Introduction to Machine Learning for Chemists:
Visualization, Data Processing, Analysis,
Molecular Design

# Bibliography: review articles
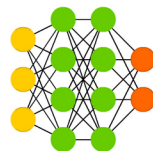
LPCNO

A mostly complete chart of

# Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen    asimovinstitute.org

## Review articles



B. Sanchez-Lengeling & A. Aspuru-Guzik (**2018**)
Inverse molecular design using machine learning: Generative models for matter engineering
*Science* **361**, 360-365

K. T. Butler, D. W. Davies, H. Cartwright, O. Isayev & A. Walsh (**2018**)
Machine learning for molecular and materials science
*Nature* **559**, 547-555

P. Schlexer Lamoureux, K. T. Winther, J. A. Garrido Torres, V. Streibel, M. Zhao, M. Bajdich, F. Abild-Pedersen & T. Bligaard (**2019**)
Machine Learning for Computational Heterogeneous Catalysis
*ChemCatChem* **11**, 3581-3601

J. A. Keith, V. Vassilev-Galindo, B. Cheng, S. Chmiela, M. Gastegger, K.-R. Müller & A. Tkatchenko (**2021**)
Combining Machine Learning and Computational Chemistry for Predictive Insights Into Chemical Systems
*Chem. Rev.* **121**, 9816–9872

## Review articles



B. Huang & O. A. von Lilienfeld (**2021**)
Ab Initio Machine Learning in Chemical Compound Space
*Chem. Rev*. **121**, 10001-10036



Nandy, A.; Duan, C.; Taylor, M. G.; Liu, F.; Steeves, A. H.; Kulik, H. J. (**2021**)
Computational Discovery of Transition-Metal Complexes: From High-Throughput Screening to Machine Learning.
*Chem. Rev*. **121**, 9927–10000

## Review articles

### ACS Catalysis

Review

Cite This: *ACS Catal.* 2020, 10, 2260−2297

pubs.acs.org/acscatalysis

**Machine Learning for Catalysis Informatics: Recent Applications and Prospects**

Takashi Toyao,[†,‡] Zen Maeno,[†] Satoru Takakusagi,[†] Takashi Kamachi,[‡,§] Ichigaku Takigawa,[*,‖,⊥] and Ken-ichi Shimizu[*,†,‡]

[†]Institute for Catalysis, Hokkaido University, N-21, W-10, Sapporo 001-0021, Japan
[‡]Elements Strategy Initiative for Catalysts and Batteries, Kyoto University, Katsura, Kyoto 615-8520, Japan
[§]Department of Life, Environment and Materials Science, Fukuoka Institute of Technology, 3-30-1Wajiro-Higashi, Higashi-ku, Fukuoka 811-0295, Japan
[‖]RIKEN Center for Advanced Intelligence Project, 1-4-1 Nihonbashi, Chuo-ku, Tokyo 103-0027, Japan
[⊥]Institute for Chemical Reaction Design and Discovery (WPI-ICReDD), Hokkaido University, Kita 21 Nishi 10, Kita-ku, Sapporo, Hokkaido 001-0021, Japan

**ABSTRACT:** The discovery and development of catalysts and catalytic processes are essential components to maintaining an ecological balance in the future. Recent revolutions made in data science could have a great impact on traditional catalysis research in both industry and academia and could accelerate the development of catalysts. Machine learning (ML), a subfield of data science, can play a central role in this paradigm shift away from the use of traditional approaches. In this review, we present a user's guide for ML that we believe will be helpful for scientists performing research in the field of catalysis and summarize recent progress that has been made in utilizing ML to create homogeneous and heterogeneous catalysts. The focus of the review is on the design, synthesis, and characterization of catalytic materials/compounds as well as their applications to catalyzed processes. The ML technique not only enhances ways to discover catalysts but also serves as a powerful tool to establish a deeper understanding of relationships between the properties of materials/compounds and their catalytic activities, selectivities, and stabilities. This knowledge facilitates the establishment of principles employed to design catalysts and to enhance their efficiencies. Despite such advantages of ML, it is noteworthy that the current ML-assisted development of real catalysts remains in its infancy, mainly because of the complexity of catalysis associated with the fact that catalysis is a time-dependent dynamic event. In this review, we discuss how seamless integration of experiment, theory, and data science can be used to accelerate catalyst development and to guide future studies aimed at applications that will impact society's need to produce energy, materials, and chemicals. Moreover, the limitations and difficulties of ML in catalysis research originating from the complex nature of catalysis are discussed in order to make the catalysis community aware of challenges that need to be addressed for effective and practical use of ML in the field.

