

Performing research with publicly-available data ORD I Lausanne

Philippe Schwaller (he/him/his)

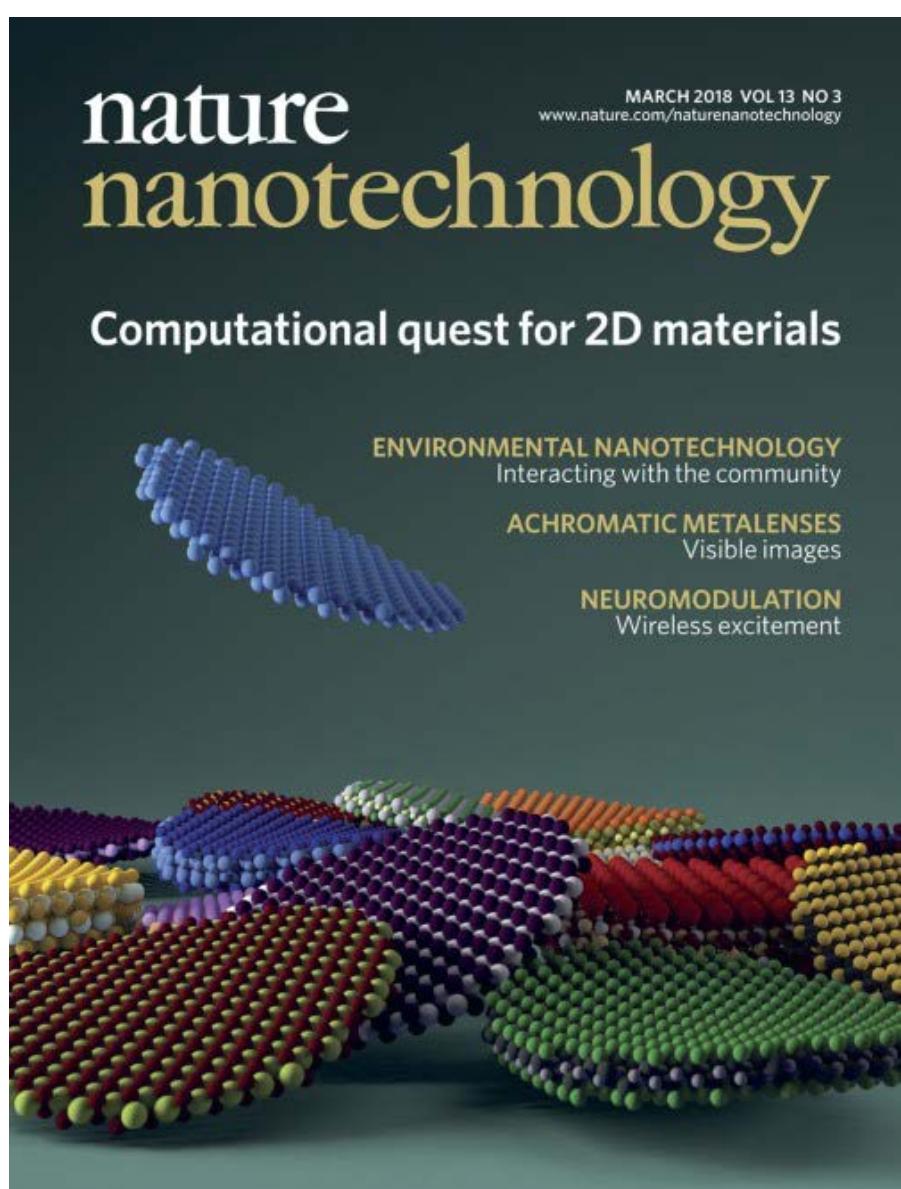
Laboratory of Artificial Chemical
Intelligence (LIAC) – SB-ISIC

@SchwallerGroup

philippe.schwaller@epfl.ch

Grew up in Fribourg,
Switzerland
- French
- Swiss German / German

EPFL



- MaX prize for frontier HPC applications ('17)
- PRACE HPC Excellence Award ('22)



TT Assistant Professor
in Digital Chemistry
since Feb 2022



Materials Science &
Engineering

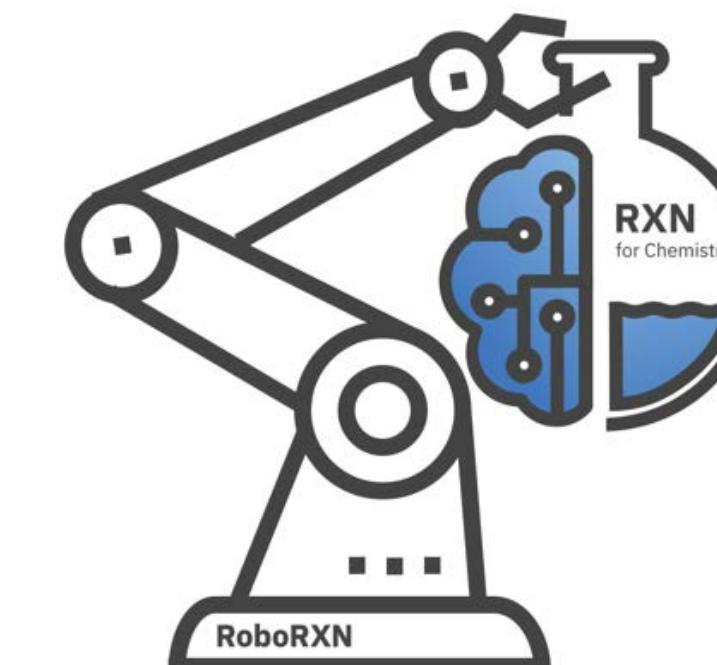
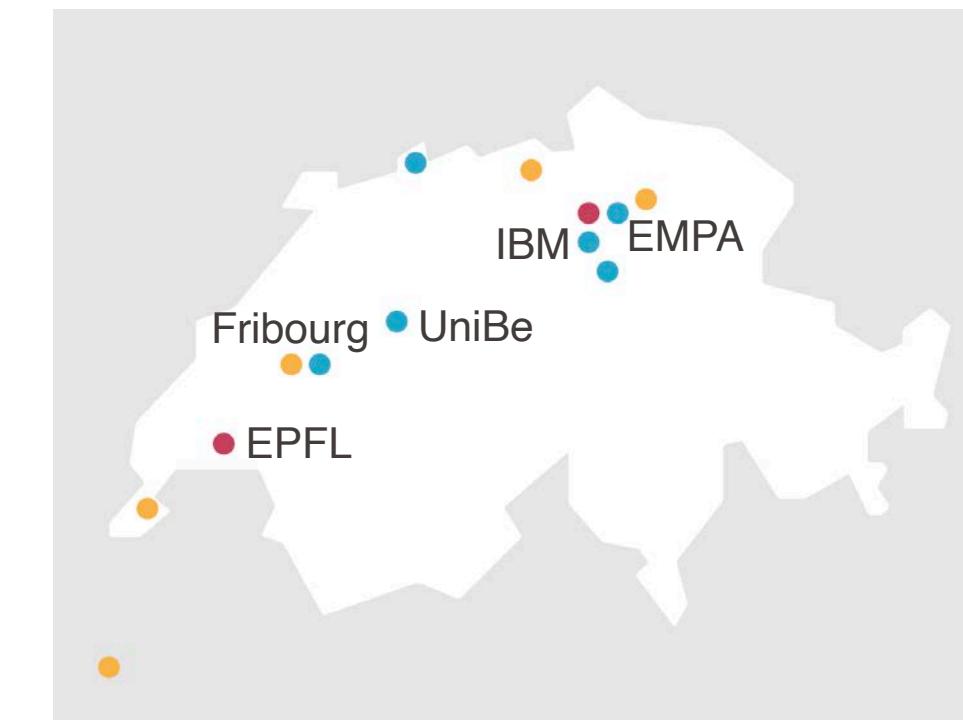
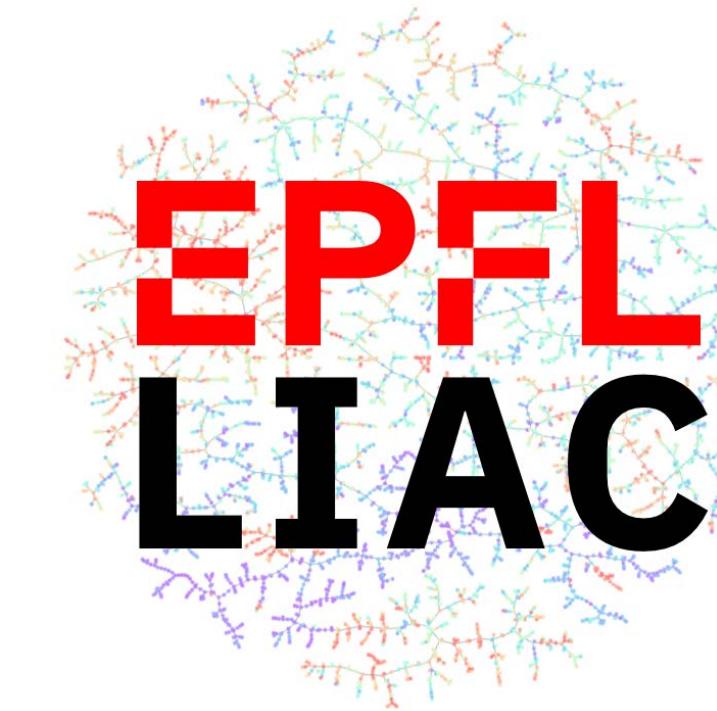
Virtual screening &
simulation workflows
Prof Nicola Marzari

BSc ('14)/MSc ('16)



Materials Science and Technology

Lab work on ternary polymer
blends for organic solar cells
Prof Frank Nüesch

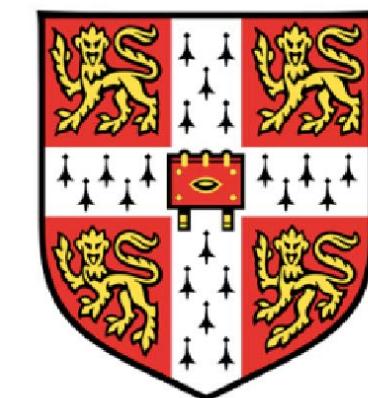


Machine learning for chemical synthesis
Intern/PhD/Postdoc
Dr Teodoro Laino

u^b

^b
UNIVERSITÄT
BERN

PhD in Chemistry
and Molecular Sciences ('21)
Prof Jean-Louis Reymond



UNIVERSITY OF
CAMBRIDGE



IBM Research Europe

PhD students:

Bojana Ranković

Oliver Schilter

Andres CM Bran

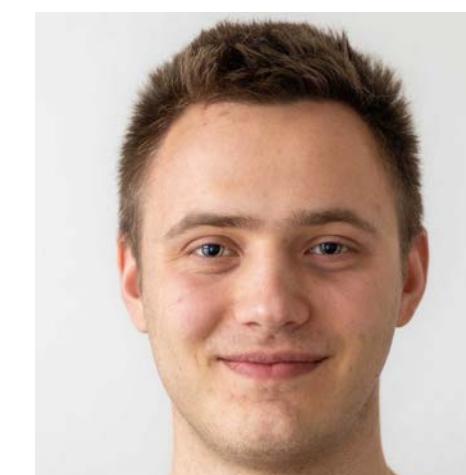
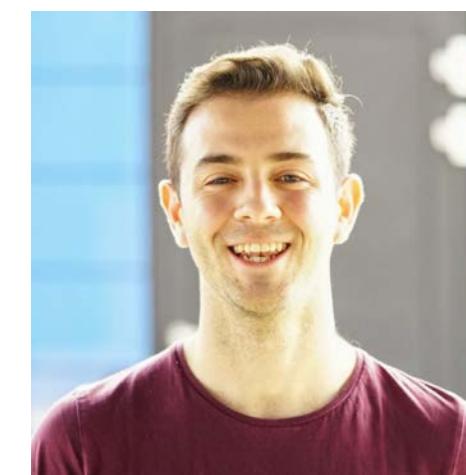
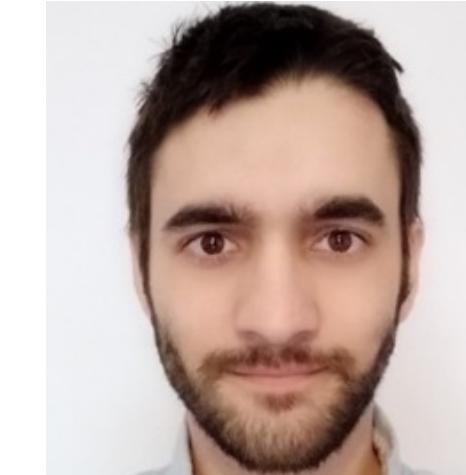
Junwu Chen

Jeff Guo

Victor Sabanza Gil

Paulo Neves

Rebecca Neeser



Stéphane D'Ascoli

Cheng-Hua Huang

Artur Stefaniuk

Malte Franke

Xu Huang

Lucien Brey

Annick Delmonaco

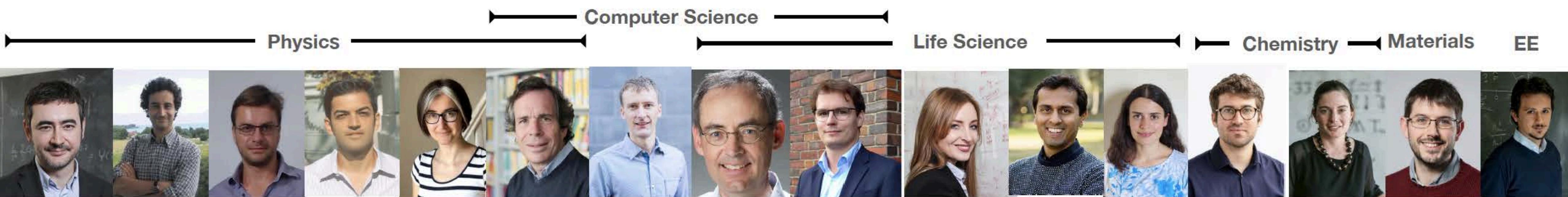
Working group:



Goal:

Identify questions that re-appear when applying ML to accelerate various areas of science. Join forces to solve them.

