset)
Melting Point for le	onic Liquids (RFR, based on ECFP4 descriptors)	trained using the dataset MP ILs FULL set (Zview dataset profile or export the dataset) validated using dataset MP test set Version2 (Zview dataset profile or export the

Open article editor Attach files to this article

OCHEM modelling

- Comprehensive modeling
- Multitask learning (up to 100 properties)
- >20 descriptors blocks
- GPU + CPU modern methods Supports models
 - >1,000,000 compounds
 - >1,000 servers
 - up to 1GB in size (Java limit)
- Model private/publishing
- Conditions, external descriptors
- ToxAlerts
- Consensus models

Predicted property: LLNA skin sensitization Training set: TRAINING-SARpy-SKIN-SENS-giugno20 OK.xlsx

letrics AUC	ᅌ for Training set	ᅌ Valida	tion: Cros	s-Validat	ion (84 m
		LSSVMG	ASNN	PLS	KNN
ALog	gPS, OEstate	0.74	0.68	0.61	0.64
	CDDD	0.8	0.74	0.75	0.71
CDK2 (cons,topol,g	eom,elec,hybrid) 3D:corina	0.75	0.71	0.56	0.71
ChemaxonDescrip	tors (pH 0 - 14:1) 3D:corina	0.76	0.7	0.59	0.68
Dragon6	(2D blocks: 1 28)	0.64	0.66	0.59	0.65
Dragon6 (3D b	locks: 1-29) 3D:corina	0.76	0.72	0.57	0.65
Fragmer	ntor (length:2 - 4)	0.72	0.7	0.59	0.63
GS	Frag (F + L)	0.69	0.69	0.61	0.61
InductiveDe	escriptors 3D:corina	0.69	0.71	0.57	0.67
	JPlogP	0.73	0.74	0.59	0.67
	MAP4	0.71	0.65	0.59	0.67
MORDRE	D (All) 3D:corina	0.77	0.73	0.57	0.68
Mera, M	lersy 3D:corina	0.73	0.69	0.55	0.67
	OEstate	0.74	0.67	0.63	0.68
PyDesc	riptor 3D:corina	0.71	0.71	0.7	0.67
QNPF	R (length:1 - 3)	0.68	0.62	0.58	0.58
RDKIT (3D block	(s: 1-11 15-16) 3D:corina	0.77	0.72	0.56	0.65
SIRMS (labels:cha	arge+logp+hb+refractivity)	0.76	0.73	0.59	0.67
Spectrophores	accuracy=20) 3D:corina	0.68	0.6	0.52	0.6
Stru	cturalAlerts	0.67	0.64	0.58	0.51
alvaDesc (3D bloc	:ks: (only) 1-30) 3D:corina	0.75	0.71	0.57	0.68

* Sparse format, DOI:10.1186/s13321-016-0113-y

Consensus modelling

- Best method(s) are defined
- Average prediction of models is used
- The consensus prediction is more accurate and stable

	ironmental Protection Agency TECHNOLOGY I LAWS & REGULATIONS I A	BOUT EPA	ALL EPA THIS AREA Advanced Search
Computational Toxico You are here: EPA Home » Res	logy Research search & Development » CompTox » Chemi	cal Data Challenges & Release	⊠ Contact Us
CompTox Home Basic Information Organization EPA Exposure Research	Research Projects Chemical Databases ToxCast Stakeholder Events EPA Chemical Safety Research	Research Publications Scientific Reviews Communities of Practice ToxCast Data Challenges	Staff Profiles CompTox Partners Jobs and Opportunities

EPA's high-throughput screening data on 1,800 chemicals is accessible through the interactive Chemical Safety for Sustainability Dashboards (iCSS dashboard). The iCSS dashboard provides user-friendly and customizable access to toxicity data from ToxCast and Tox21 high-throughput chemical screening technologies.

Using the **TopCoder** and **InnoCentive** crowd-sourcing platform, EPA invited the science and technology community to work with the data and provide solutions for how the new toxicity data can be used to predict potential health effects. The ToxCast data challenges focused on using this data and other publicly available data to predict the lowest effect level from traditional toxicity studies using laboratory animals. Challenge winners received awards for solving this challenge.

Key Links

- Lowest Effect Level Challenge Results (PDF, 497KB, 18pp)
- · Chemical Safety for Sustainability Dashboards
- Complete ToxCast Phase II Data & Files
- TopCoder Challenge
- InnoCentive Challenge
- Stakeholder Workshops

Novotarskyi, S. et al. Chem. Res. Toxicol. 2016, 29, 768-75.