**KEYWORDS:** *machine learning, catalysis informatics, high-throughput experiments/computations, data mining, structure−activity relationships*

Toyao, T.; Maeno, Z.; Takakusagi, S.; Kamachi, T.; Takigawa, I.; Shimizu, K.-i. (**2020**)
Machine Learning for Catalysis Informatics: Recent Applications and Prospects
*ACS Catal*. **10** (3), 2260–2297.

### MACHINE LEARNING

### REVIEWS

Check for updates

### Nanoparticle synthesis assisted by machine learning

*Huachen Tao[1], Tianyi Wu[1], Matteo Aldeghi[1,2,3], Tony C. Wu[1,3], Alán Aspuru-Guzik[1,2,3,4✉] and Eugenia Kumacheva[1,5,6✉]*

Abstract | Many properties of nanoparticles are governed by their shape, size, polydispersity and surface chemistry. To apply nanoparticles in chemical sensing, medical diagnostics, catalysis, thermoelectrics, photovoltaics or pharmaceutics, they have to be synthesized with precisely controlled characteristics. This is a time-consuming, laborious and resource-intensive task, because nanoparticle syntheses often include multiple reagents and are conducted under interdependent experimental conditions. Machine learning (ML) offers a promising tool for the accelerated development of efficient protocols for nanoparticle synthesis and, potentially, for the synthesis of new types of nanoparticles. In this Review, we discuss ML algorithms that can be used for nanoparticle synthesis and highlight key approaches for the collection of large datasets. We examine ML-guided synthesis of semiconductor, metal, carbon-based and polymeric nanoparticles, and conclude with a discussion of current limitations, advantages and perspectives in the development of ML-assisted nanoparticle synthesis.

H. Tao, T. Wu, M. Aldeghi, T. C. Wu, A. Aspuru-Guzik & E. Kumacheva (**2021**)
Nanoparticle synthesis assisted by machine learning
*Nat. Rev. Mater*. **6**: 701-716

**AIChem ← is also currently used to reduce the cost of computational [quantum] chemistry**

## Calculation of energies & forces (MD / geometry optimization)



**Figure 5.** (a) "Magic cube"[59] depiction of hierarchies of correlated wavefunction approaches. (b) "Jacob's Ladder"[60] depiction of hierarchies of Kohn−Sham density functional theory (DFT) approaches. (c) Hierarchies of atomistic potentials. (d) Overall hierarchies in predictive atomic scale modeling methods.

J. A. Keith, V. Vassilev-Galindo, B. Cheng, S. Chmiela, M. Gastegger, K.-R. Müller & A. Tkatchenko (**2021**)
Combining Machine Learning and Computational Chemistry for Predictive Insights Into Chemical Systems. *Chem. Rev.* **121**, 9816–9872

# Short selection of [simple] applications of supervised learning to chemistry

## ML is only as good as the data it learned from

**quality**

**diversity**

**volume**

### when ML is applied for prediction purpose (ML-QSPR)

## Exploring chemical compound space with quantum-based machine learning

*O. Anatole von Lilienfeld, Klaus-Robert Müller and Alexandre Tkatchenko*

Abstract | Rational design of compounds with specific properties requires understanding and fast evaluation of molecular properties throughout chemical compound space — the huge set of all potentially stable molecules. Recent advances in combining quantum-mechanical calculations with machine learning provide powerful tools for exploring wide swathes of chemical compound space. We present our perspective on this exciting and quickly developing field by discussing key advances in the development and applications of quantum-mechanics-based machine-learning methods to diverse compounds and properties, and outlining the challenges ahead. We argue that significant progress in the exploration and understanding of chemical compound space can be made through a systematic combination of rigorous physical theories, comprehensive synthetic data sets of microscopic and macroscopic properties, and modern machine-learning methods that account for physical and chemical knowledge.