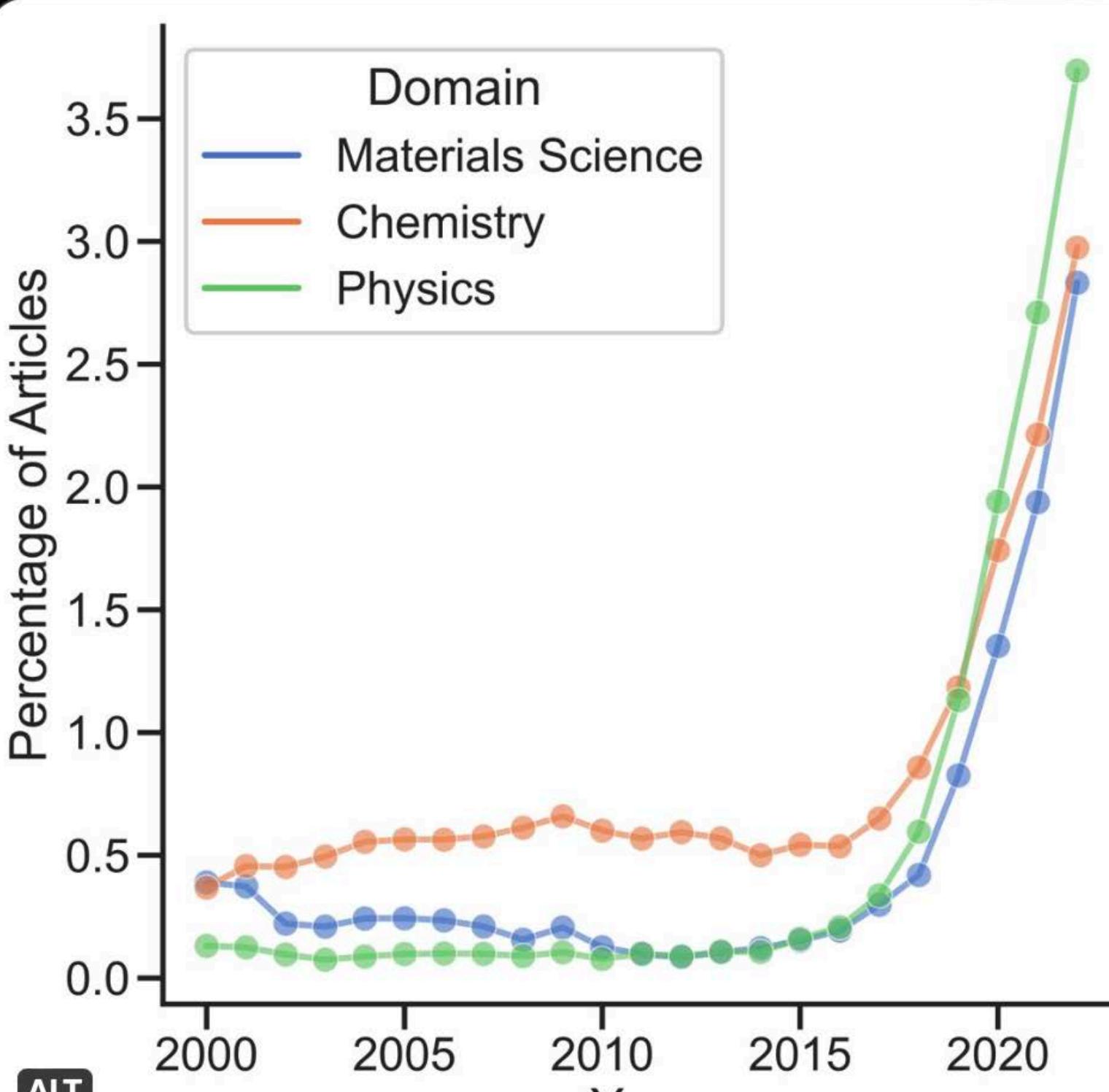
Serendipity is our friend.



 **Ben Blaiszik** @BenBlaiszik · 45m

Does it feel like you are seeing more impactful #ML and #AI for science publications? This probably explains it. 🚀

We've seen strong continued growth in #AI and #ML for science across a broad set of domains including materials science, chemistry, physics and more.



Year	Materials Science (%)	Chemistry (%)	Physics (%)
2000	0.4	0.4	0.15
2005	0.25	0.55	0.1
2010	0.15	0.6	0.1
2015	0.15	0.55	0.15
2020	1.4	2.2	1.9
2022	2.8	3.0	3.6

ALT

1 2 9 309

 **Ben Blaiszik** @BenBlaiszik · 45m

Replies to [@BenBlaiszik](#)

Computing the YoY growth rates and CAGR (details available in the repo) shows the following.

Percentage Gains in # of matching articles for 2022:

- Materials Science: 39% more articles than 2021
- Chemistry: 27% more articles than 2021
- Physics: 29% more articles than 2021

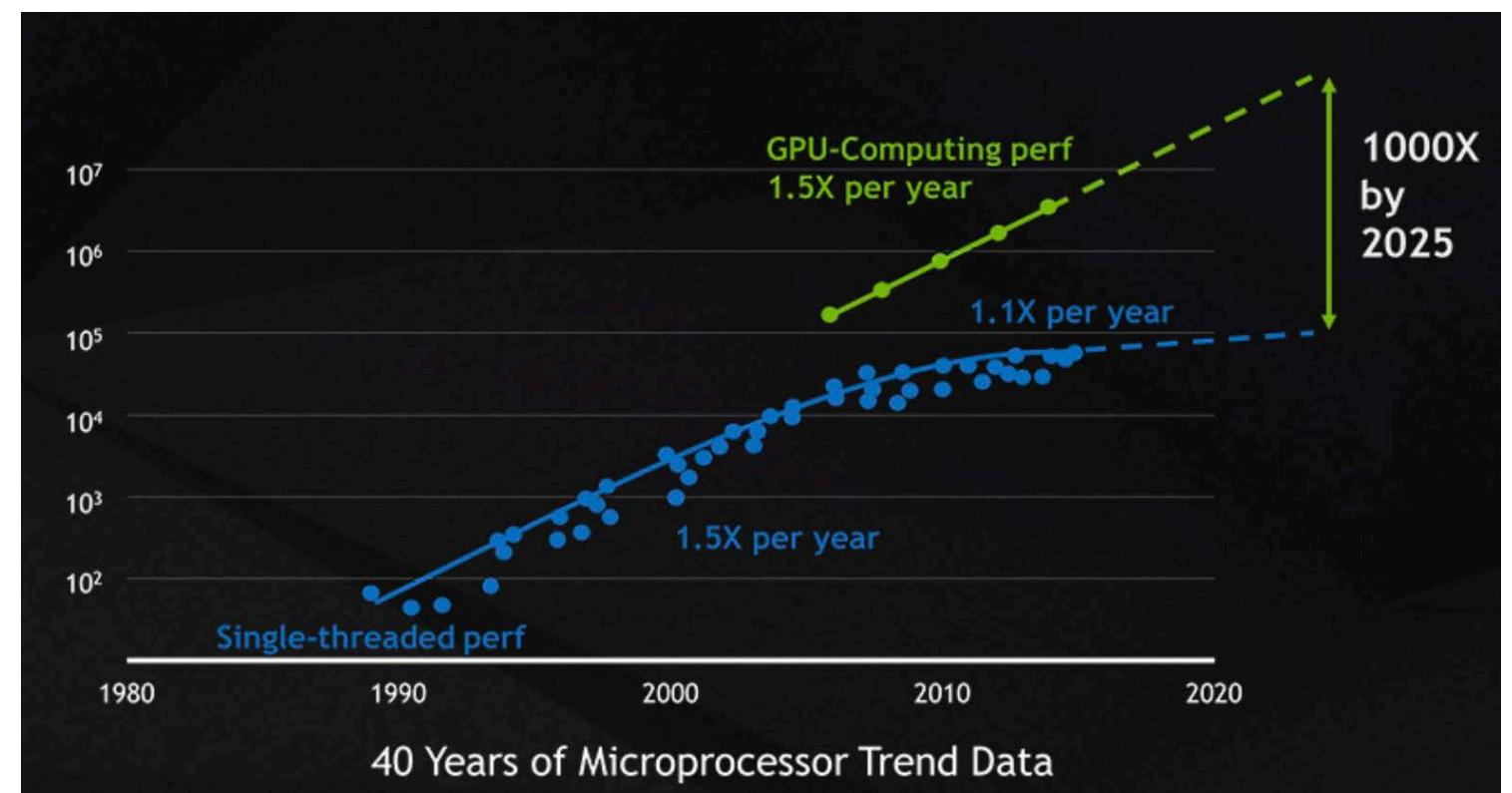


Domain	year	count	CAGR-1 (%)
Materials Science	2022	6180	39.1
Chemistry	2022	8842	27.4
Physics	2022	6829	29.3

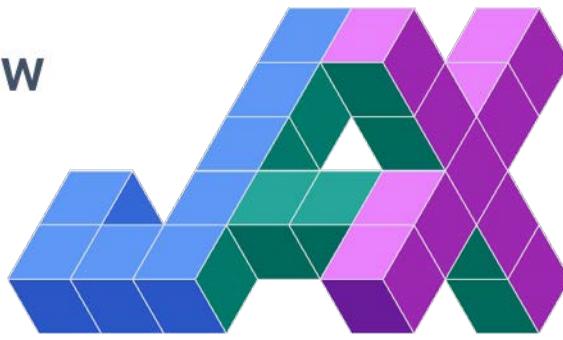
ALT

1 58

Why now?



TensorFlow



PyTorch



Keras



IMAGENET



PDB
PROTEIN DATA BANK

201,979 Structures from the PDB
1,068,577 Computed Structure Models (CSM)

The Pile An 800GB Dataset of Diverse Text for Language Modeling

PyTorch Lightning



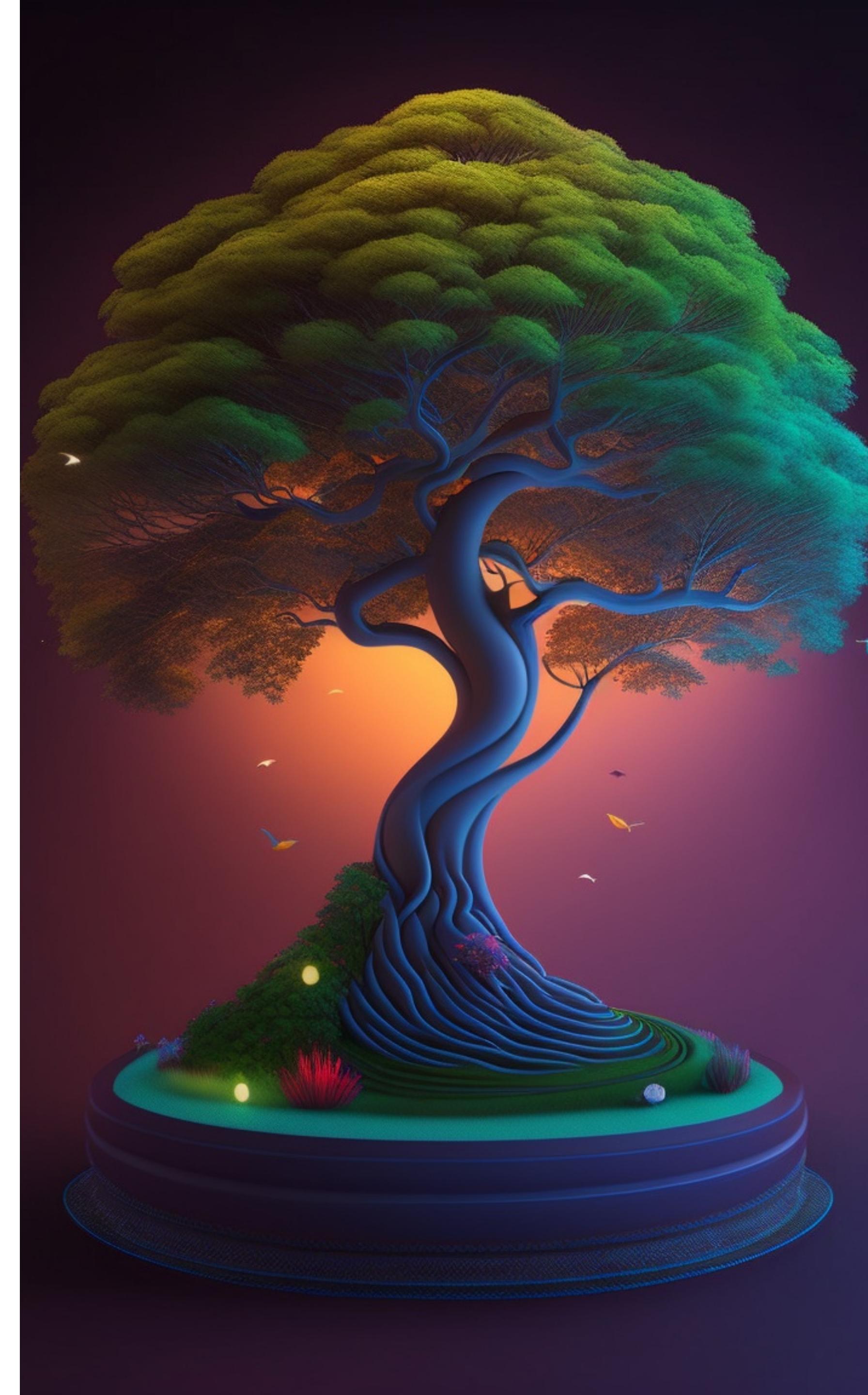
Better performing hardware (GPUs)

Accessible ML frameworks

Open data & code

What is public data?