Model to predict Lowest Effect Level (training set)

Predicted property: LEL Training set: LEL (training) Duplicate models: start

Metrics RMSE - Root Mean Square Error v for Training set v Validation: All v

	ASNN (CHEMAXON)
Adriana (3D by Adriana) 3D:corina	0.93
CDK (cons,topol,geom,elect,hybr) 3D:corina	0.93
ChemaxonDescriptors (pH 0 - 14:1) 3D:corina	0.93
Dragon6 (3D blocks: (only) 1-29) 3D:corina	0.93
Fragmentor (length: 2-4)	0.98
GSFrag (F + L)	0.97
InductiveDescriptors 3D:corina	0.94
Mera, Mersy 3D:corina	0.93
OEstate	0.96
QNPR (length:1 - 3)	0.95
Conser	sus:AVERAGE (CHEMAX
Final EPA ToxCast challenge model	0.88

About

Tox21 Data Challenge 2014

Contact Us

» Home

Registration

Data/Resources

Submissions

Discussion

Leaderboard

Survey

About the Data 🔹

The Challenge

The 2014 Tox21 data challenge is designed to help scientists understand the potential of the chemicals and compounds being tested through the Toxicology in the 21st Century initiative to disrupt biological pathways in ways that may result in toxic effects.

The goal of the challenge is to "crowdsource"

All challenge winners will receive the opportunity to submit a paper for publication in a special thematic issue of

Frontiers in Environmental Science and recognition on the NCATS website and via social media.

Best Balanced accuracy - Abdelaziz, A. et al. *Front. Environ. Sci.* 2016, 4, 2.

Multi-task learning

Multi-task learning

Problem:

- prediction of tissue-air partition coefficients
- small datasets 30-100 molecules (human & rat data)

Results:

simultaneous prediction of several properties increased the accuracy of models

Analysis of toxicity of chemical compounds

*RTECS: Registry of Toxic Effects of Chemical Substances

Sosnin, S. et al. J. Chem. Inf. Model., 2018, 59:1062-1072.

RMSE for different toxicities using CDK descriptors and DNN

Sosnin, S.; Karlov, D.; Tetko, I.V.; Fedorov, M.V. A comparative study of prediction of multitarget toxicity for a broad chemical space. *J Chem Inf Model*. **2018**, **59**, 1062-1072. 40

Machine Learning directly from chemical structures

Saccharin: c1ccc2c(c1)C(=O)NS2(=O)=O

Text processing: convolutional neural networks, transformers, LSTM Graph processing: message passing neural networks

Image augmentation

https://github.com/aleju/imgaug

Machine Learning to canonise chemical structures

SMILES canonization can be done by machine learning!

Machine Learning directly from chemical structures

P. Karpov, G. Godin, I. V. Tetko, J. Cheminform. 2020, 12, 17.

https://github.com/bigchem/transformer-cnn

Convolutional vs. Descriptor-based Neural Neural Networks

Coefficient of determination, r². Transformer CNN provides similar or better accuracy compared to traditional methods based on descriptors <u>even for small datasets (few hundrends compounds!)</u>. P. Karpov, G. Godin, I. V. Tetko, *J. Cheminform.* **2020**, *12*, 17.

Winning model: OCHEM-generated consensus model

Andrea Kopp SLAS Europe 2023

25.05.2023

HELMHOLTZ MUNICH Team of Igor Tetko with Peter Hartog, Martin Šícho and Guillaume Godin

LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

Kopp at al, DOI: <u>10.26434/chemrxiv-2023-p8qcv</u>

Challenge set-up

- Experimentally: Nephelometer measures undissolved sediment
- Classification into *low, medium* and *high* soluble with phenytoin and amiodarone as thresholds
- 70k training datapoints, 15k public leaderboard, 15k private leaderboard
- Stratified random sampling

Workflow with OCHEM

Molecular representation

@Peter Hartog with BioRender.com

Kopp at al, DOI: <u>10.26434/chemrxiv-2023-p8qcv</u>

Quadratic kappa metric scores

@Peter Hartog with BioRender.com

Kopp at al, DOI: <u>10.26434/chemrxiv-2023-p8qcv</u>

openOCHEM

Elias James Corey Nobel Prize 1990

TM Transformations:

FGA – function group additionFGI - function group interconversionFGR – function group removing

Cyclisation, Fisher indole, Mannich, Michael, Oxidation, Rearrangement, and many others.