### but it is also a remarkable tool to find correlations between data

## Control of an organic synthesis robot

**space of chemical reactions explored quickly**

**organic synthesis robot performs chemical reactions and analysis faster than a human**

**prediction of the reactivity of possible reagent combinations after conducting a small number of experiments**

**decision making by ML**

robot equipped with real-time sensors to record the spectra of the reaction mixtures:

- flow benchtop NMR system
- mass spectrometer
- IR system

reactions mixtures classified as reactive (**1**) or non reactive (**0**) by the robot (uses a **support vector machine (SVM) algorithm** with a **linear kernel model**)

**Extended Data Fig. 1 | Reaction space explored.** The chemical inputs (**1–18**) used in the platform to search for new transformations and to evaluate the performance of the algorithm.

J. M. Granda, L. Donina, V. Dragone, D.-L. Long & L. Cronin (**2018**). Controlling an organic synthesis robot with machine learning to search for new reactivity *Nature* **559**, 377-381

**flattening of the data →**

**Fig. 4 | Exploring the Suzuki–Miyaura reaction using machine learning. a**, The reaction space of the Suzuki–Miyaura reaction. Shown are the identity of reactants, ligand, base and solvent, and the vector representation of the reaction for machine learning. **b**, Validation of the predictive power of the model for a test set of 30% of the reactions (1,728 reactions). RMSE, root-mean-square error. **c**, Simulation of the machine-learning-controlled exploration of this reaction space. The yellow bar shows the initial random choice of 10% of reaction space (576 reactions). The green bars show the next batches of 100 reactions chosen by the machine learning algorithm. The error bars represent the standard deviation within individual batches for Suzuki–Miyaura coupling.

J. M. Granda, L. Donina, V. Dragone, D.-L. Long & L. Cronin (**2018**). Controlling an organic synthesis robot with machine learning to search for new reactivity
*Nature* **559**, 377-381

10

## Prediction of spectroscopic parameters



F. M. Paruzzo, A. Hofstetter, F. Musil, S. De, M. Ceriotti & L. Emsley (**2018**)
Chemical shifts in molecular solids by machine learning
*Nat. Commun.* **9**, 4501

**Test**



**predicted values as a function of the actual values**

Deep Convolutional Network (DCN)

## Image recognition



**Feature extraction**

**Classification**

**Fig. 3.** Schematic detailing the architecture of the CNNs. (A colour version of this figure can be viewed online.)

### Molecular Dynamics + STEM image simulation



(a) (13,0)

(b) (13,2)

(c) (13,4)

Potential energy (eV/atom)

**Fig. 1.** Some examples of final snapshots from the molecular dynamics simulations. Defects (higher potential energy, appear in yellow) have been introduced by high-temperature annealing. The visualizations are generated using the OVITO software [51]. (A colour version of this figure can be viewed online.)

Dr. Probe - Software

DOWNLOAD    DOCUMENTATION    EXAMPLES    DEVELOP

J. Barthel, *Dr. Probe: A software for high-resolution STEM image sim*

NVIDIA. CUDA.