- openly accessible
- free of charge
- and can be used, reused, and redistributed by anyone
- often with few or no restrictions
- <https://creativecommons.org/choose/>



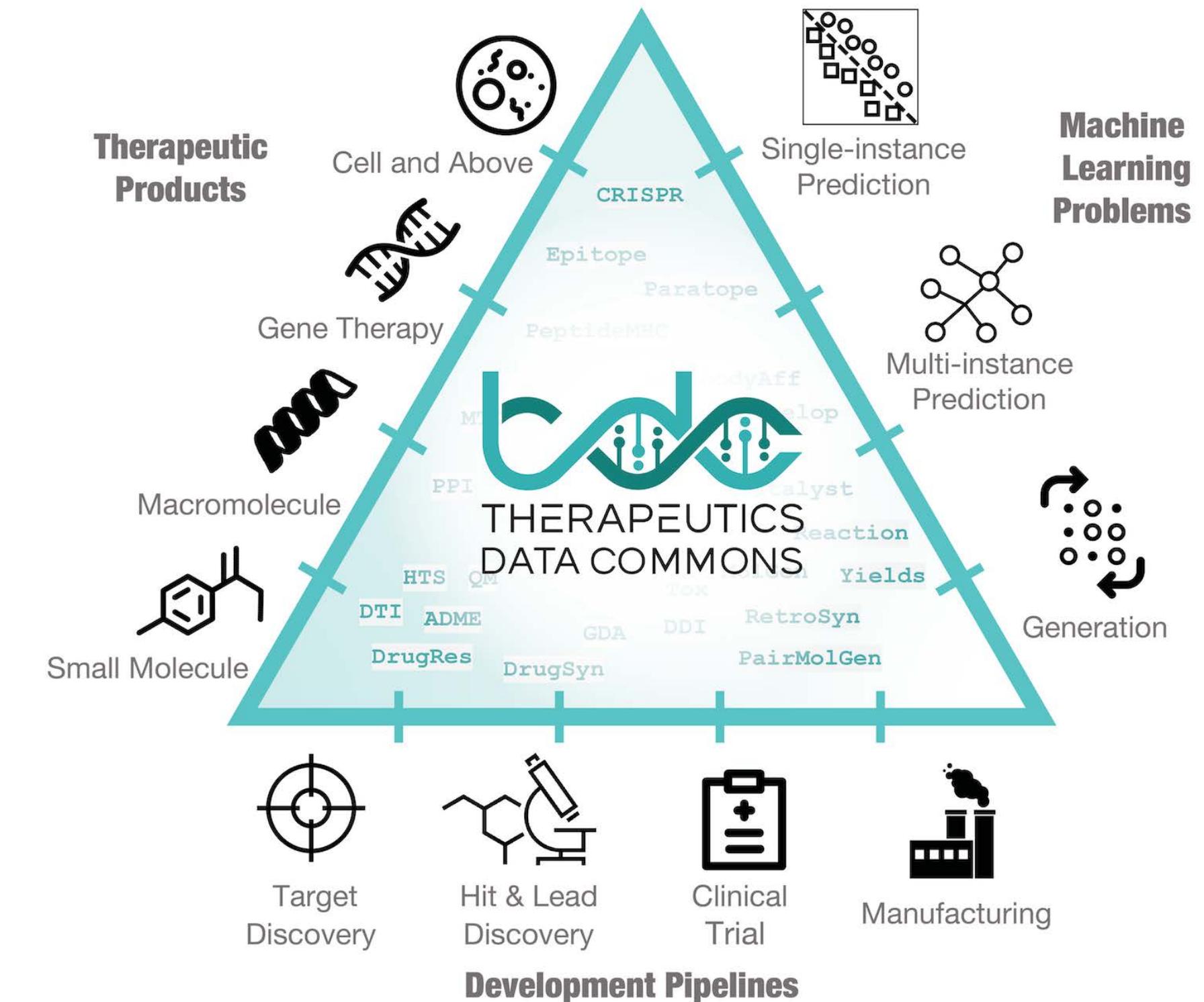
Open research data's significance

- promotes transparency and trust
- accountability in research
- fosters collaboration among researchers
- enabling researchers to build upon each other's work
- accelerating research & scientific advancements

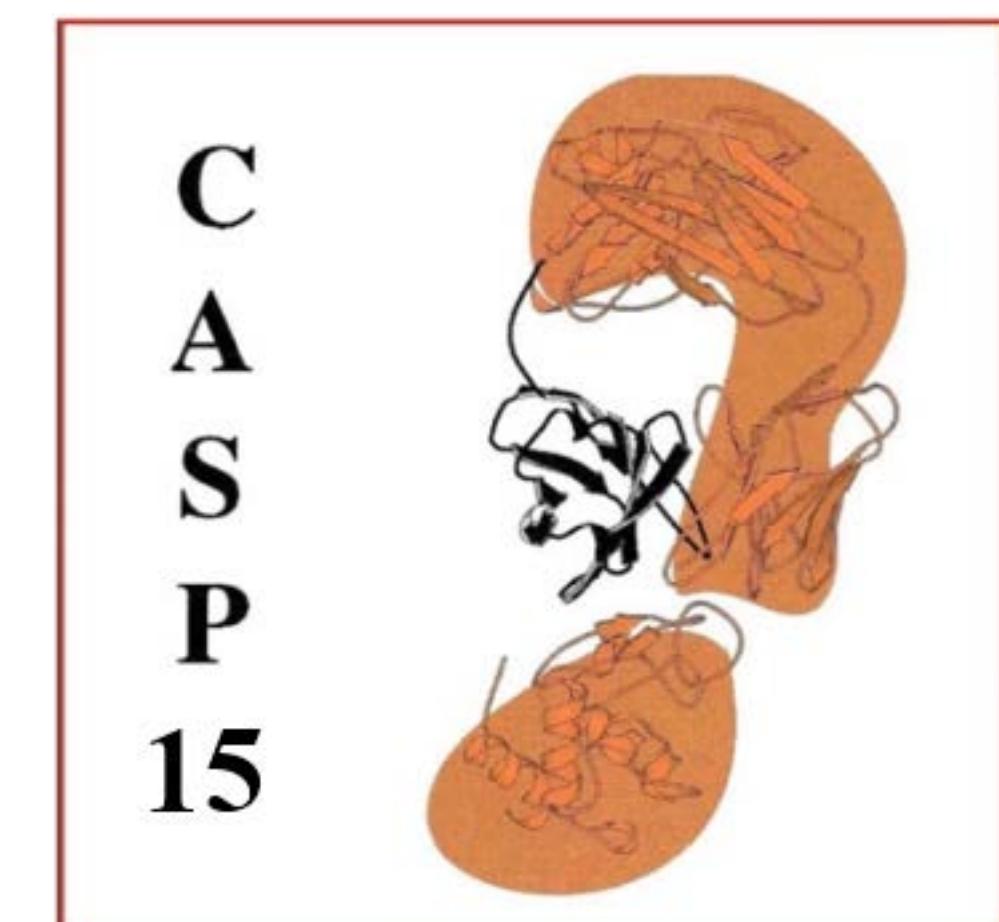


Standardized benchmarks

- compare the performance of different methods, algorithms, or models
- curated datasets and evaluation metrics
- assess strengths and weaknesses of approaches



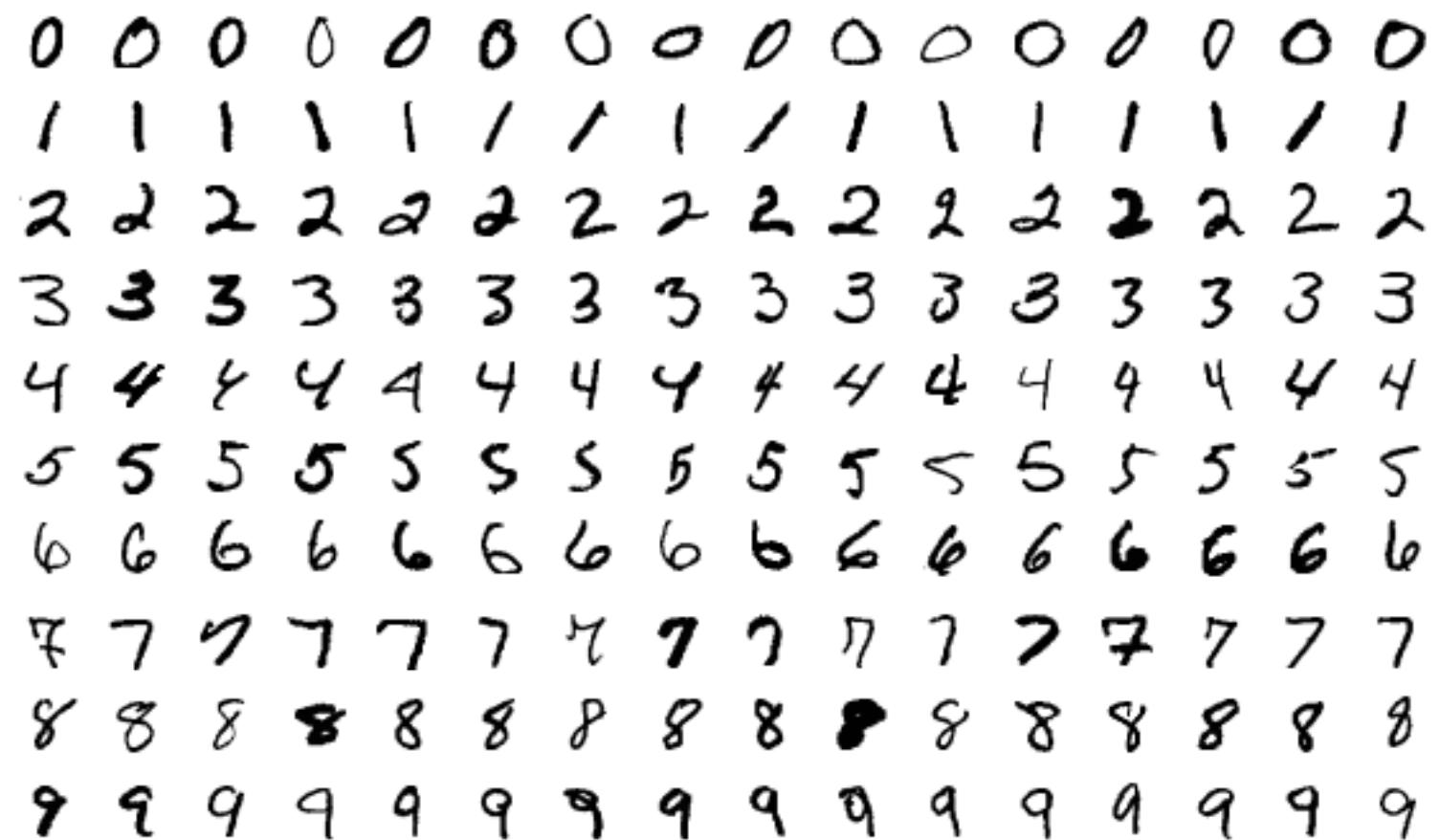
- Important to be as close as possible to reality



QUANTUM-MACHINE.ORG
[Home](#) | [Datasets](#) | [Publications](#) | [Software](#)

QM9

Open data that initiated deep learning wave¹⁰



THE MNIST DATABASE

of handwritten digits

[Yann LeCun](#), Courant Institute, NYU

[Corinna Cortes](#), Google Labs, New York

[Christopher J.C. Burges](#), Microsoft Research, Redmond

1994, 60k images
6k citations



ImageNet: A Large-Scale Hierarchical Image Database

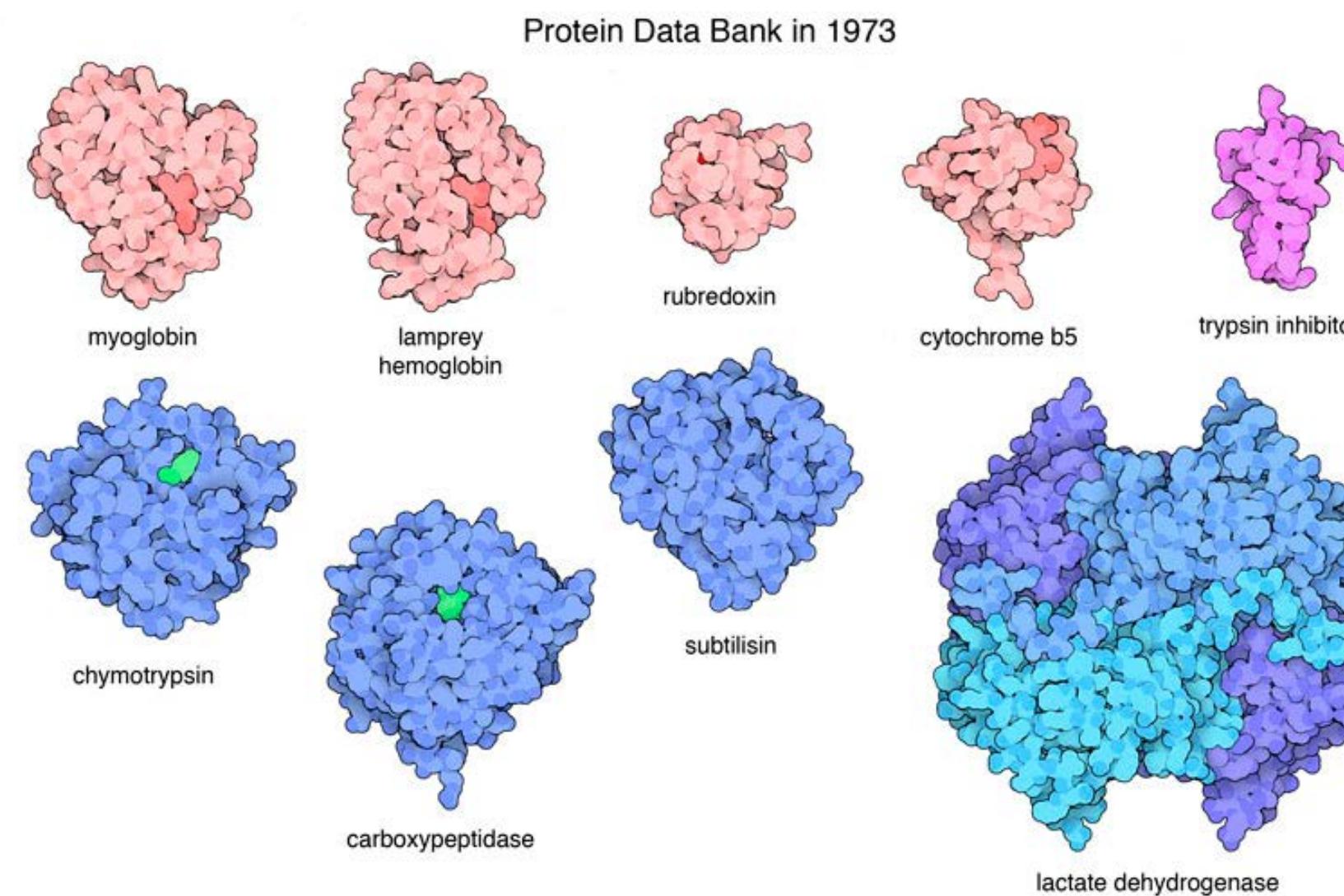
Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li and Li Fei-Fei
Dept. of Computer Science, Princeton University, USA

