

## **AI in Drug Discovery – An Overview**

### **Andrea Volkamer and Pat Walters**

September 16, 2024

# **Session 2 - Data is all you need!**

# Who we are!

---

**Pat Walters**  
Relay Therapeutics



**Raquel López-Ríos de Castro**  
Charité Berlin, MSKCC NYC



**Michael Backenköhler**  
Saarland University



**Andrea Volkamer**  
Saarland University



# What we will do today

---

## Session 1 - 1:30 - 2:30 pm

- An introduction to Artificial Intelligence (AI) and Machine Learning (ML)
- Molecular representations
- AI architectures

## Session 2 - 3:00 - 4:00 pm

- Importance of data quality for AI and ML
- Data (pre)processing
- Exploratory data analysis
- Applicability domains

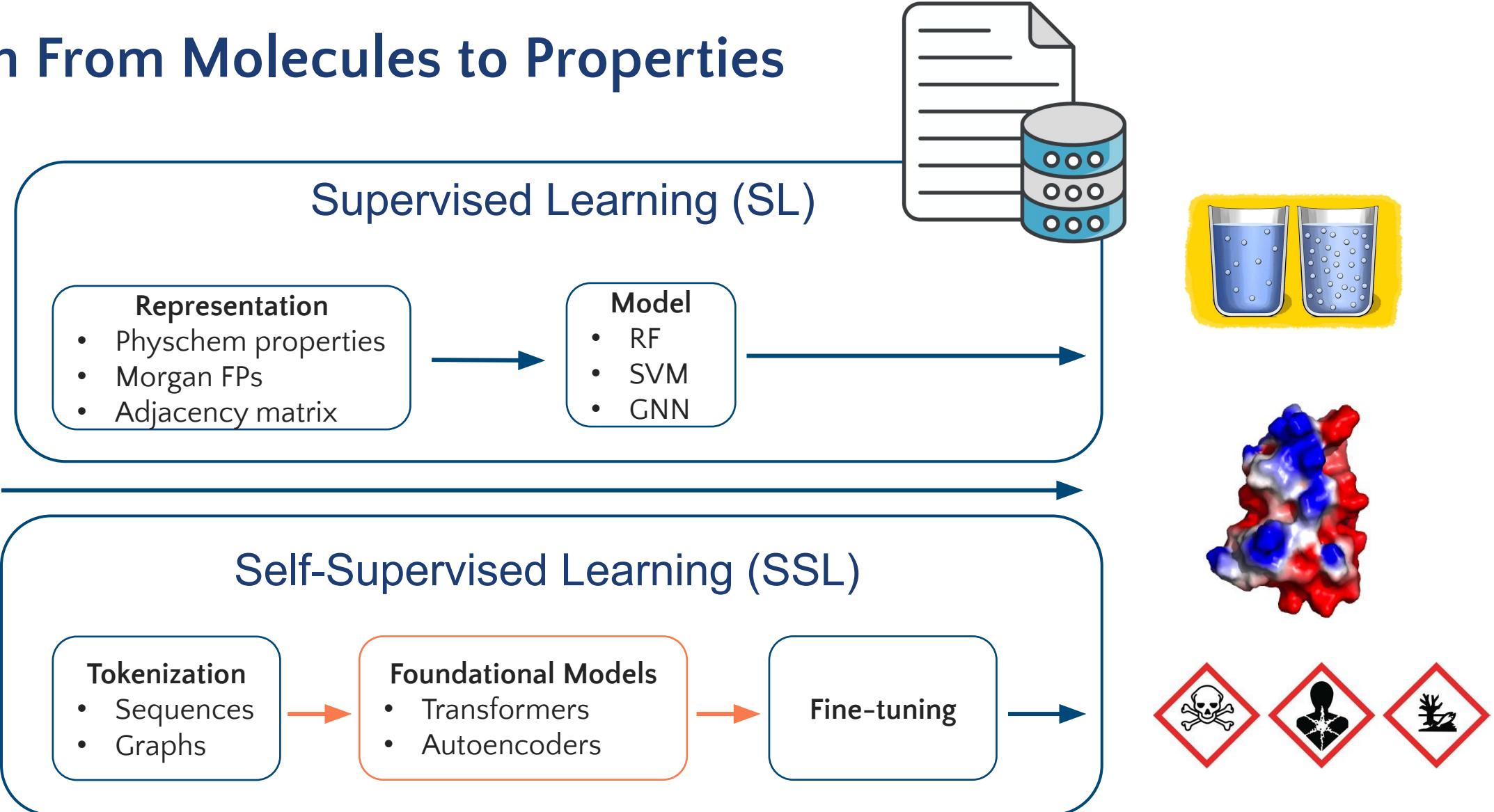
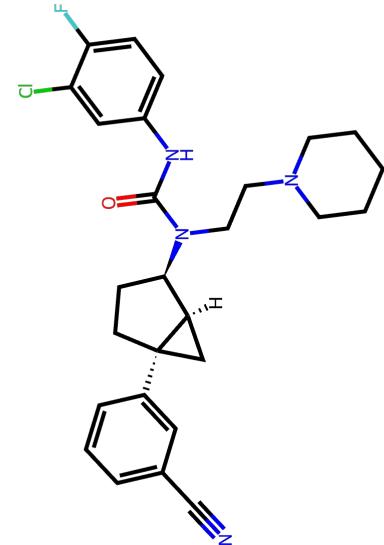
## Session 3 - 4:30 - 5:30 pm