Hybrid parallel CPU & GPU computing

[BIG] DATA

5000 images per chirality

~1.3 10⁶ images in total



(a) True chirality: (8, 0) — classified as: (8, 0), 100.0%
(b) True chirality: (11, 0) — classified as: (11, 0), 100.0%
(c) True chirality: (13, 1) — classified as: (13, 1), 100.0%
(d) True chirality: (24, 5) — classified as: (25, 5), 99.4%
(e) True chirality: (24, 5) — classified as: (25, 5), 98.8%
(f) True chirality: (14, 2) — classified as: (13, 2), 98.5%
(g) True chirality: (25, 4) — classified as: (6, 1), 72.6%
(h) True chirality: (28, 1) — classified as: (9, 0), 55.5%
(i) True chirality: (22, 4) — classified as: (4, 3), 92.6%

G. D. Förster, A. Castan, A. Loiseau, J. Nelayah, D. Alloyeau, F. Fossard, C. Bichara & H. Amara (**2020**)
A deep learning approach for determining the chiral indices of carbon nanotubes from high-resolution transmission electron microscopy images
*Carbon* **169**, 465-474

Version: Sunday, October 1, 2023

**Reaction outcome prediction: Prediction of enantiomeric excess**

## Chemical Science

ROYAL SOCIETY OF CHEMISTRY

**EDGE ARTICLE**

View Article Online
View Journal | View Issue

Check for updates

# Reaction-based machine learning representations for predicting the enantioselectivity of organocatalysts†

Simone Gallarati, ‡[a] Raimon Fabregat, ‡[a] Rubén Laplaza, [ab]
Sinjini Bhattacharjee, [ac] Matthew D. Wodrich [ab]
and Clemence Corminboeuf [*abd]

Hundreds of catalytic methods are developed each year to meet the demand for high-purity chiral compounds. The computational design of enantioselective organocatalysts remains a significant challenge, as catalysts are typically discovered through experimental screening. Recent advances in combining quantum chemical computations and machine learning (ML) hold great potential to propel the next leap forward in asymmetric catalysis. Within the context of quantum chemical machine learning (QML, or atomistic ML), the ML representations used to encode the three-dimensional structure of molecules and evaluate their similarity cannot easily capture the subtle energy differences that govern enantioselectivity. Here, we present a general strategy for improving molecular representations within an atomistic machine learning model to predict the DFT-computed enantiomeric excess of asymmetric propargylation organocatalysts solely from the structure of catalytic cycle intermediates. Mean absolute errors as low as 0.25 kcal mol$^{-1}$ were achieved in predictions of the activation energy with respect to DFT computations. By virtue of its design, this strategy is generalisable to other ML models, to experimental data and to any catalytic asymmetric reaction, enabling the rapid screening of structurally diverse organocatalysts from available structural information.

$$\text{e.e.} = \frac{|A - B|}{A + B} \times 100$$

**Computational chemistry**

$$\Delta\Delta E^{\ddagger} \updownarrow$$



$$\text{e.e.} = \frac{1 - e^{\Delta\Delta E^{\ddagger}/RT}}{1 + e^{\Delta\Delta E^{\ddagger}/RT}} \times 100$$

**With $\Delta\Delta E^{\ddagger} \lesssim 5$ kcal.mol$^{-1}$**

S. Gallarati, R. Fabregat, R. Laplaza, S. Bhattacharjee, M. D. Wodrich & C. Corminboeuf (2021)
Reaction-based machine learning representations for predicting the enantioselectivity of organocatalysts, *Chem. Sci.* **12**, 6879-6889

13

## Chemical context

**Target reaction**



Bipyridine N,N'-dioxide organocatalyst

**Enantiodetermining TS**

Scheme 2 Catalytic cycle for the propargylation of benzaldehyde with allenyltrichlorosilane, showing the rate-limiting and stereocontrolling transition state. Adapted from ref. 85.

**KRR** (not an ANN)



**(1) Training**

1 Database of intermediates — Int 2, Int 3 — 754 pairs of DFT-optimized intermediates

2 Generate reaction representation — a) Int 2; b) Int 3 - Int 2 — Molecular representation (e.g., SLATM)

3 Map representation to target property with KRR — TS, $E_a$, Reaction progress — Target: DFT-computed $E_a$

4 Hyperparameter optimization and cross-validation — Train, Test — 100 cross-validation train/test splits (678/76)

**76 Lewis base organocatalysts**



1 (Y = H)
2 (Y = Ph)
3 (Y = $^t$Bu)
4
5
6

a: X = H
b: X = F
c: X = Cl
d: X = CH$_3$
e: X = CF$_3$
f: X = $^i$Pr
g: X = $^t$Bu
h: X = CCH
i: X = CN
j: X = Ph

[BIG] DATA

Scheme 1 Library of axially chiral bipyridine N,N'-dioxide organo-catalysts. R = H or Me. Adapted from ref. 74.