{jiadeng, wdong, rsocher, jial, li, feifeili}@cs.princeton.edu

2009, now 14M images
53k citations

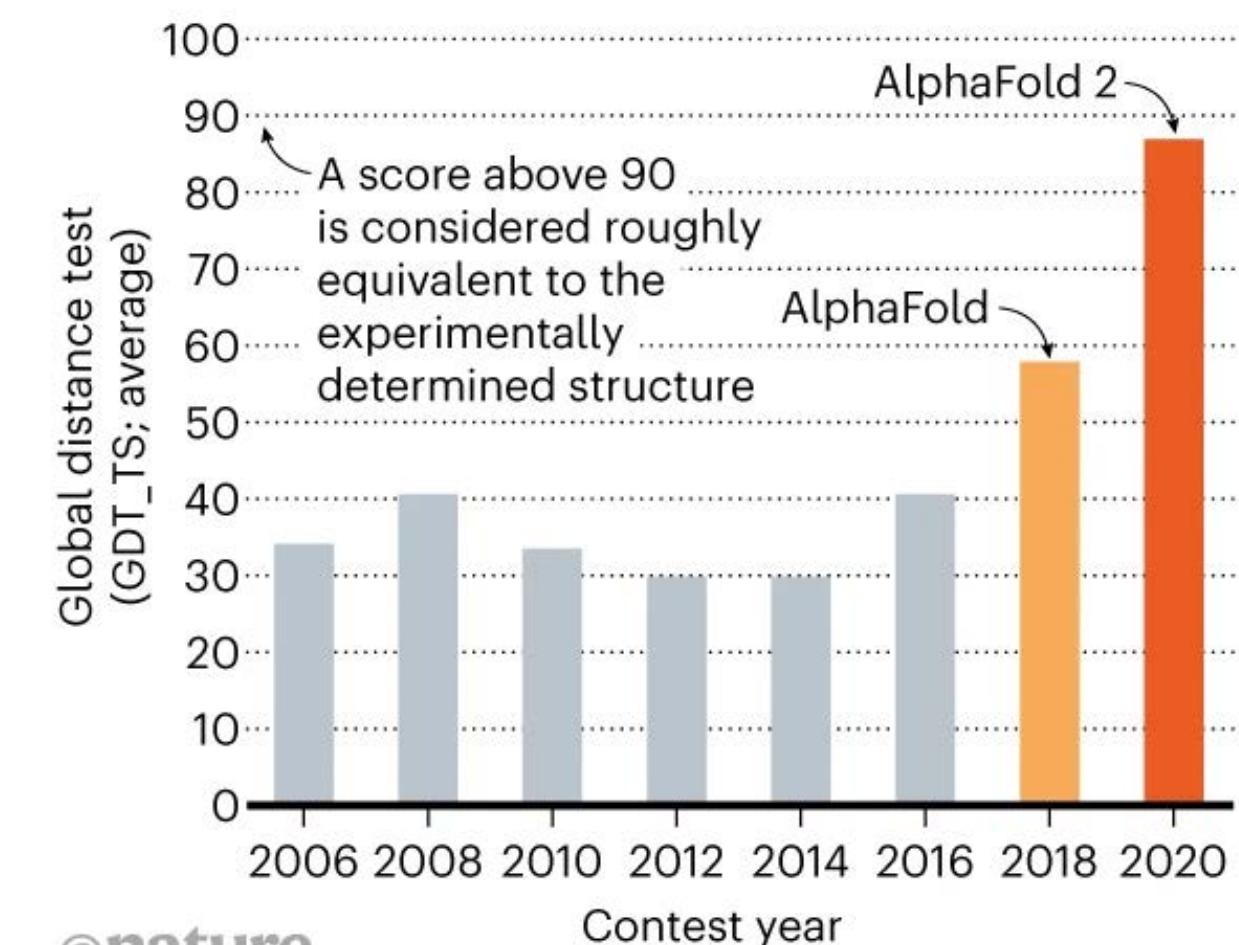
Deepmind's AlphaFold

- Won CASP protein-solving contest
- Made possible by the PDB database (established in 1971)



STRUCTURE SOLVER

DeepMind's AlphaFold 2 algorithm significantly outperformed other teams at the CASP14 protein-folding contest — and its previous version's performance at the last CASP.



Open datasets

- Paper with Code
<https://paperswithcode.com/datasets>
- HuggingFace datasets
<https://huggingface.co/datasets>

Best match

8101 dataset results

CIFAR-10
The CIFAR-10 dataset (Canadian Institute for Advanced Research, 10 classes) is a subset of the Tiny Images dataset and consists of 60000 32x32 color images. The images are labelled...
11,117 PAPERS • 69 BENCHMARKS

ImageNet
The ImageNet dataset contains 14,197,122 annotated images according to the WordNet hierarchy. Since 2010 the dataset is used in the ImageNet Large Scale Visual Recognition Chal...
10,726 PAPERS • 104 BENCHMARKS

COCO (Microsoft Common Objects in Context)
The MS COCO (Microsoft Common Objects in Context) dataset is a large-scale object detection, segmentation, key-point detection, and captioning dataset. The dataset consists of...
7,691 PAPERS • 80 BENCHMARKS

MNIST
The MNIST database (Modified National Institute of Standards and Technology database) is a large collection of handwritten digits. It has a training set of 60,000 examples, and a test se...
6,103 PAPERS • 51 BENCHMARKS

Filter by Modality

Images	2330
Texts	2149
Videos	756
Audio	468
...	...

Filter by Task

Question Answering	334
Semantic Segmentation	264
Object Detection	235
Image Classification	216

Tasks Sizes Sub-tasks Languages Licenses Other

Filter Tasks by name

Datasets 32,846 Filter by name new Full-text search Sort: Most Downloads

allenai/nllb	Preview Updated Sep 29, 2022 ↓ 1.96M 46
glue	Preview Updated about 1 month ago ↓ 1.02M 158
google/MusicCaps	Preview Updated Mar 8 ↓ 984k 49
piva	Updated Jan 25 ↓ 939k 20
EleutherAI/lambada_openai	Preview Updated Dec 16, 2022 ↓ 873k 17
sciq	Updated about 1 month ago ↓ 810k 21
wikitext	Preview Updated about 1 month ago ↓ 506k 124
super_glue	Preview Updated about 1 month ago ↓ 454k 86
red_caps	Preview Updated Jan 25 ↓ 363k 29
kilt_tasks	Preview Updated Nov 18, 2022 ↓ 354k 18
openwebtext	Preview Updated about 1 month ago ↓ 349k 119
imdb	Preview Updated about 1 month ago ↓ 255k 79
Helsinki-NLP/tatoeba_mt	Preview Updated Oct 21, 2022 ↓ 183k 25
squad	Updated about 1 month ago ↓ 176k 81

Multimodal

- Feature Extraction Text-to-Image
- Image-to-Text Text-to-Video
- Visual Question Answering
- Graph Machine Learning

Computer Vision

- Depth Estimation Image Classification
- Object Detection Image Segmentation
- Image-to-Image
- Unconditional Image Generation
- Video Classification
- Zero-Shot Image Classification

Natural Language Processing

- Text Classification Token Classification

**Without open data I would not
stand in front of you today.**

A wide-angle photograph of a dark night sky filled with stars. The Milky Way galaxy is prominent, its bright central band and darker, more diffuse outer regions stretching across the frame. The horizon is visible at the bottom, showing a silhouette of mountain peaks against a yellowish-orange glow from a nearby town or city lights.

Exploring the nearly
*endless chemical
space*

10^{60} drug-like
molecules

What molecule to make?

How to make it?
Reaction prediction
Synthesis planning

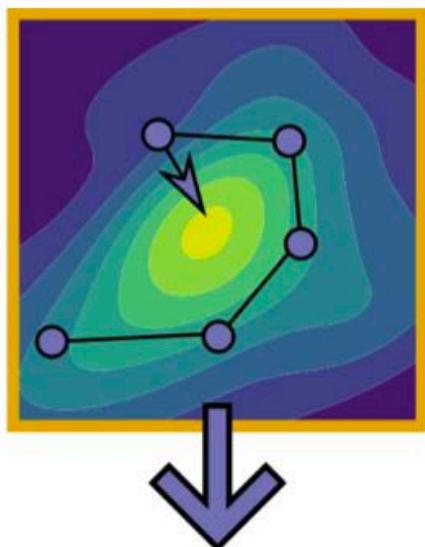
Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules

Rafael Gómez-Bombarelli[#] , Jennifer N. Wei[#] , David Duvenaud^{†#} , José Miguel Hernández-Lobato^{§#} , Benjamín Sánchez-Lengeling[‡] , Dennis Sheberla[‡] , Jorge Aguilera-Iparraguirre[†] , Timothy D. Hirzel[†] , Ryan P. Adams^{VI} , and Alán Aspuru-Guzik^{*†‡} 

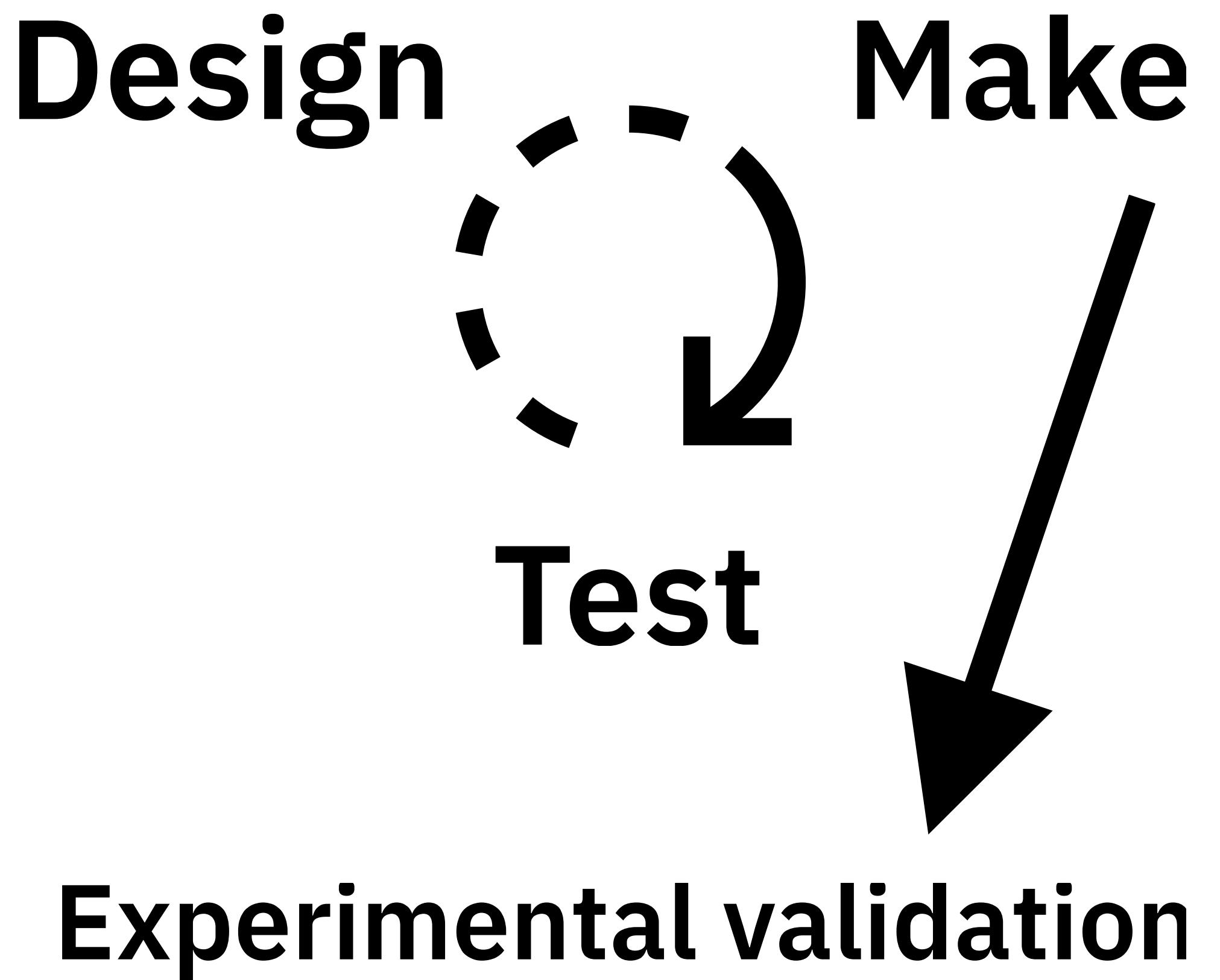
[Submitted on 5 Jan 2017]