- AI in practice
- Molecule generation
- Active learning

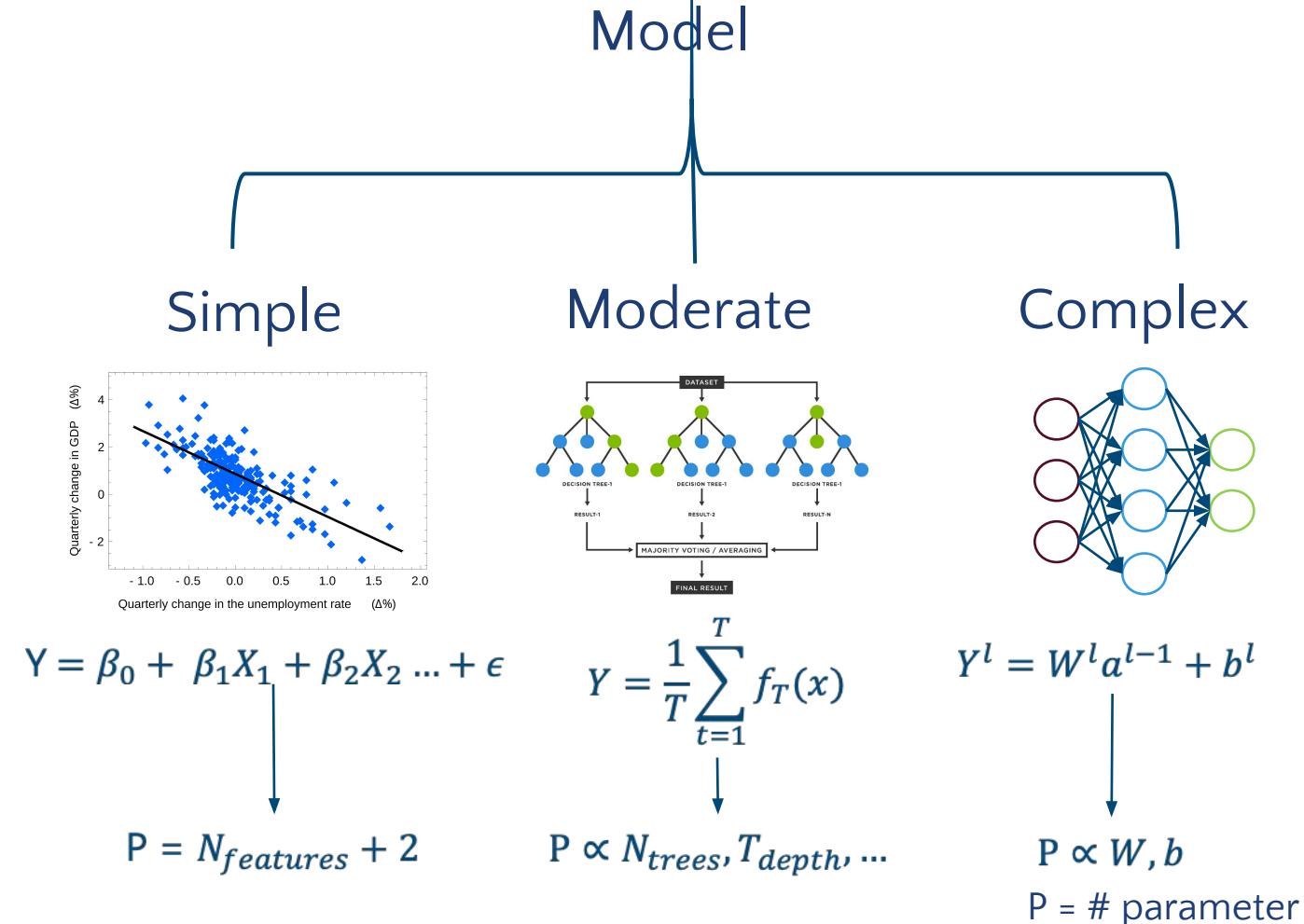
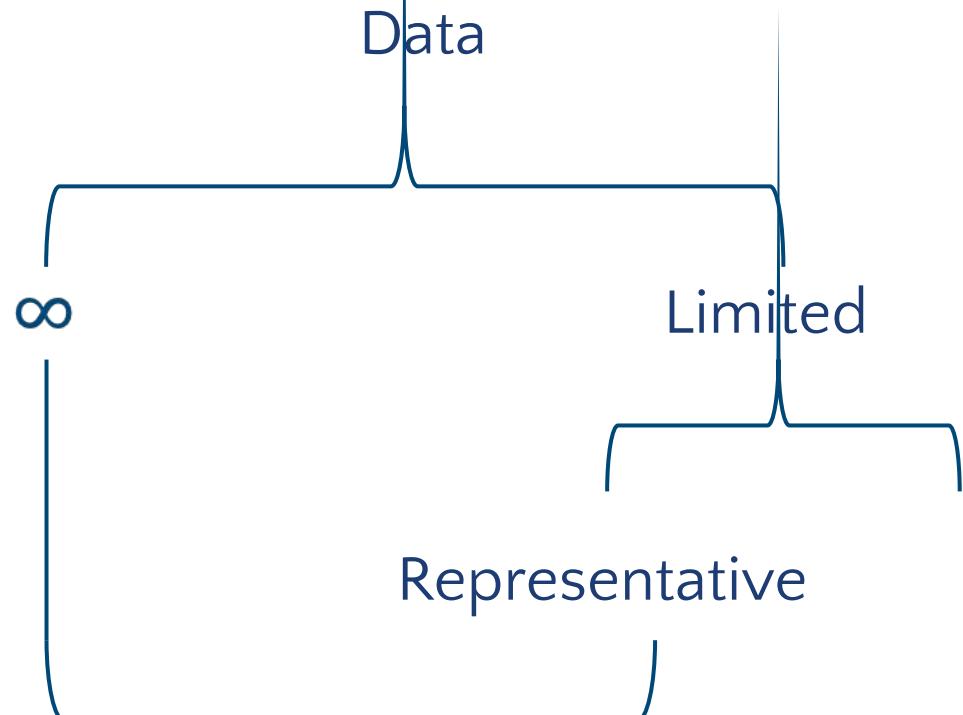
Lectures supported by hands-on sessions ...

The screenshot shows a Google Colab interface. At the top, there's a toolbar with icons for file operations (Datei, Bearbeiten, Anzeige, Einfügen, Laufzeit, Tools, Hilfe), a code editor tab (+ Code), a text editor tab (+ Text), and a link to copy the notebook to Google Drive. Below the toolbar, a message says: "This notebook is a quick test to determine if you can run a notebook on Google Colab. Click the return. At this point, you'll see a dialog that looks like this." A dark modal dialog box is displayed, containing a warning message: "Warning: This notebook was not authored by Google". It continues: "This notebook is being loaded from GitHub. It may request access to your data stored with Google, or read data and credentials from other sessions. Please review the source code before executing this notebook." At the bottom right of the dialog are "Cancel" and "Run anyway" buttons. Below the dialog, a text input field shows the command "[ ] 1+1" and a status bar at the bottom says "[ ] Beginnen Sie mit dem Programmieren oder generieren Sie Code mit KI."