**5 possible ligand arrangements around the hexacoordinate silicon in the TS**
**+**
**the alkyl nucleophile can add to either face of benzaldehyde**

**760 DFT TS geometries and activation energies**

14

## Main results

**Mean Absolute Error**

$$MAE = \frac{1}{n} \sum_{i=1}^{n} \left| E_a^{(ref.)} - \hat{E}_a \right|$$

**SLATM = structural descriptors**



**Learning curves:**
**MAE in test sets predictions of $E_a$**

**0.54 ± 0.06 kcal mol$^{-1}$ (insufficient accuracy)**

**0.31 ± 0.20 kcal mol$^{-1}$**

**0.25 ± 0.40 kcal mol$^{-1}$**



**Fig. 2** Predictions of $\Delta\Delta E^{\ddagger}$ vs. DFT reference for the three approaches discussed. Mean Absolute Errors (MAE) are reported in kcal mol$^{-1}$. These predictions are obtained by averaging the predictions obtained from the cross-validation scheme with 100 different random train/test splits. The error bars indicate the standard deviation of ML $\Delta\Delta E^{\ddagger}$, derived from the standard deviations in the $E_a$ prediction of the 100 different random train/test splits.

## Prediction of CO$_2$ solubility in ionic liquids

Contents lists available at ScienceDirect

**Chemical Engineering Science**

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

**CHEMICAL ENGINEERING SCIENCE**

ELSEVIER

Prediction of CO$_2$ solubility in ionic liquids using machine learning methods

Zhen Song [a], Huaiwei Shi [a,b], Xiang Zhang [c], Teng Zhou [a,b,*]

[a] Process Systems Engineering, Otto-von-Guericke University Magdeburg, Universitätsplatz 2, D-39106 Magdeburg, Germany
[b] Process Systems Engineering, Max Planck Institute for Dynamics of Complex Technical Systems, Sandtorstr. 1, D-39106 Magdeburg, Germany
[c] Department of Chemical and Biological Engineering, The Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong, China

### Database = 10 116 experimental data points

[BIG] DATA

**124 different ILs**

$T$: [243.2 K – 453.15 K]

$p$: [0.00798 bar – 499.9 bar]

**training set, 80% of the data = 8093 values**

**test set, 20% of the data = 2023 values**

*cations*: imidazolium, pyrrolidinium, pyridinium, piperidinium, ammonium, phosphonium, sulfonium

*anions*: tetrafluoroborate [BF4], chloride [Cl], dicyanamide [DCA], nitrate [NO3], hexafluorophosphate [PF6], thiocyanate [SCN], tri-cyanomethanide [C(CN)3], hydrogen sulfate [HSO4], bis(trifluoromethylsulfonyl)amide [Tf2N], methylsulfate [MeSO4], etc

### possible modelling approaches

**Thermodynamic models**

**DFT studies in solvent**

} **not accurate enough**

### other possibility
**quantitative structure-property relationship (QSPR)**

**Machine Learning
(surrogate modeling)**

Version: Sunday, October 1, 2023

LPCNO

Z. Song, H. Shi, X. Zhang & T. Zhou (**2020**)
Prediction of CO2 solubility in ionic liquids using machine learning methods, *Chem. Eng. Sci.* **223**, 115752

16

## three-layer feed forward ANN



Deep Feed Forward (DFF)