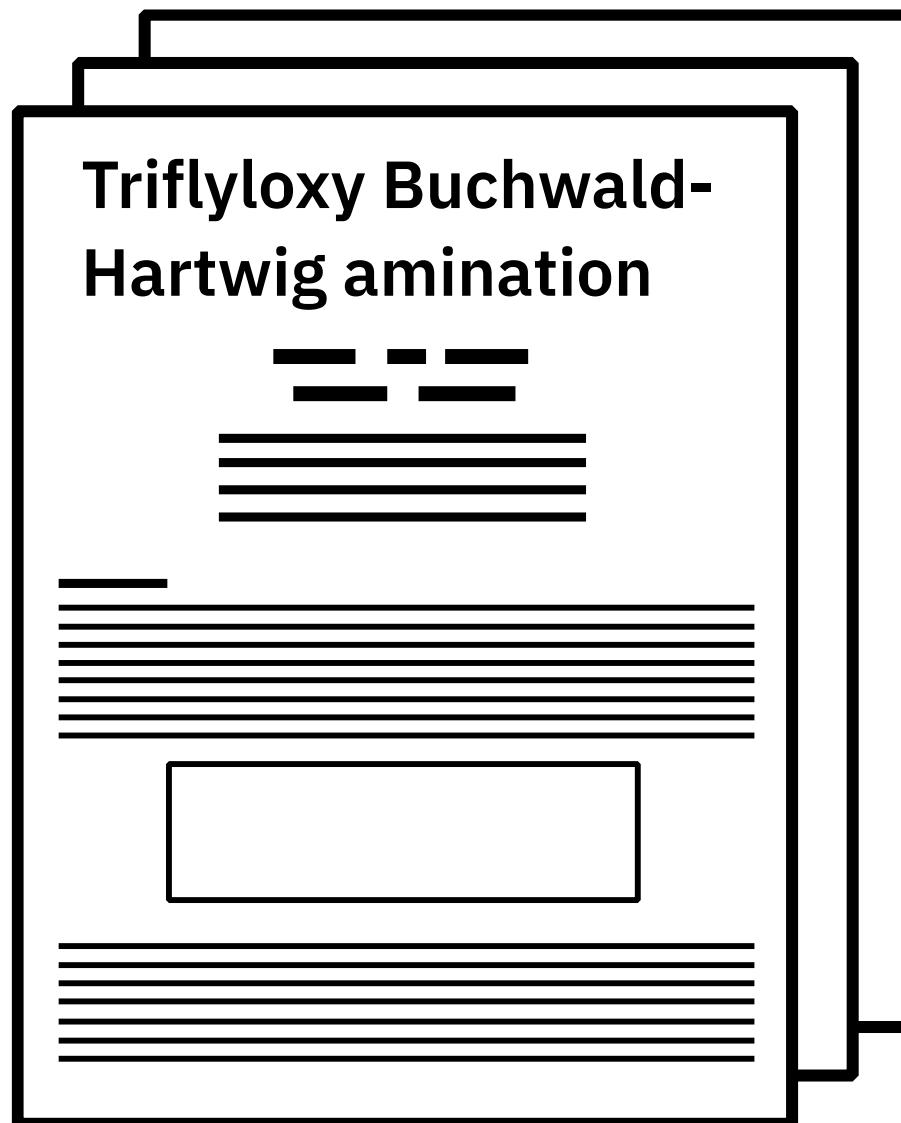
Generating Focussed Molecule Libraries for Drug Discovery with Recurrent Neural Networks

Marwin H.S. Segler, Thierry Kogej, Christian Tyrchan, Mark P. Waller

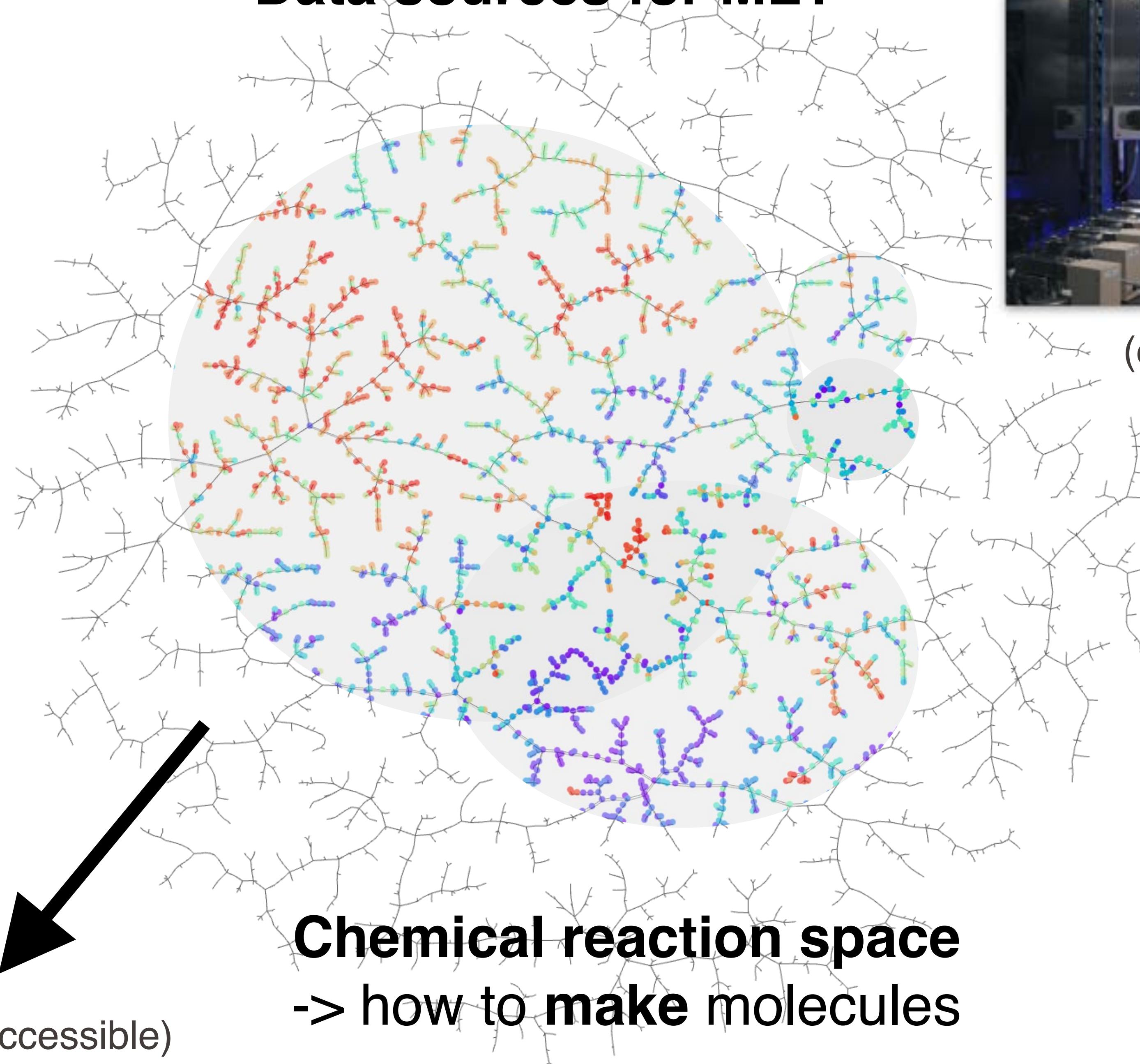


Inverse:
Optimization,
evolutionary
strategies,
generative models





Data sources for ML?



- Patents (broad, accessible)



(e.g. ORD, Kearnes et al.)

$$\hat{H}\Psi = E\Psi$$



- Simulations (narrow)

Open Data for Chemical Reactions

Chemical reactions from US patents (1976-Sep2016)

[Cite](#)[Download all \(1.39 GB\)](#)[Share](#)[Embed](#)[+ Collect](#)

Dataset posted on 2017-06-13, 18:49 authored by [Daniel Lowe](#)

Reactions extracted by text-mining from United States patents published between 1976 and September 2016. The reactions are available as CML or reaction SMILES. Note that the reactions SMILES are derived from the CML. The files can be unzipped using a program like 7-Zip.

The reactions were extracted using an enhanced version of the reaction extraction code described in
<https://www.repository.cam.ac.uk/handle/1810/244727>

with LeadMine (<https://www.nextmovesoftware.com/leadmine.html>) used for chemical entity recognition.

General tips:

Duplicate reactions are frequent due to the same or highly similar text occurring in multiple patents, this is especially true

USAGE METRICS

37653
views

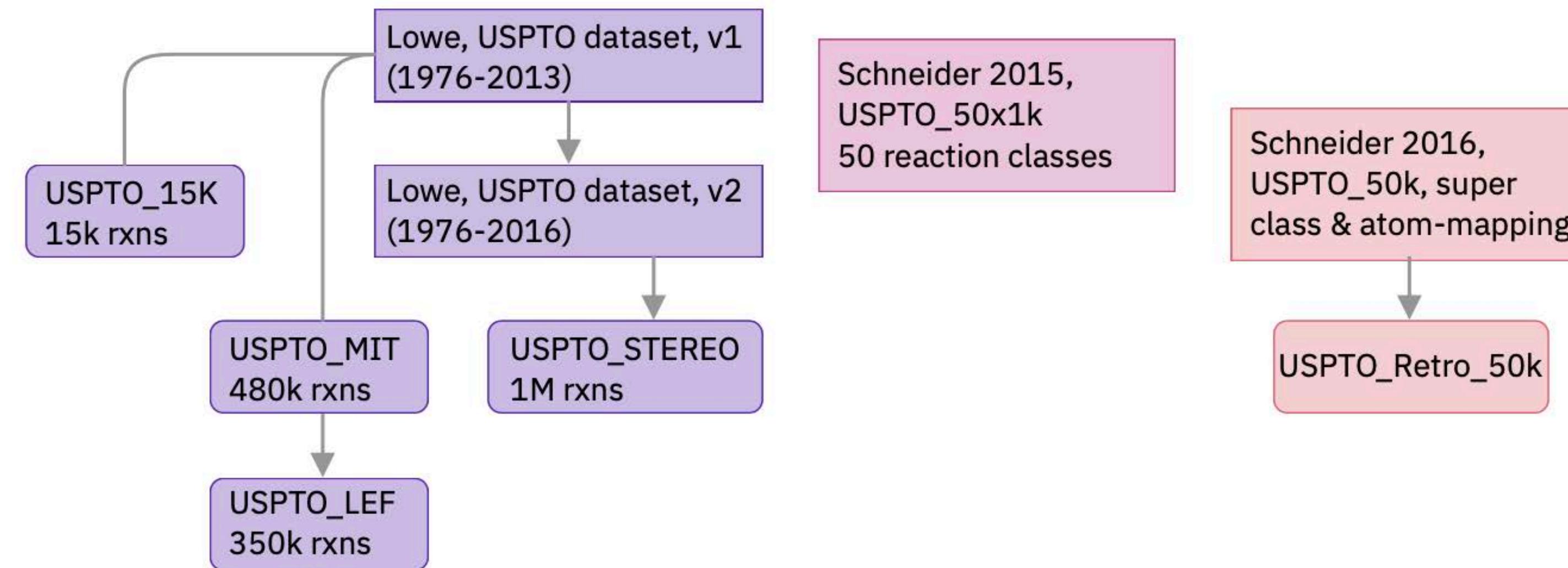
20303
downloads

19
citations



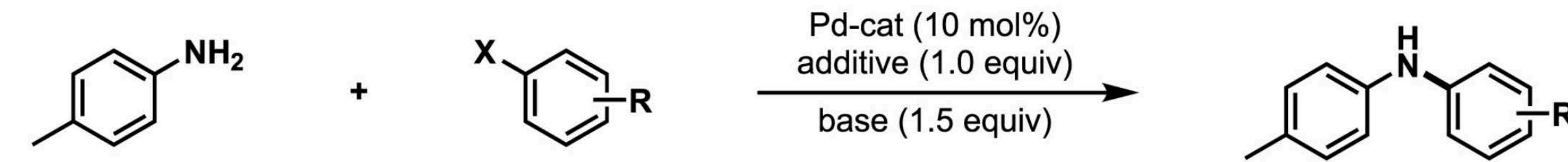
CATEGORIES

- Cheminformatics and quantitative structure-activity relationships
- Organic chemical synthesis



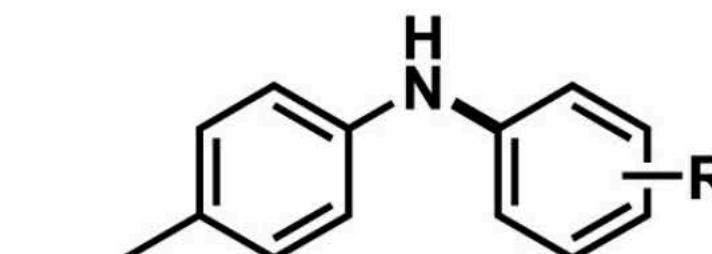
EPFL Open Data for Chemical Reactions

18



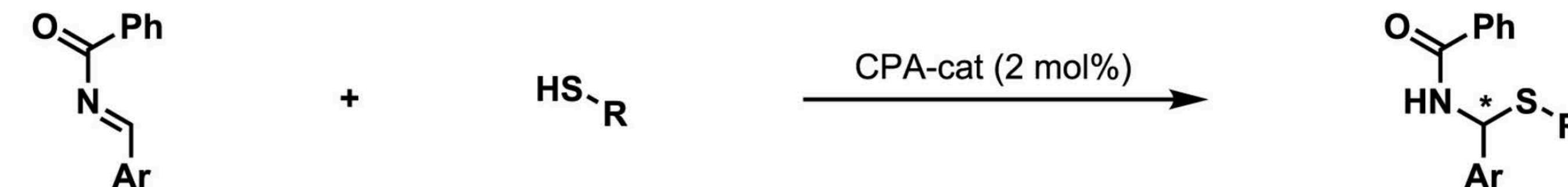
- 15 aryl halides
- 3 bases

- 4 Pd-catalysts
- 23 isoxazole additives



total reactions: 4140

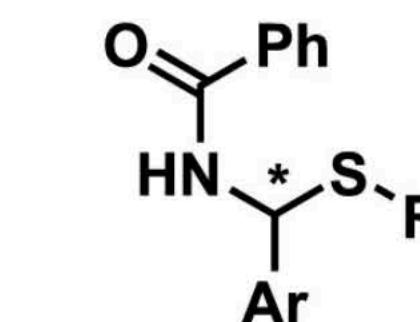
C–N cross coupling reactions of 4-methylaniline with various aryl halides by Doyle and co-workers.