# The Path From Molecules to Properties



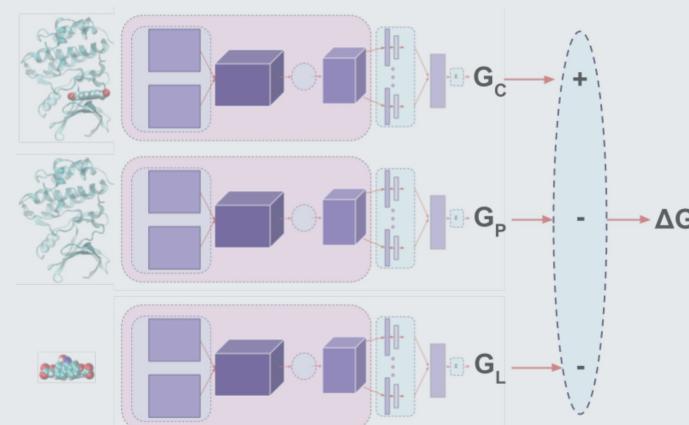
# The Two Limiting Factors



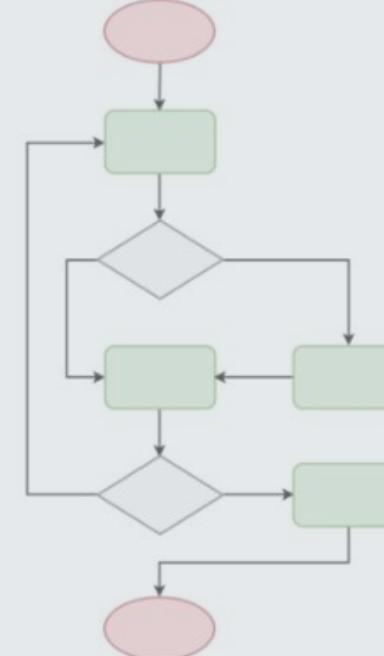
# Key Components of Machine Learning in Drug Discovery (or Anything Else)

# Data

# Representation



# Algorithms



# Where Machine Learning Excels $\longleftrightarrow$ Pharmaceutical Data Properties

## (Ultra)Large amounts of data

- Pharmaceutical data is minuscule compared to many other fields

## Responses are definitive

## Samples are independently distributed

- No cases where the same example falls into 2 different categories

## Samples are identically distributed

- Equal distributions of positive and negative examples

## Training data is representative of what is being predicted

## **Data is sparse**

- Rarely have a complete data matrix

## **Data is truncated**

- Many assay values report as “<1” or “>30”
- Difficult to know the true value

## **Data has a limited dynamic range**

- Often spans only 2 or 3 logs
- Small dynamic range combined with experimental error makes significant correlations difficult

## **Even data from the “same” assay can be heterogeneous**

- PK data measured with different doses, formulations
- Response can vary with operator, equipment, lab

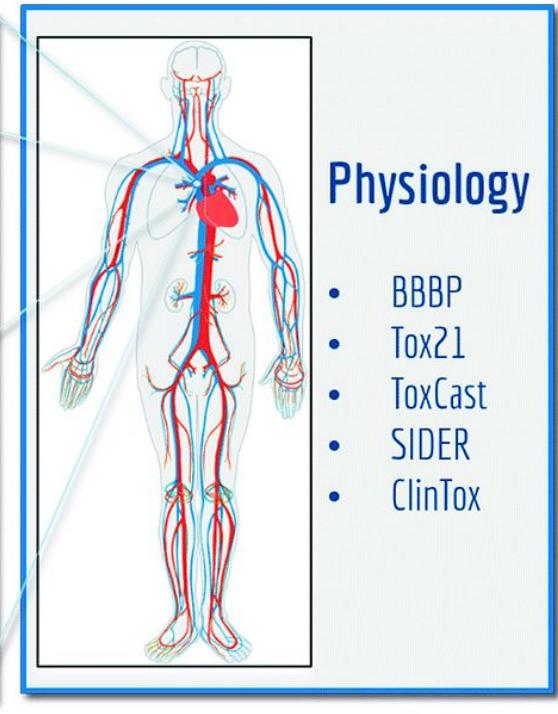
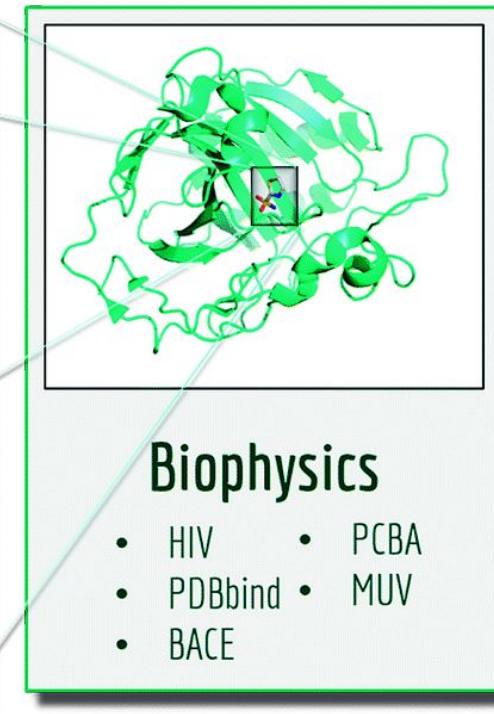
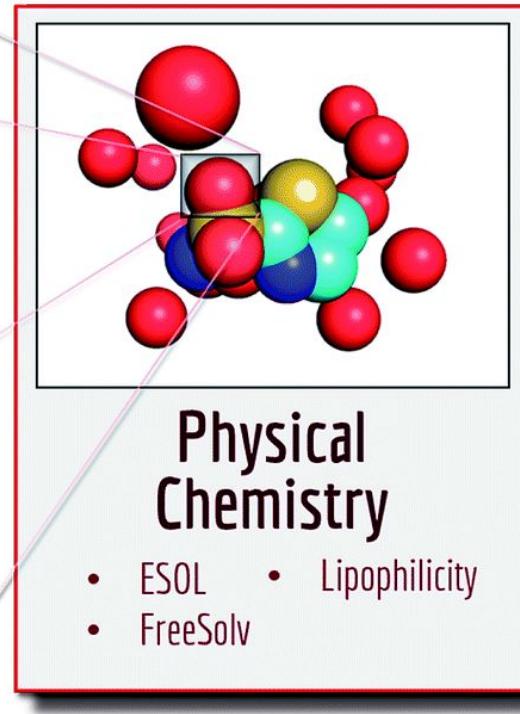
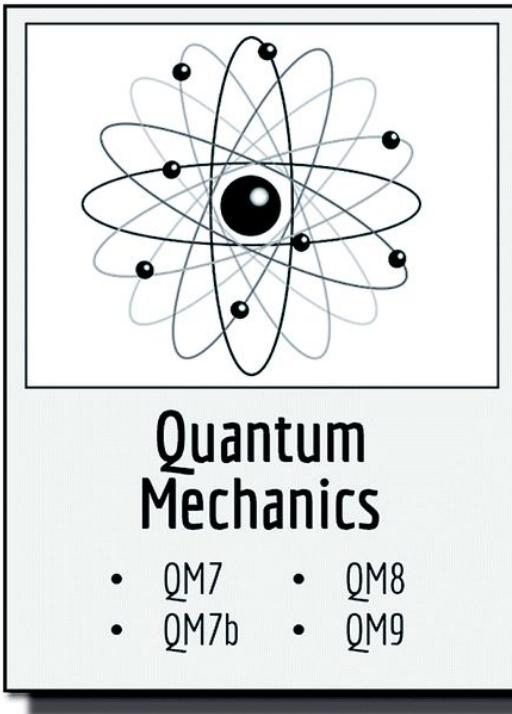
## **Data covers a limited chemical space**