Input p (53×1) — IL structure information

$W_1$ (7×53), $b_1$ (7×1), $a_1$ (7×1), $a_1 = W_1 \times p + b_1$

$f_1(a_1)$

$W_2$ (1×7), $b_2$ (1×1), $a_2$ (1×1), $a_2 = W_2 \times f_1(a_1) + b_2$

Output = $f_2(a_2)$ — $CO_2$-in-IL solubility

**Fig. 1.** Schematic structure of the employed three-layer ANN (the dimensions of $W_1$, $W_2$, $b_1$, and $b_2$ are given in the brackets).

| T (K) | P (bar) | [BETA] | [DMPO4] | [HSO4] | [DBPO4] | [methide] | [C3F7CO2] | [NH] | [TOS] | [MPip] | [S] | CH=CH2 | CH=CH | [CH2] | [CH] | [OCH2] | [OCH3] | [CF2] | [OH] | [MIm] | [MPyrro] | [N] | [P] | [BF4] | [Cl] | [DCA] | x_CO2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 313.15 | 50 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0.583 |
| 298.25 | 0.4624 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0.014334 |

**data organized as vectors**
**IL structure information = occurrences of functional groups in the IL**



training set: MAE = 0.0200, R² = 0.9842
test set: MAE = 0.0202, R² = 0.9836

**Fig. 2.** Comparison between the experimental and ANN-GC predicted $CO_2$ solubility.

17

**Two notebooks**

Python in the
[Physical] Chemistry Lab

[ PytChem ]

**DS4B-CO2_solubility-ANN.ipynb**
**DS4B-CO2_solubility-SVR.ipynb**

## General context: steam methane reforming

**WGSR: CO + H₂O ⟶ CO₂ + H₂**

steam methane reforming
**+206,2 kJ·mol⁻¹**

**−41,1 kJ·mol⁻¹**



Fig. 1 – Illustration showing the overall process of methane being converted to CO₂ and H₂ on a supported metal catalyst. Many catalysts follow this same general mechanism.

Feed Forward (FF)



*neuralnet* package available in the R software environment

Fig. 2. Three-Layer Feedforward Neural Network employed in this work.

F. M. Cavalcanti, M. Schmal, R. Giudici & R. M. B. Alves (**2019**) A catalyst selection method for hydrogen production through Water-Gas Shift Reaction using artificial neural networks
*J. Environ. Manage.* **237**, 585-594

## Data and main outcome of the trained ANN

**Table 1**
Variables selected for the ANN and their ranges.

| Variable | Unit | Minimum value | Maximum value |
|---|---|---|---|
| Temperature | °C | 200 | 450 |
| Pressure | bar | 0.8 | 27.6 |
| Catalyst mass | g | 0.02 | 2.86 |
| Gas hourly space velocity (GHSV) | $h^{-1}$ | 795 | $(1,200,000)$[b] |
| CO feed composition | vol% | 1.30 | 37.2 |
| $H_2O$ feed composition | vol% | 1.50 | 69.2 |
| $CO_2$ feed composition | vol% | 0 | 96.0 |
| $H_2$ feed composition | vol% | 0 | 62.5 |
| Inert feed composition ($N_2$ or He) | vol% | 0 | 96.50 |
| $CH_4$ feed composition | vol% | 0 | 0.70 |
| Active phase composition[a] | wt% | Co, Ni, Cu, Ru, Pd, Ag, Ir, Pt, Au, Cr, Zn | |
| Support type[a] | – | $Fe_2O_3$, AC, CNT, $Mo_2C$, $CeO_2$, $La_2O_3$, $ZrO_2$, MgO, $Al_2O_3$, $TiO_2$ | |
| Promotor/dopant concentration[a] | wt% | Na, K, Mg, Ba, B, Al, Si, Pb, S, Hg, Y, Ti, Zr, La, Ce, Fe | |
| Surface area | $m^2/g$ | 1.1 | $(1487)$[c] |
| Calcination temperature | °C | 25 | 800 |
| Calcination time | h | 0 | 10 |
| CO conversion | dimensionless | 0 | 1 |

[a] Categorical or categorical-quantitative variables.
[b] This maximum value reported for GHSV of $1.2 \times 10^6$ ml gas/ml catalyst/h (Rhodes et al., 2002) is very unusual in catalytic experiment ranges ($\approx 10^4$ $h^{-1}$), leading to a rather small residence time (0,003 s).
[c] This maximum value reported for the surface area of 1487 $m^2/g$ (Buitrago et al., 2012) is well above the normally found catalyst surface area values (100–300 $m^2/g$), since this catalyst support is a special industrial activated carbon prepared from olive stones by direct steam activation.