- 5 imines

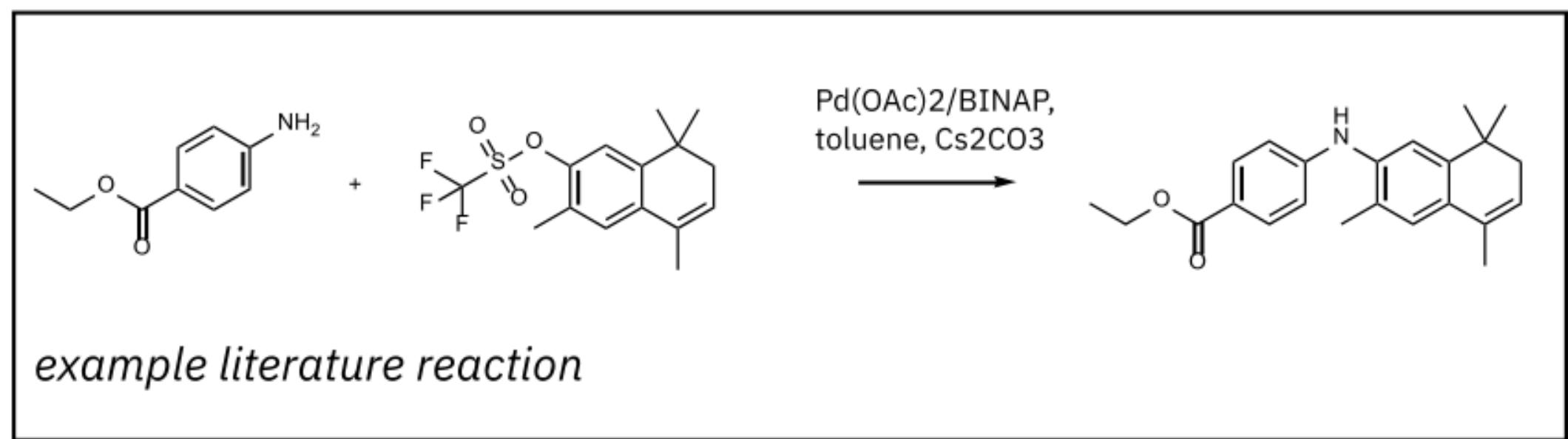
- 5 thiols

- 43 CPA catalysts



total reactions: 1075

Asymmetric N,S -acetal formation using CPA catalysts by Denmark et al.³⁸



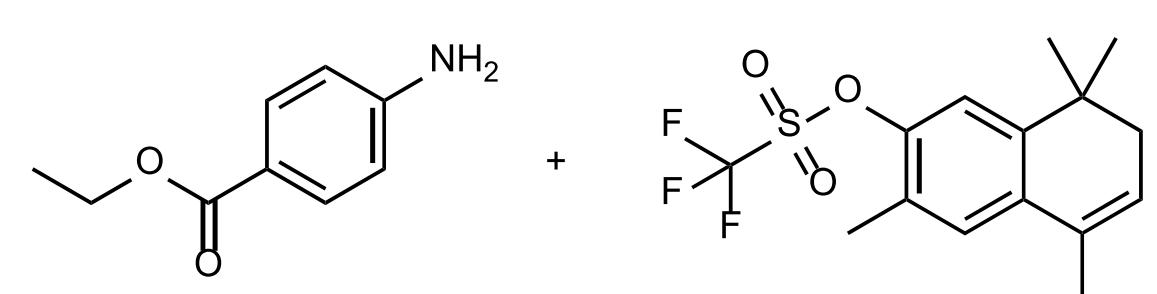
General

- USPTO: https://figshare.com/articles/dataset/Chemical_reactions_from_US_patents_1976-Sep2016_/_5104873
- Train/valid/test splits: https://github.com/coleygroup/Graph2SMILES/blob/main/scripts/download_raw_data.py
- 50k with classes and atom-mapping: <https://pubs.acs.org/doi/full/10.1021/acs.jcim.6b00564>
- 50k with 50 reaction classes: <https://pubs.acs.org/doi/full/10.1021/ci5006614>
- 1000 template classes: <https://rxn4chemistry.github.io/rxnfp>
- Open Reaction Database: <https://open-reaction-database.org/client/browse>

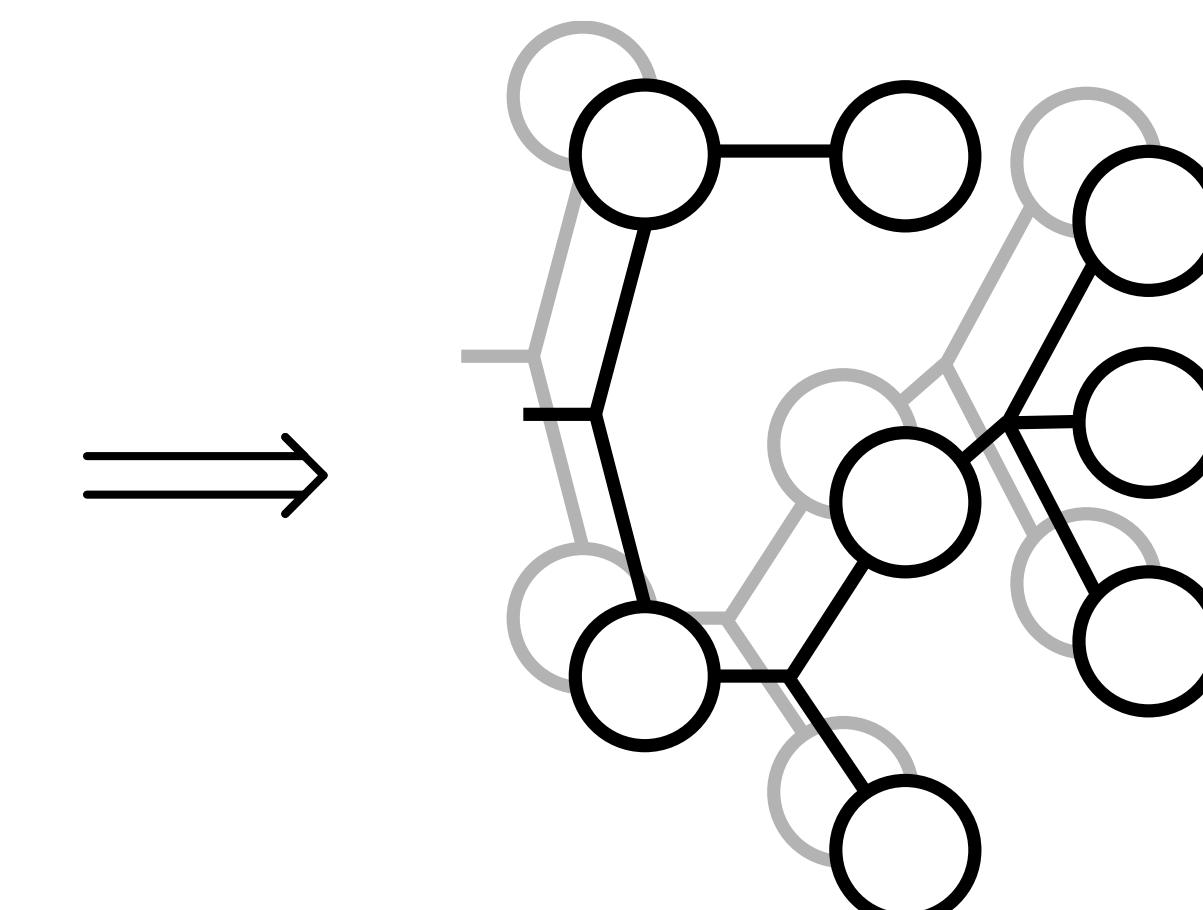
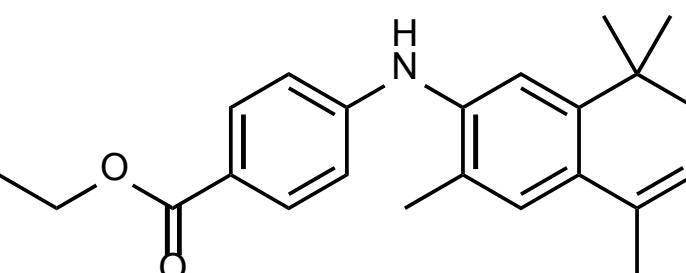
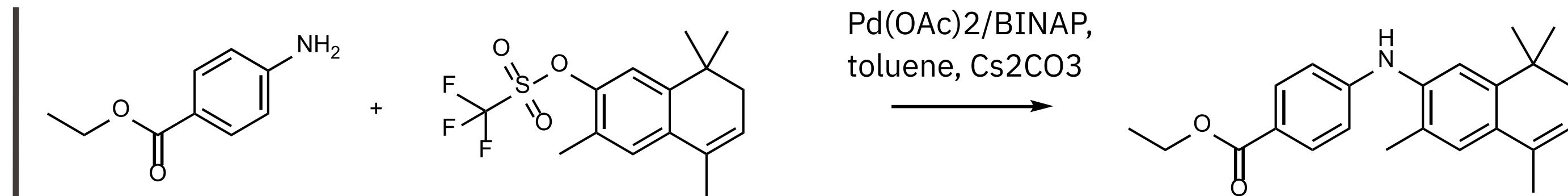
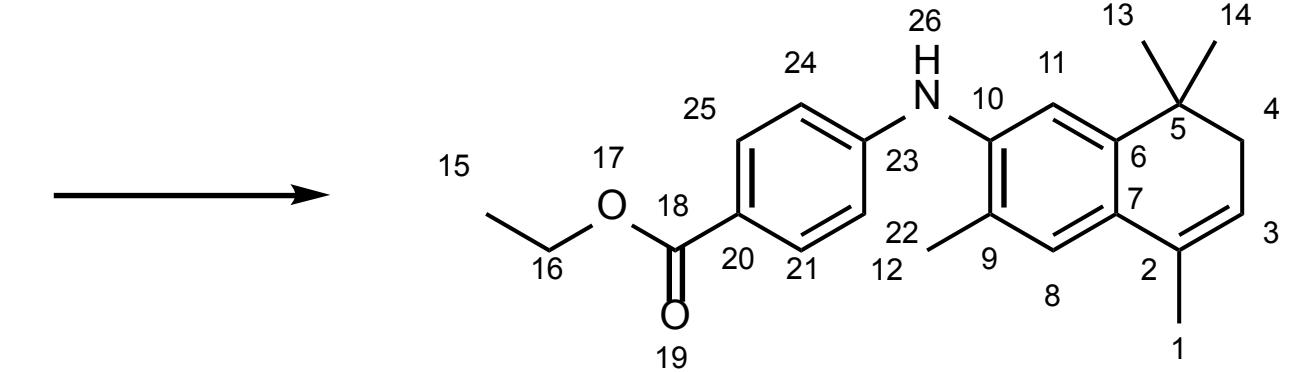
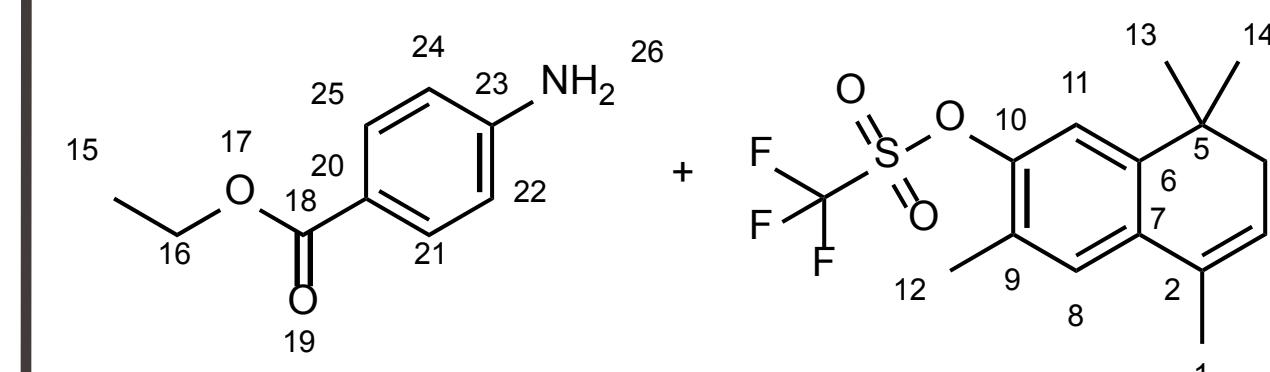
Specific reactions

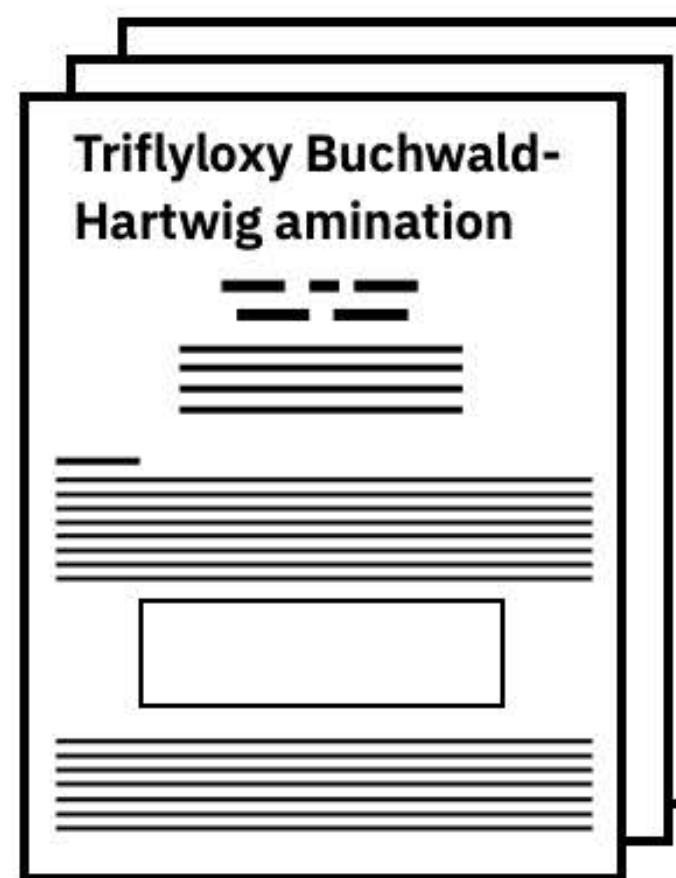
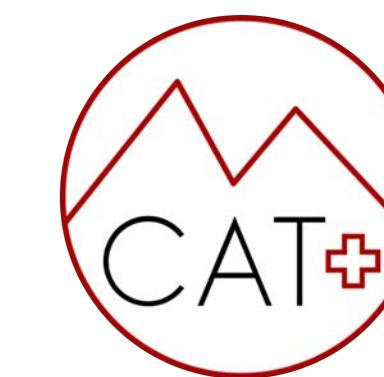
- 4k Buchwald-Hartwig reactions (yield): <https://www.science.org/doi/10.1126/science.aar5169>
- 1k Asymmetric N,S-acetal formation using CPA catalysts (enantioselectivity):
<https://www.science.org/doi/abs/10.1126/science.aau5631>

There are many more...