- Even global models can be local

# Small Molecule Benchmark Datasets: MoleculeNet

Common molecule benchmarks for molecular property prediction:

16 datasets divided into 4 categories.



<https://moleculenet.org/>

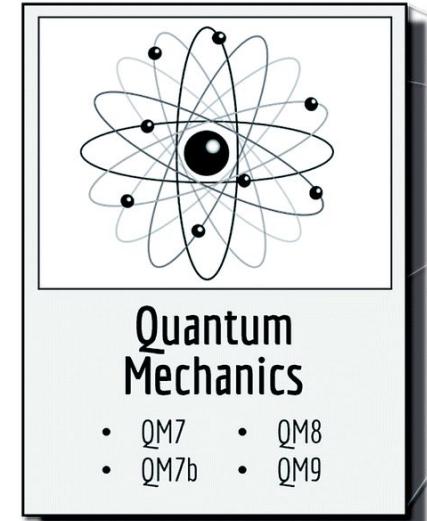
# Small Molecule Benchmark Datasets: e.g., QM9 and Tox21

## Quantum Mechanics data: QM9 → regression tasks

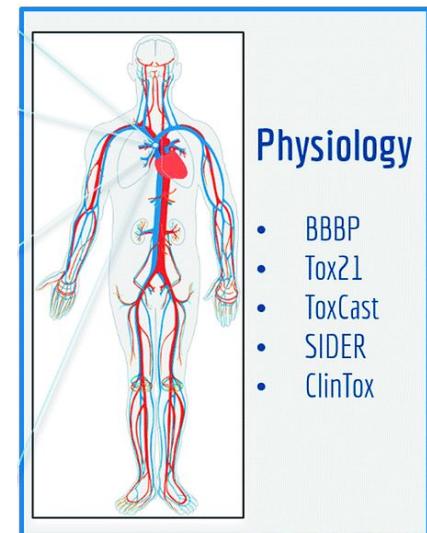
- ~134k molecules (SMILES, 3D) with <= 9 heavy atoms
- 12 properties: Geometric, energetic, electronic, and thermodynamic properties of DFT-modelled small molecules
- Commonly used due to its comparably large size

## Tox21 → classification

- Toxicity data with molecules measured against toxicity-relevant endpoints
- ~8k molecules
- 12 tasks, here biological targets, including nuclear receptors and stress response pathways



<https://moleculenet.org/>



## CAREFUL! See Pat's BlogPost

<https://practicalcheminformatics.blogspot.com/2023/08/we-need-better-benchmarks-for-machine.html>

# Benchmark Data Sets: Therapeutic Data Commons (TDC)

 Check for updates

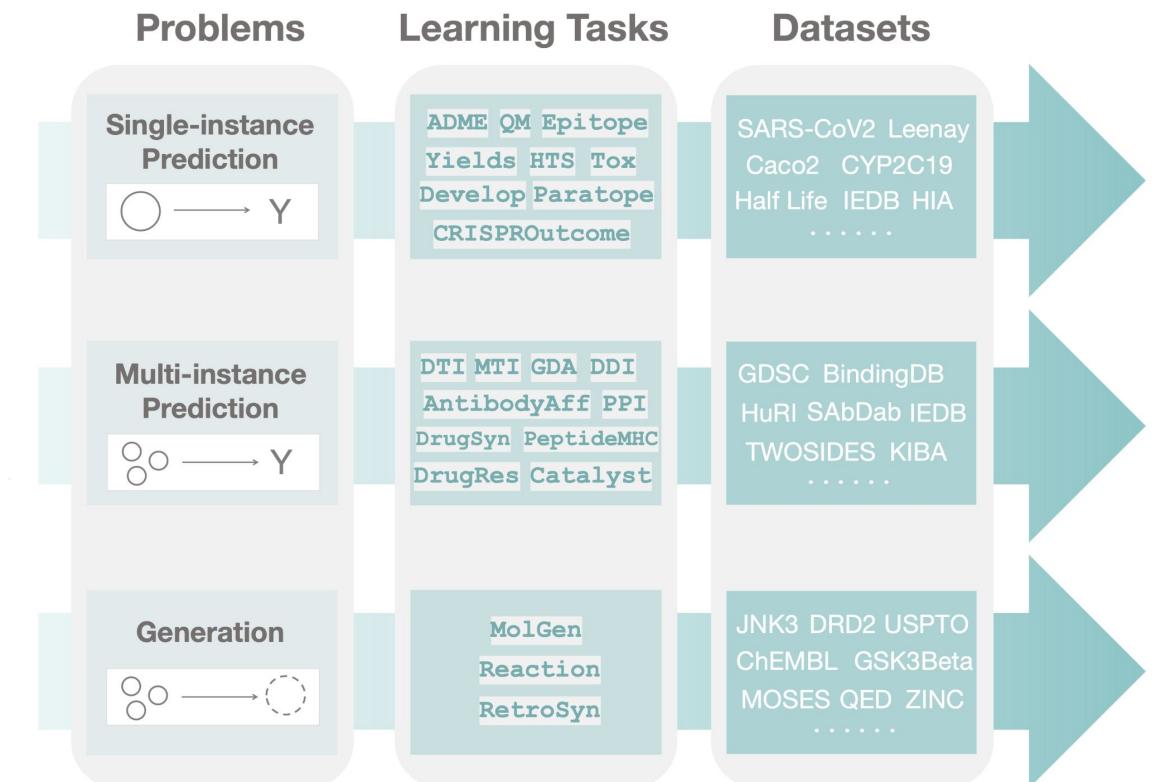
comment

## Artificial intelligence foundation for therapeutic science

Artificial intelligence (AI) is poised to transform therapeutic science. Therapeutics Data Commons is an initiative to access and evaluate AI capability across therapeutic modalities and stages of discovery, establishing a foundation for understanding which AI methods are most suitable and why.

Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W. Coley, Cao Xiao, Jimeng Sun and Marinka Zitnik

Nat Chem Biol (2022). <https://doi.org/10.1038/s41589-022-01131-2>

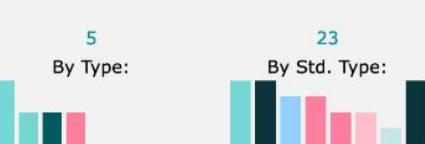


## CAREFUL! See Pat's BlogPost

<https://practicalcheminformatics.blogspot.com/2023/08/we-need-better-benchmarks-for-machine.html>

# Small Molecule Bioactivity Datasets: ChEMBL or PubChem

Database for collecting binding affinities

ChEMBL ID	Search Hit	Name	Synonyms	Type	Max Phase	Molecular Weight	Targets	Bioactivities
CHEMBL4558324		LAZERTINIB	C-18112003-G, GNS1480, GNS-1480, JNJ-73841937-AAA, Lazertinib, Yh25448, YH25448, YH-25448	Small molecule	3	554.66		
CHEMBL4650319		MOBOCERTINIB	AP32788, AP-32788, Mobocertinib, Tak-788, TAK-788	Small molecule	4	585.71		
CHEMBL4761468		AUMOLERTINIB	Almonertinib, Ameile, Amerol, Aumolertinib, Egfr t790m inhibitor hs-10296, EQ143, EQ-143, HS-10206, Hs 10296, Hs-10296, HS-10296	Small molecule	3	525.66		



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

# Small Molecule Datasets (unlabelled): ZINC and GuacaMol

## ZINC (<https://pubs.acs.org/doi/full/10.1021/acs.jcim.5b00559>)

- Ultralarge collection of purchasable, “drug-like” molecules without target value
- ZINC15: ~120M molecules
- ZINC20: ~1.4 billion molecules,

## GuacaMol (<https://pubs.acs.org/doi/10.1021/acs.jcim.8b00839>)

- Large collection of molecules from ChEMBL 24 without target value
- ~1.5M molecules
- Published training, test and validation set → for benchmarking

→ Commonly used benchmark for pretraining (LLMs) and generative models

# (Benchmark) Data Quality - Technical Issues

Inspired by Pat's BlogPost  
[we-need-better-benchmarks-for-machine-learning/](https://patcc.org/we-need-better-benchmarks-for-machine-learning/)

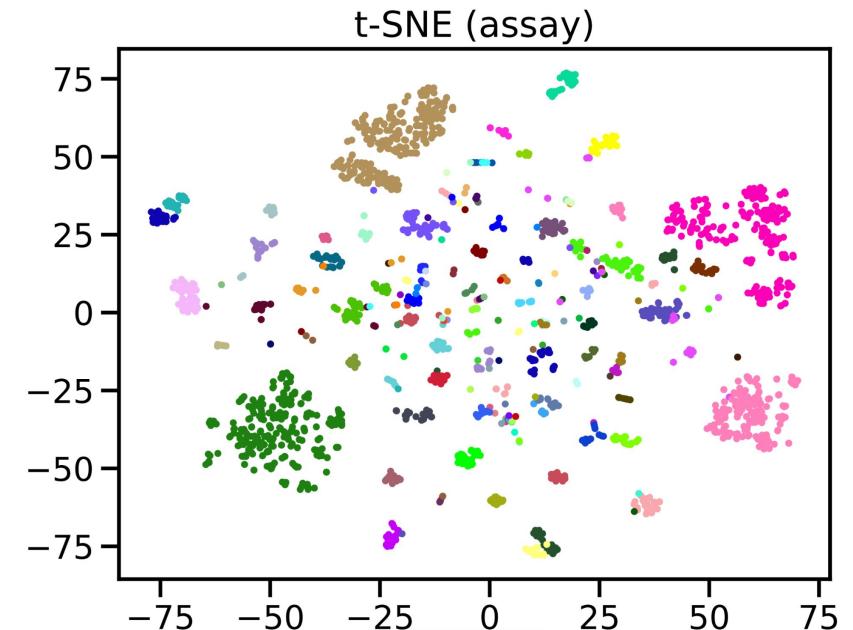
- **Valid structures**
  - Invalid SMILES, not be parsed by standard cheminformatics tools (e.g. rdkit)
- **Consistent chemical representations**
  - Standardized based on accepted conventions
- **Stereochemistry**
  - We need to know the correct structure
- **Consistent measurements**
  - Ideally experiments performed in same lab, but at least same standards
- **Realistic dynamic range and cutoffs**
  - Typical range: solubilities between 1 and 500  $\mu\text{M}$  (cut at both ends ‘<’ an ‘>’)
  - Sensible cut-off choices for classification
- **Clear definitions of training, validation, and test sets**
  - Data leakage: random, scaffold and cluster splits
  - IDs need to be given in benchmark
- **Data curation errors**
  - Duplicate structures with differing labels

## *Not so good examples from MoleculeNet*

- 11 SMILES with uncharged tetravalent nitrogen atom in BBB
- 59 beta-lactam antibiotics in BBB, carboxylic acid in three different forms
- 71% of the molecules in BACE have at least one undefined stereocenter
- BACE data collected from 55 papers
- ESOL spans 13 logs, easy to separate

# Typical Challenges in (Public) Chemical Datasets

- Different experimental protocols
- Inconsistent values
- Missing documentation
- Data tends to be heavily clustered



- kinodata-EGFR ligand activities
- t-SNE on 2048 bit Morgan fingerprints
- colors = assays

# Applying ML to Drug Design with Trust



Do I have the right data?



Can I trust this prediction?



Is the model interpretable?