**Taken from**: T. L. LeValley, A. R. Richard & M. Fan (**2014**)
The progress in water gas shift and steam reforming hydrogen production technologies – A review
*Int. J. Hydrogen Energy* **39**, 16983-17000

**[BIG] DATA**

**283 experimental data points**


**Training set**


**Test set**

## Best conditions? 1. ANN at work



**Fig. 2.** Three-Layer Feedforward Neural Network employed in this work.

**Pd** and **Co**: inconsistencies in the ANN model
- lack of data for these conditions?
- or inconsistent kinetic data?



**Fig. 8.** CO conversion *versus* temperature for different active phases (NH = 12, Metal/CeO$_2$, Metal = 2 wt%, P = 1 bar, m$_{cat}$ = 0.1 g, GHSV = 1000 h$^{-1}$, surface area = 100 m$^2$/g, T$_{calc}$ = 300 °C, t$_{calc}$ = 4 h, feed composition: 2% CO, 10% H$_2$O, 88% N$_2$).

**Best metals = Ru, Ni, and Cu**

## Best conditions? 2. Sensitivity analysis using ANN predictions

if...
- relationships between inputs and outputs are poorly understood
- poor or partial understanding of the driving forces and mechanisms

⇒ limits our confidence in the output of the model

⇒ sensitivity analysis provides a kind of "quality assurance"

**Table 3**
CO conversion sensitivities related to the considered input variables.

| Type of variable | Variable | Sensitivity |
|---|---|---|
| Catalyst design and texture | Cu composition | 0.0389 |
| | Surface area | **0.424** |
| | Calcination Temperature | 0.0917 |
| | Calcination Time | 0.0197 |
| Operating conditions | Temperature | **1.14** |
| | GHSV | $-0.00740$ |
| | CO feed composition | 0.0365 |
| | $H_2O$ feed composition | 0.0180 |
| | Inert feed composition | $-0.0777$ |

Is such analysis <u>really</u> necessary to demonstrate that surface area and temperature are important to the development of industrial catalysts for the WGS reaction?

→ actually, sensitivity analysis rather orders by importance the strength and relevance of the inputs in determining the variation in the output

Best conditions and evaluation of the most relevant variables? → "standard" statistics is needed

Version: Sunday, October 1, 2023

LPCNO

THE JOURNAL OF
**PHYSICAL CHEMISTRY** A

## Improved Chemical Prediction from Scarce Data Sets via Latent Space Enrichment

*Published as part of The Journal of Physical Chemistry virtual special issue "Young Scientists".*

Nicolae C. Iovanac and Brett M. Savoie*

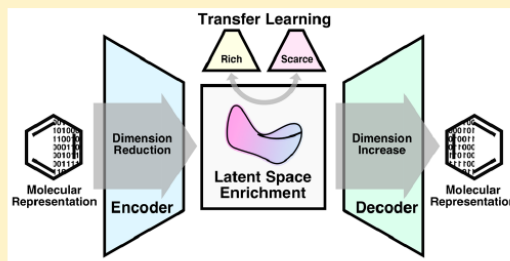Charles D. Davidson School of Chemical Engineering, 480 Stadium Mall Drive, Purdue University, West Lafayette, Indiana 47906, United States
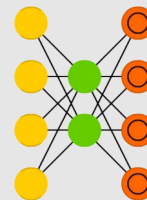
**S** Supporting Information