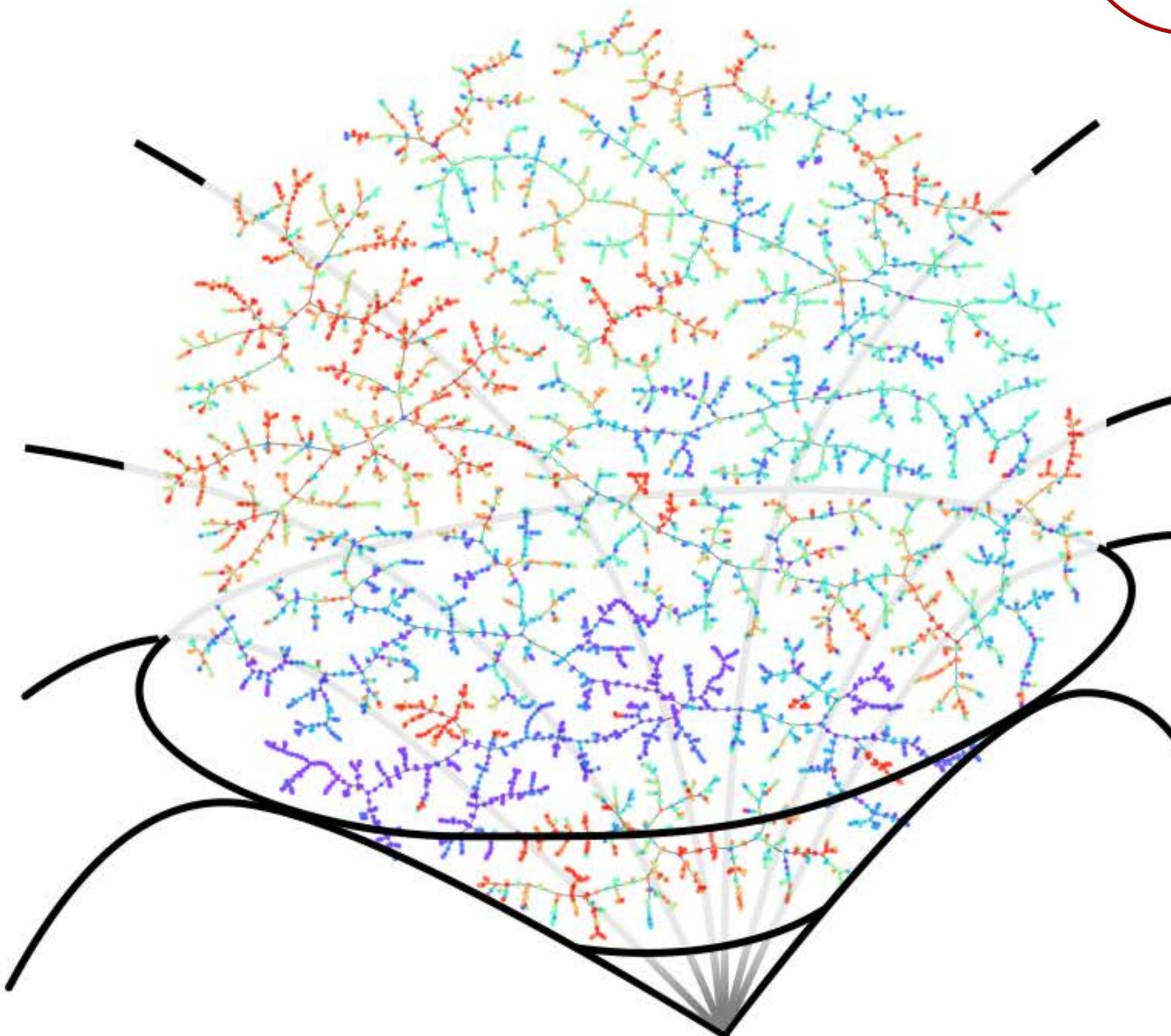


?

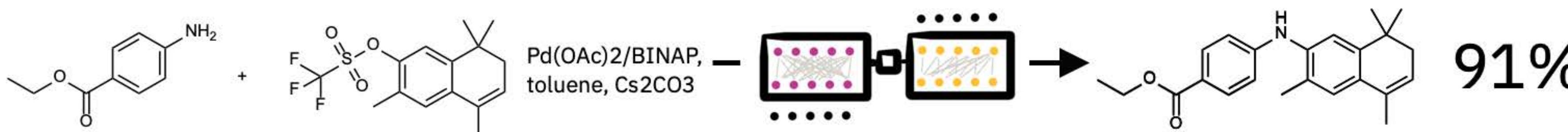
Reaction prediction:Schwaller et al. *ACS Cent. Sci.* 2019Pesciullesi et al. *Nature Comm.* 11, 4874
(2020)**Synthesis planning:** Schwaller et al.
Chem. Sci., 2020, 11, 3316-3325**Classification/fingerprints:** Schwaller et al. *Nature Mach. Intell.*, 2021**Reaction yields:** Schwaller et al. *Mach. Learn.: Sci. Technol.*, 2021; Schwaller et al. NeurIPS 2020 MI4Mol workshop**Atom-mapping:** Schwaller et al. *Science Advances*, 2021



US20030166932A1: General Procedure
H A solution of trifluoromethanesulfonic acid 3,5,8,8-tetramethyl-7,8-dihydronaphthalen-2-yl ester (Compound 35, 0.41 g, 1.2 mmol), Pd(OAc)₂ (0.027 g, 0.12 mmol), BINAP (0.11 g, 0.18 mmol), Cs₂CO₃ (0.56 g, 1.72 mmol), ethyl 4-aminobenzoate (0.25 g, 1.5 mmol) and 5 mL of toluene was flushed with argon for 10 min, then stirred at 100° C. in a sealed tube for 48 h. ...



$$\hat{H}\Psi = E\Psi$$



NCCR
Catalysis

Machine intelligence for chemical reaction space

Philippe Schwaller , Alain C. Vaucher, Ruben Laplaza, Charlotte Bunne, Andreas Krause, Clemence Corminboeuf, Teodoro Laino

Open Data in the NCCRs

- NCCR Catalysis collects data sets on a Zenodo page: <https://zenodo.org/communities/nccr-catalysis/?page=1&size=20>

November 30, 2022 (1.0.0) Dataset Open Access

View

Dataset for Reaction-Induced Formation of Stable Mononuclear Cu(I)Cl Species on Carbon for Low-Footprint Vinyl Chloride Production

View

Faust Akl, Dario; Giannakakis, Georgios; Ruiz-Ferrando, Andrea; Agrachev, Mikhail; Medrano-García, Juan D.; Guillén-Gosálbez, Gonzalo; Jeschke, Gunnar; Clark, Adam H.; Safonova, Olga V.; Mitchell, Sharon; López, Núria; Pérez-Ramírez, Javier;

This dataset complements the publication entitled "Reaction-Induced Formation of Stable Mononuclear Cu(I)Cl Species on Carbon for Low-Footprint Vinyl Chloride Production" by Dario Faust Akl, Georgios Giannakakis, Andrea Ruiz-Ferrando, Mikhail Agrachev, Juan D. Medrano-García, G

Uploaded on April 20, 2023



Romain Graux

February 27, 2023 (1.0.0) Dataset Open Access

View

Activation in the rate of oxygen release of Sr_{0.8}Ca_{0.2}FeO_{3-δ} through removal of secondary surface species with thermal treatment in a CO₂-free atmosphere

View

Luongo, Giancarlo; H. Bork, Alexander; M. Abdala, Paula; Wu, Yi-Hsuan; Kountoupi, Evgenia; Donat, Felix; R. Müller, Christoph;

Here we report the data presented in the article titled: "Activation in the rate of oxygen release of Sr_{0.8}Ca_{0.2}FeO_{3-δ} through removal of secondary surface species with thermal treatment in a CO₂-free atmosphere", Luongo, G., Bork., A., Abdala, P. A., Wu, Y., Kountoupi, E., Donat, F.

Uploaded on April 20, 2023

- Raw data is dumped on Zenodo
- No direct link to the compounds
- No machine-actionable reaction data
- Question: How can we do better?

Raw data for the article " Pd(II)-Catalyzed Aminoacetoxylation of Alkenes via Tether Formation"

Raw NMR and MS data for the article "Pd(II)-Catalyzed Aminoacetoxylation of Alkenes via Tether Formation" published in Organic Letters, DOI:

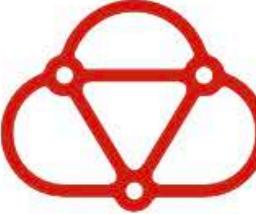
<https://doi.org/10.1021/acs.orglett.2c01838>

The number of the folders correspond to compounds numbers in the article. All details concerning conditions and equipment for measurements can be found in the supporting information of the article.

Name	Date Modified
> 1a	Today at 14:07
> 1a'	Today at 14:07
> 1b	Today at 14:07
> 1c	Today at 14:07
> 1d	Today at 14:07
> 1e	Today at 14:07
> 1f	Today at 14:07
> 1g	Today at 14:07
> 1h	Today at 14:07
> 1i	Today at 14:07
> 1j	Today at 14:07

EPFL Materials Cloud from NCCR Marvel

26



MATERIALSCLOUD



LEARN WORK DISCOVER EXPLORE ARCHIVE

Built for seamless sharing of resources in computational materials science.

LEARN Lectures and tutorials in computational materials science

WORK Simulation tools and services - in the cloud or on your computer

DISCOVER Curated research data with tailored visualizations

EXPLORE Interactive browser for AiiDA provenance graphs

ARCHIVE An open-access, moderated repository for research data in computational materials science

Please cite [L. Talirz et al., Sci Data 7, 299 \(2020\)](#), if you use Materials Cloud in your research.

Latest records



Single-model uncertainty quantification in neural network potentials does not consistently outperform model ensembles

DOI [10.24435/materialscloud:55-sd](#)

Aik Rui Tan, Shingo Urata, Samuel Goldman, Johannes C. B. Dietschreit, Rafael Gómez-Bombarelli

Neural networks (NNs) often assign high confidence to their predictions, even for points far out-of-distribution, making uncertainty quantification (UQ) a challenge. When they are employed to model interatomic potentials in materials systems, this problem leads to unphysical structures that disrupt simulations, or to biased statistics and dynamics that do not reflect the true physics. Differentiable UQ techniques can find new informative data and drive active learning loops for robust potentials. However, a variety of UQ techniques, including newly developed ones, exist for atomistic simulations and there are no clear guidelines for which are most effective or suitable for a given case. In this work, we examine multiple UQ schemes for improving the robustness of NN interatomic potentials (NNIPs) through active learning. In particular, we compare incumbent ensemble-based methods against strategies that use single, deterministic NNs: mean-variance estimation, deep evidential ...

Latest version: v1
Publication date: May 04, 2023

MARVEL



Mechanism and prediction of hydrogen embrittlement in fcc stainless steels and high entropy alloys

DOI [10.24435/materialscloud:ct-x8](#)

Xiao Zhou, Ali Tehranchi, W.A. Curtin

The urgent need for clean energy coupled with the exceptional promise of hydrogen (H) as a clean fuel is driving development of new metals resistant to hydrogen embrittlement. Experiments on new fcc high entropy alloys present a paradox: these alloys absorb more H than Ni or SS304 (austenitic 304 stainless steel) while being more resistant to embrittlement. Here, a new theory of embrittlement in fcc metals is presented based on the role of H in driving an intrinsic ductile-to-brittle transition at a crack tip. The theory quantitatively predicts the H concentration at which a transition to embrittlement occurs in good agreement with experiments for SS304, SS316L, CoCrNi, CoNiV, CoCrFeNi, and CoCrFeMnNi. The theory rationalizes why CoNiV is the alloy most resistant to embrittlement and why SS316L is more resistant than the high entropy alloys CoCrFeNi and CoCrFeMnNi, which opens a path for the computationally guided discovery of new embrittlement-resistant alloys.

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾[nature](#) > [nature chemistry](#) > [perspectives](#) > [article](#)Perspective | [Published: 04 April 2022](#)

Making the collective knowledge of chemistry open and machine actionable

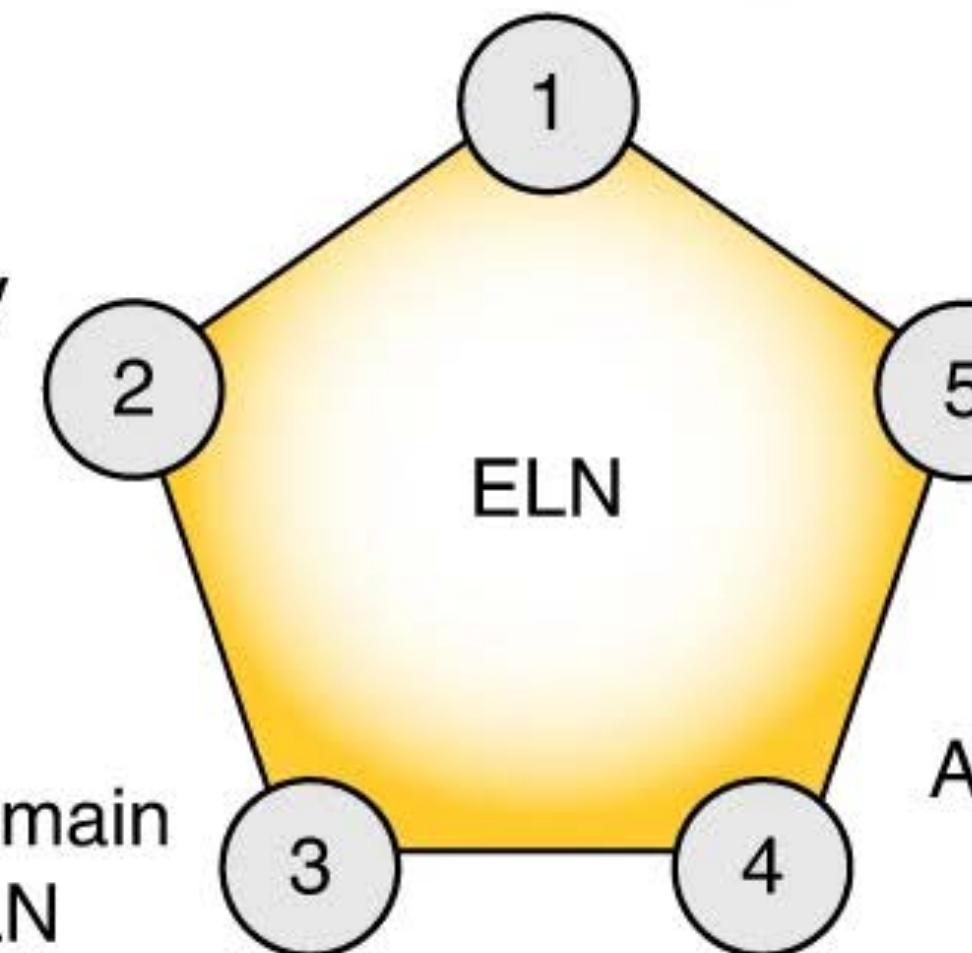
[Kevin Maik Jablonka](#), [Luc Patiny](#)✉ & [Berend Smit](#)✉[Nature Chemistry](#) 14, 365–376 (2022) | [Cite this article](#)