(Skipped due to time constraints)

- Model validation and data splits
- Model generalizability and applicability

# Is My Training Data Relevant?

Machine learning is all about labeling things using examples



If I train on this?

Can I predict this?

# Example Scenarios



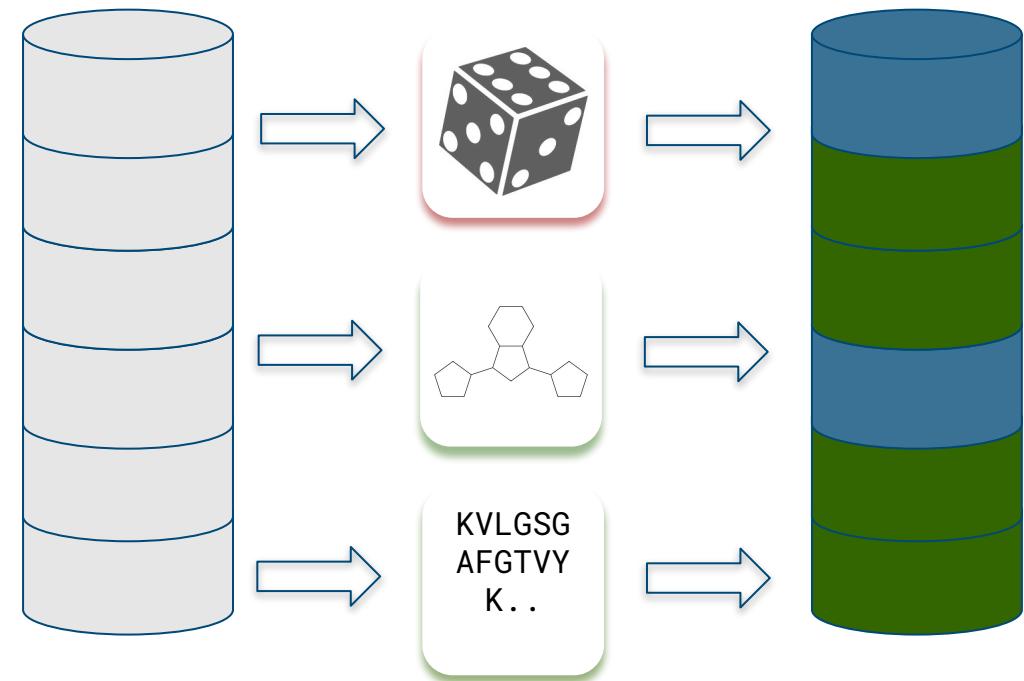
Can I predict properties for my molecules based on literature data?



Can I make predictions for a new project from my existing data?

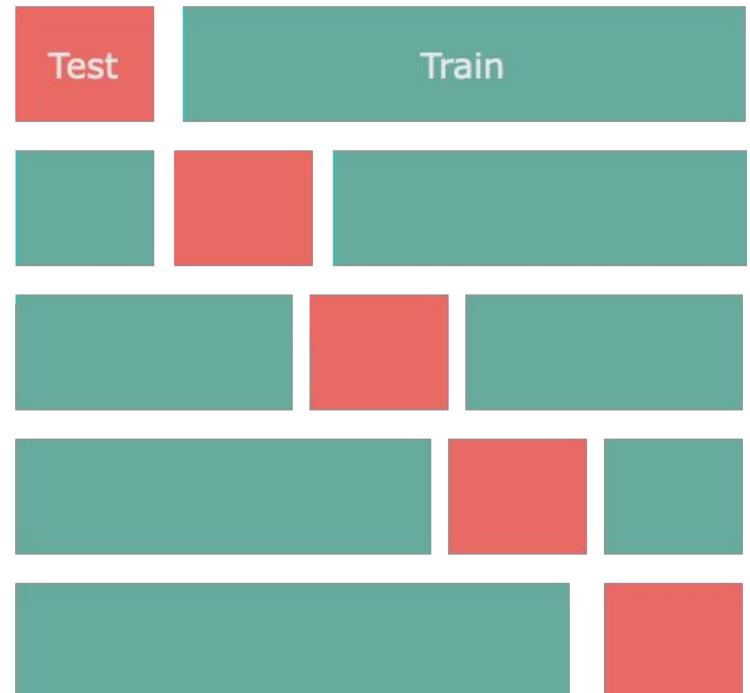
# Data Split is Crucial to Assess Generalizability

- Train/val/test split
- Random splits overestimate out-of-domain generalization
- Better: Splits based on
  - Ligand scaffold
  - Protein sequence

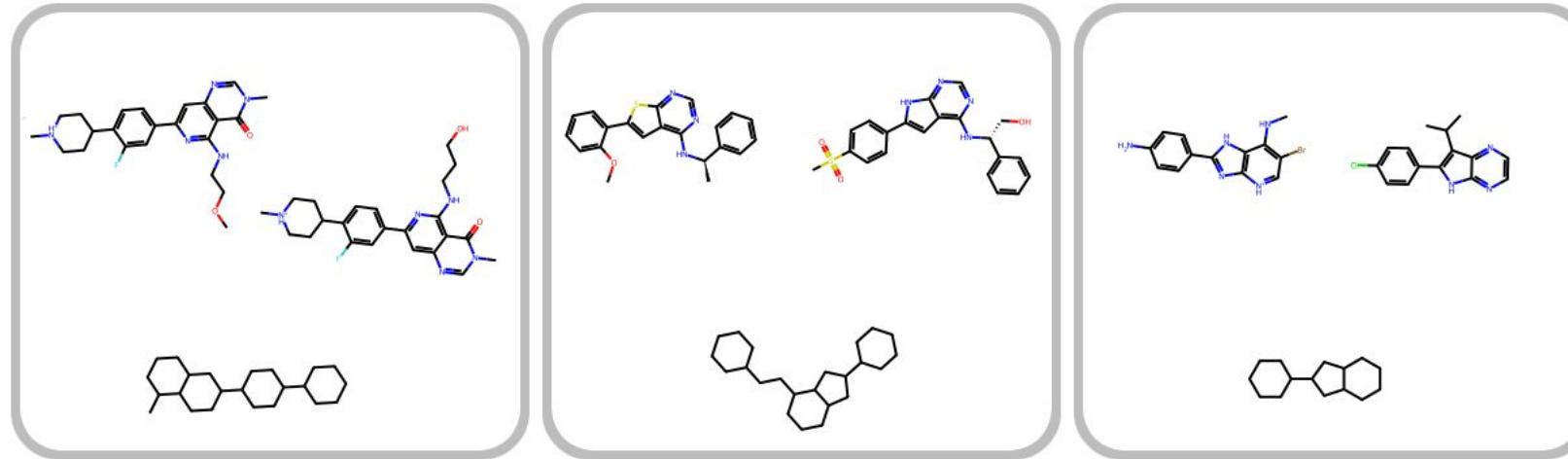


# Grouped k-Fold Split

- Folds are in accordance with groups (scaffold clusters)
- Every data point is part of the test set *exactly once*
- Folds may not be of exact equal size
- Consider labels (stratified split)