Retrosynthesis – it is a problem solving technique for transforming the structure of **synthetic target molecule (TM)** to a sequence of progressively simpler structures along the pathway which ultimately leads to simple or commercially available starting materials for a chemical synthesis.

analysis)

These rules make core of the first programs **OCSS** (organic chemical simulation of synthesis) and **LHASA** (logic and heuristics applied to synthetic

Natural Language Processing (NLP) for reaction predictions

Reaction prediction task can be treated as a neural machine translation

The Transformer architecture is currently the state of art for organic reactions

c1cccc1C(=O)O.CC(O)C

https://arxiv.org/abs/1706.03762 (Attention is all you need)

Synthesis planning (top-1)

Direct synthesis USPTO-500k

Schwaller P. (Seq2Seq) 2018	80,3
Jin. W (Weisfeiler-Lehman network) 2017	79 <i>,</i> 6
Kien D. (Policy network) 2019	82,4
Coley C. (Graph convolution) 2018	85,6
Schwaller P. (Transformer) 2018	90,4

Retro-synthesis USPTO 50k

Rule-Base	34,8
Pande V. et al 2016	37,4
ICANN 2019 (our work)	42,7

Karpov, P.; Godin, G.; Tetko, I. V. In *A Transformer Model for Retrosynthesis*, Artificial Neural Networks and Machine Learning – ICANN 2019: Workshop and Special Sessions, Münich, 17th - 19th September 2019; Tetko, I. V.; Kůrková, V.; Karpov, P.; Theis, F., Eds. Springer International Publishing: Münich, 2019; pp 817-830.

Image augmentation

https://github.com/aleju/imgaug

Synthesis planning (top-1)

Retro-synthesis (50k), USPTO-50

Rule-Base	34,8
Seq2Seq, Pande V. 2016	37,4
Graph Logic Networks, Coley C. 2019	52 <i>,</i> 5
Our work: ICANN2019	42,7
Our work: (Augmented Transformer, 2020)*	53,5

Direct synthesis (500k), USPTO-MIT

Schwaller P. (Seq2Seq) 2018 Jin. W (Weisfeiler-Lehman network) 2017	80,3 79,6
Coley C. (Graph convolution) 2018	85,6
Schwaller P. (Transformer) 2019	90,4
Our work (Augmented Transformer, 2020)*	91,9

*Tetko et al Nature Comm. 2020. 11, 11, 1-11.

Retrosynthesis planning

- Design a synthetic route for a target molecule
- Working backward to identify a sequence of reactions to reach simpler starting materials

P. Torren-Peraire et al Models Matter: The Impact of Single-Step Retrosynthesis on Synthesis Planning, <u>https://arxiv.org/abs/2308.05522</u>

Conclusions

- Computational tools are essential for drug discovery
- Use of diverse methods and descriptors can help to identify the best models for the given task
- Representation learning methods provide similar or better performance than traditional models based on descriptors
- Multitask learning can further improve model performances
- Combination of models as consensus usually contributes best performing approaches
- Publishing model on-line allows their wide promotion and acceptance by the scientific community
- Multistep reaction predictions is a new challenge to be addressed with AI

AiChemist MSC DN

https://aichemist.eu/

Home 🗸

MSC ITN Project AiChemist

Optimising biological activity and ADME properties, while minimising toxicity, are objectives when developing new compounds. Advanced machine learning methods are indispensable to this process. The project will develop and benchmark representation learning approaches, addressing their accuracy and explainability, using public and *in-house* data for endpoints ranging from chemical reactions to toxicity. The program will be done with the target users: large companies, regulatory agencies and SMEs.

6 out of 14 positions are still available

Also see Twitter: https://twitter.com/aichemist_dn

https://icann2024.org

Welcome to ICANN 2024

The 33rd International Conference on Artificial Neural Networks. A conference of the European Neural Network Society.

The 2024 edition of the ICANN will be organized by the Dalle Molle Institute for Artificial Intelligence Research (IDSIA USI-SUPSI) in Lugano, Switzerland in collaboration with AIDD and AiChemist Horizon MSCA projects.

Conference venue: USI-SUPSI Campus Est, Via la Santa 1, 6962 Lugano-Viganello, Switzerland

Conference dates: September 17 to September 20, 2024.

Special session: AI in Drug Discovery Big Data and advanced machine learning in chemistry eXplainable AI (XAI) in chemistry Chemoinformatics Use of deep learning to predict molecular properties Modeling and prediction of chemical reaction data Generative models

SUPSI

Acknowledgements

Andi Kopp Peter Hartog Fabian Krüger Paula Torren-Peraire Varvara Voinarovska Katya Ahmad Nesma Mousa Mark Embrechts

Guillaume Godin (Firmenich) Ruud van Deursen (Firmenich)

Michael Sattler (HMGU)

Alexander von Humboldt Stiftung/Foundation

AiChemist