14k Accesses | 13 Citations | 118 Altmetric |



Kevin Jablonka

Data-intensive research
needs machine-actionable FAIR data;
'Nullius in verba'



ELN should automatically
make all data
machine-actionable

Data should not remain
locked in an ELN

Continuously inventing
new standards
will not make chemistry FAIR.
Existing systems
should be made interoperable

An open science infrastructure
needs to be modular
and open source

MARVEL

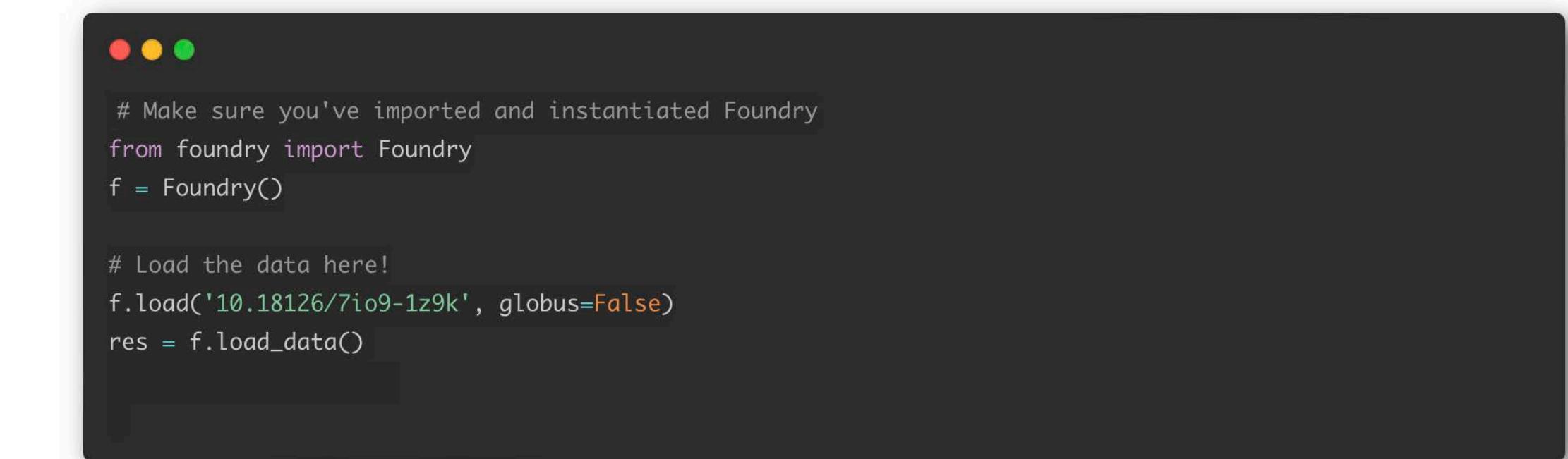
cheminfo

- Better metadata / links
- Reusability of datasets
 - Have associated code to load and process the data
 - <https://foundry-ml.org/>
 - https://github.com/MLMI2-CSSI/foundry/blob/main/examples/bandgap/bandgap_demo.ipynb

Using this dataset

First, you'll need to install the [latest version](#) of Foundry. You can do this with the command:

```
pip install foundry-ml
```



```
...  
# Make sure you've imported and instantiated Foundry  
from foundry import Foundry  
f = Foundry()  
  
# Load the data here!  
f.load('10.18126/7io9-1z9k', globus=False)  
res = f.load_data()
```

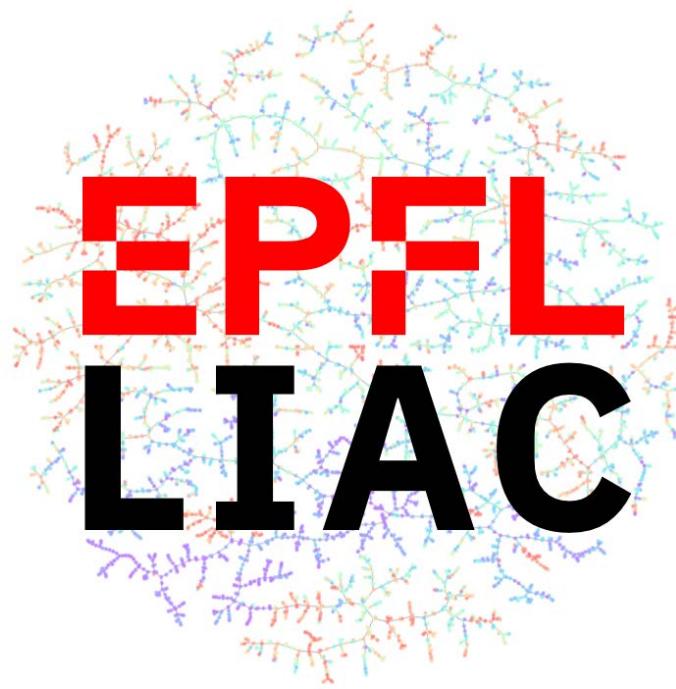
You can load this dataset with 2 lines of code if you already have Foundry set up. If you need to set up Foundry, check out our [example notebooks](#) and [documentation](#) for how to get started.

[GET DATA WITH GLOBUS](#)

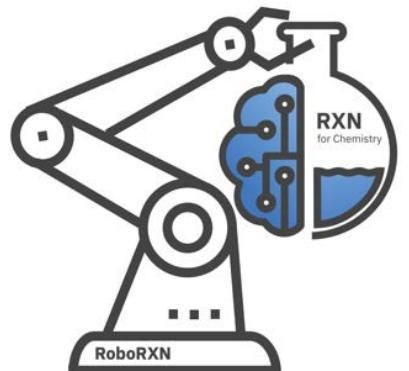
Take home messages

- Open Research Data is crucial
 - Scientific advancement
 - Reliability of results
 - Visibility and citations
 - Greater impact
- There is progress, but we should improve
- Make it accessible through
 - Well-documented linked code
 - Examples
 - Metadata / links





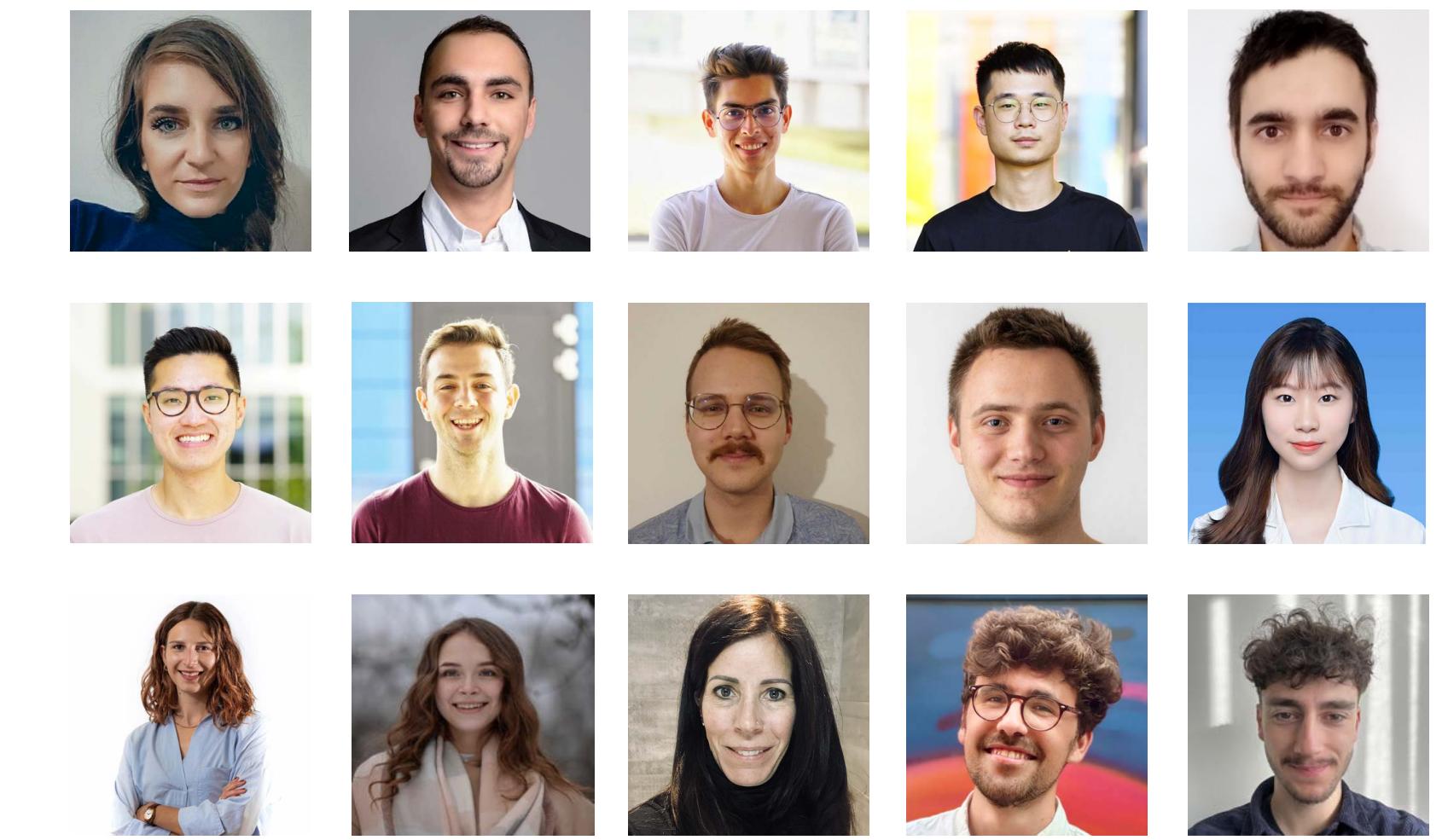
IBM Research AI



u^b
UNIVERSITÄT
BERN



Bojana Rankovic, Andres CM Bran,
Junwu Chen, Jeff Guo,
Oliver Schilter (with IBM),
Victor S Gil (with Luterbacher),
Paulo Neves (with Janssen),
Stéphane D'Ascoli,
Rebecca Neeser (with Correia)



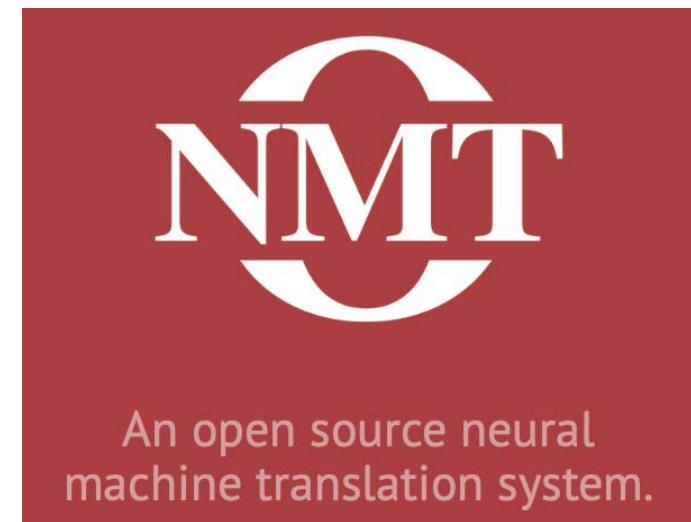
*Teodoro Laino, Alain Vaucher,
Alessandra Toniato, Theophile
Gaudin, Daniel Probst, Matteo
Manica and the RoboRXN team, Ben
Hoover, Hendrik Strobelt*

*Jean-Louis Reymond, Giorgio
Pesciullesi, David Kreutter, Alice
Capecci, Amol Thakkar, and the
Reymond group*

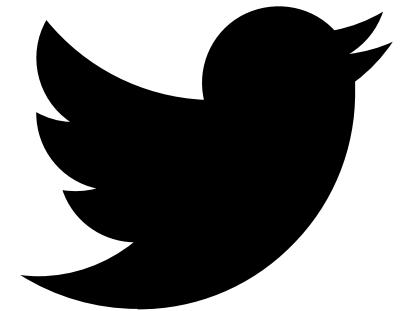
Alpha Lee, Peter Bolgar and Ryan-
Rhys Griffiths

Clemence Corminboeuf, Andreas
Krause, Ruben Laplaza, Charlotte
Bunne, Jeremy Luterbacher

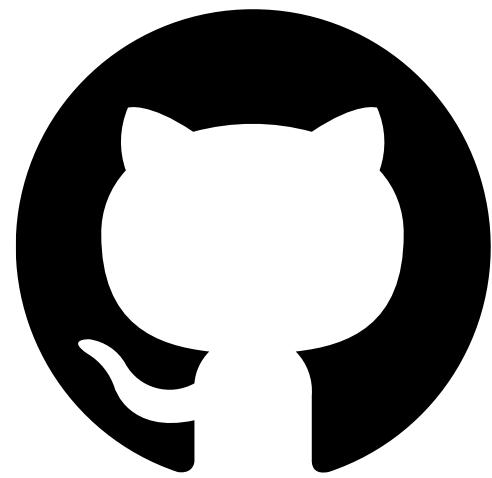
Kevin Jablonka and Berend Smit (EPFL)
Sam Cox and Andrew D White
(University of Rochester)



Many thanks for your attention!



@pschwllr @SchwallerGroup



<https://github.com/pschwllr>

<https://github.com/schwallergroup>



philippe.schwaller@epfl.ch

<https://schwallergroup.github.io>

**Time for questions
and discussion.**