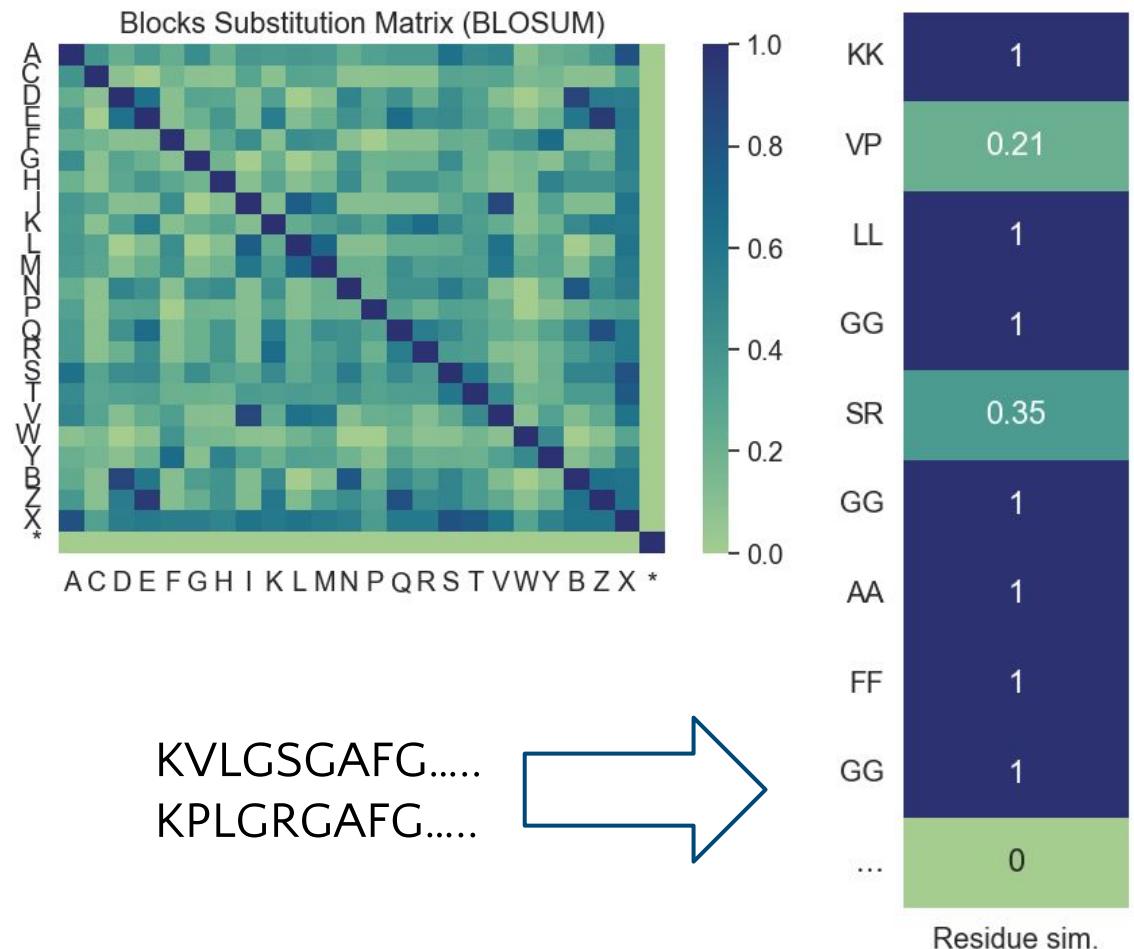
# Data Splits – Scaffold Split



- Avoid data leakage over structurally similar ligands
- Group by generic scaffold (*or similarity*)
- Ensure each group only in one of {train, test, val}

# Data Splits - Target-Based Split

- Cluster pocket sequences based on BLOSUM-similarity
- Generate splits by randomly assigning clusters to {train, val, test}



# DataSAIL - Data Splitting Against Information Leakage

Developed by Roman Joeres at Prof. Kalinina's lab, Saarland University  
uses mathematical optimization to identify the most difficult split

<https://github.com/kalininalab/DataSAIL>



# Applicability and Uncertainty of Machine Learning Models



Do I have the right data?



Can I trust this prediction?



Is the model interpretable?

Conformal Prediction

Bootstrap

Gaussian Process Regression

Jackknife

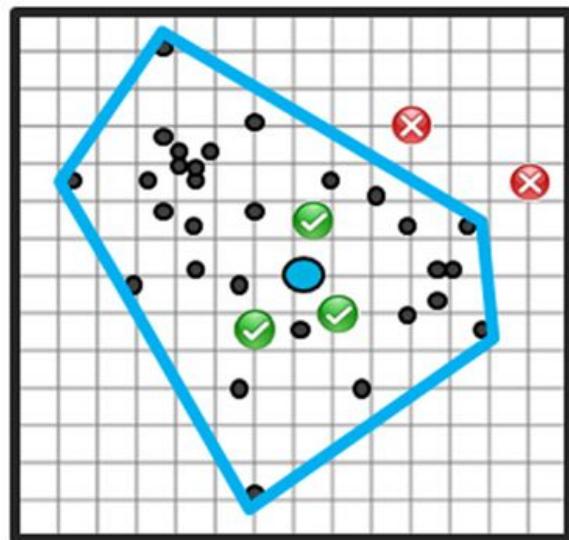
Monte Carlo Dropout

# Applicability Domain

**Definition:** The applicability domain of a (Q)SAR/ML model is the response and chemical structure space in which the model makes predictions with a given reliability.