**ABSTRACT:** Modern machine learning provides promising methods for accelerating the discovery and characterization of novel chemical species. However, in many areas experimental data remain costly and scarce, and computational models are unavailable for targeted figures of merit. Here we report a promising pathway to address this challenge by using chemical latent space enrichment, whereby disparate data sources are combined in joint prediction tasks to enable improved prediction in data-scarce applications. The approach is demonstrated for $pK_a$ prediction of moderately sized molecular species using a combination of experimentally available $pK_a$ data and density functional theory-based characterizations of the (de)protonation free energy. A novel autoencoder framework is used to create a continuous chemical latent space that is then used in single and joint training tasks for property prediction. By combining these two data sets in a jointly trained autoencoder framework, we observe mutual improvement in property prediction tasks in the scarce data limit. We also demonstrate an enrichment mechanism that is unique to latent space training, whereby training on excess computational data can mitigate the prediction losses associated with scarce experimental data and advantageously organize the latent space. These results demonstrate that disparate chemical data sources can be advantageously combined in an autoencoder framework with potential general application to data-scarce chemical learning tasks.

**scarce experimental data are supplemented in learning tasks with more abundant computational properties**

**exp p$K_a$ ⟷ DFT deprotonation energies**

Auto Encoder (AE)

Generative Adversarial Network (GAN)

Auto Encoder (AE)

# Entangled Conditional Adversarial Autoencoder for de Novo Drug Discovery

Daniil Polykovskiy,[*,†,‡] Alexander Zhebrak,[†] Dmitry Vetrov,[‡] Yan Ivanenkov,[†,⊥,∥] Vladimir Aladinskiy,[†,∥] Polina Mamoshina,[†] Marine Bozdaganyan,[†] Alexander Aliper,[†] Alex Zhavoronkov,[†] and Artur Kadurin[†,§]

[†]Insilico Medicine, Rockville, Maryland 20850, United States
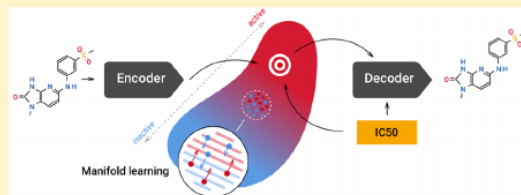[‡]National Research University Higher School of Economics, Moscow 101000, Russia
[§]Insilico Taiwan, Taipei City 115, Taiwan R.O.C
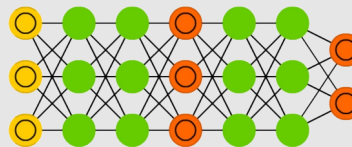[⊥]Institute of Biochemistry and Genetics Russian Academy of Science, Ufa, 450054, Russia
[∥]Moscow Institute of Physics and Technology (State University), Moscow Region, 141700, Russia

**ABSTRACT:** Modern computational approaches and machine learning techniques accelerate the invention of new drugs. Generative models can discover novel molecular structures within hours, while conventional drug discovery pipelines require months of work. In this article, we propose a new generative architecture, entangled conditional adversarial autoencoder, that generates molecular structures based on various properties, such as activity against a specific protein, solubility, or ease of synthesis. We apply the proposed model to generate a novel inhibitor of Janus kinase 3, implicated in rheumatoid arthritis, psoriasis, and vitiligo. The discovered molecule was tested in vitro and showed good activity and selectivity.
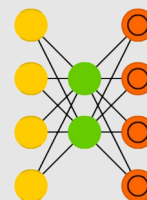
**KEYWORDS:** *adversarial autoencoders, disentanglement, conditional generation, Janus kinase*

# ML in chemistry (MLChem)

ML methods are becoming less understood while they are also more regularly used as black box tools.

Many publications show inadequate technical expertise in ML (e.g. inappropriate splitting of training, testing, and validation sets)

It can be difficult to compare different ML methods and know which is the best for a particular application or whether ML should even be used at all

Data quality and context are often missing from ML modeling, and data sets need to be made freely available and clearly explained

Version: Sunday, October 1, 2023