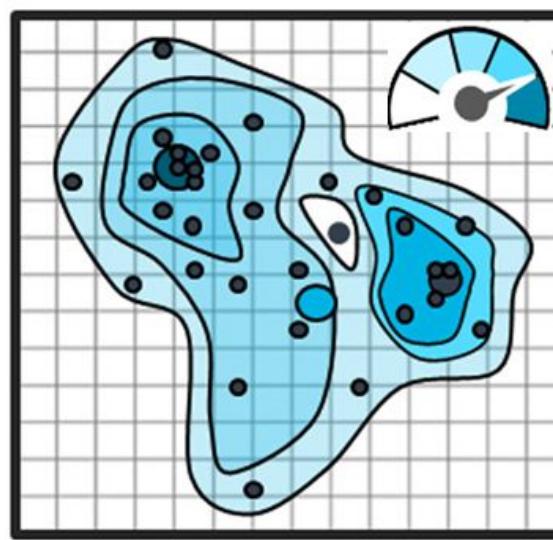
## 1. Applicability:

Interpolation within the chemical structure space



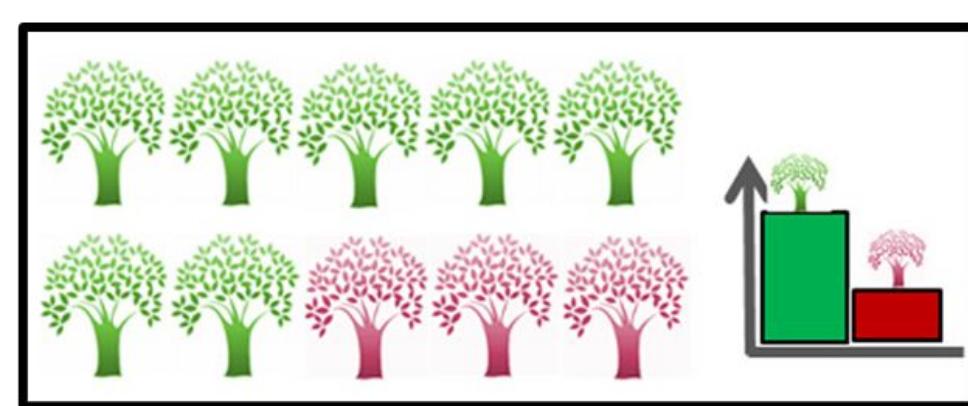
## 2. Reliability:

Density of knowledge around the query chemical structure



## 3. Decidability:

Distance of the query chemical compound to the decision boundary of model



Figures taken from: Hanser, et al., SAR and QSAR in Environmental research, **2016**, 27:11, 865-881

# Using Similarity to Assess Applicability



0.84



0.63

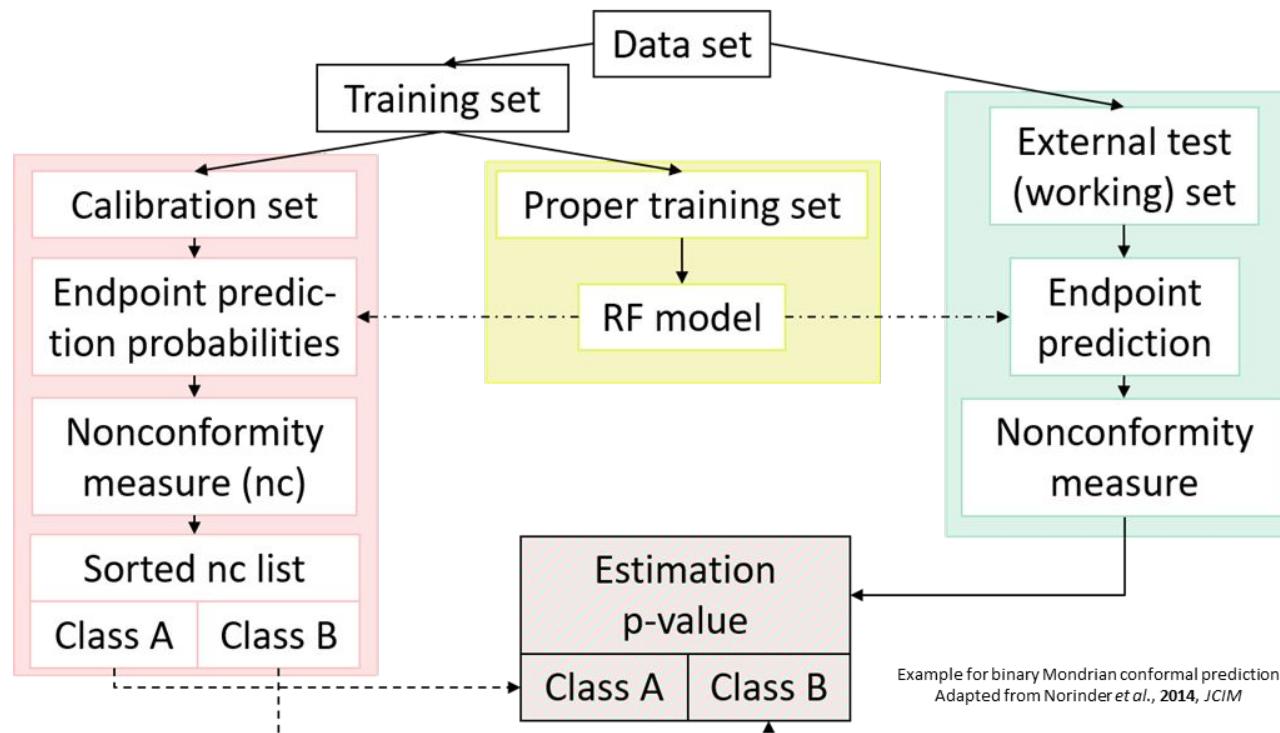


0.15



# Conformal Prediction Framework Based on ML

- Statistically valid on a given confidence level (if exchangeability is given)
- Additional calibration step: Compare predictions to those previously seen (p-value)
- Output prediction sets in binary classification: {A}, {B}, {A,B} or {}
- Validity = % of correct classifications
- Efficiency = % of single class predictions



		<chem>O=[N+]([O-])c1ccc(cc1)[C@H](C)NC(=O)C(C)C</chem>	prob(A) = 0.7	prob(B) = 0.3		
			Class A	Class B	0.99	
0.83					0.70	
0.67					0.63	
0.47	2/5				0.37	
0.43					0.33	
0.41	0/5				0.31	
					0.12	
			p-value(A) = 0.8	p-value(B) = 0.14		Significance level 0.2: {A}

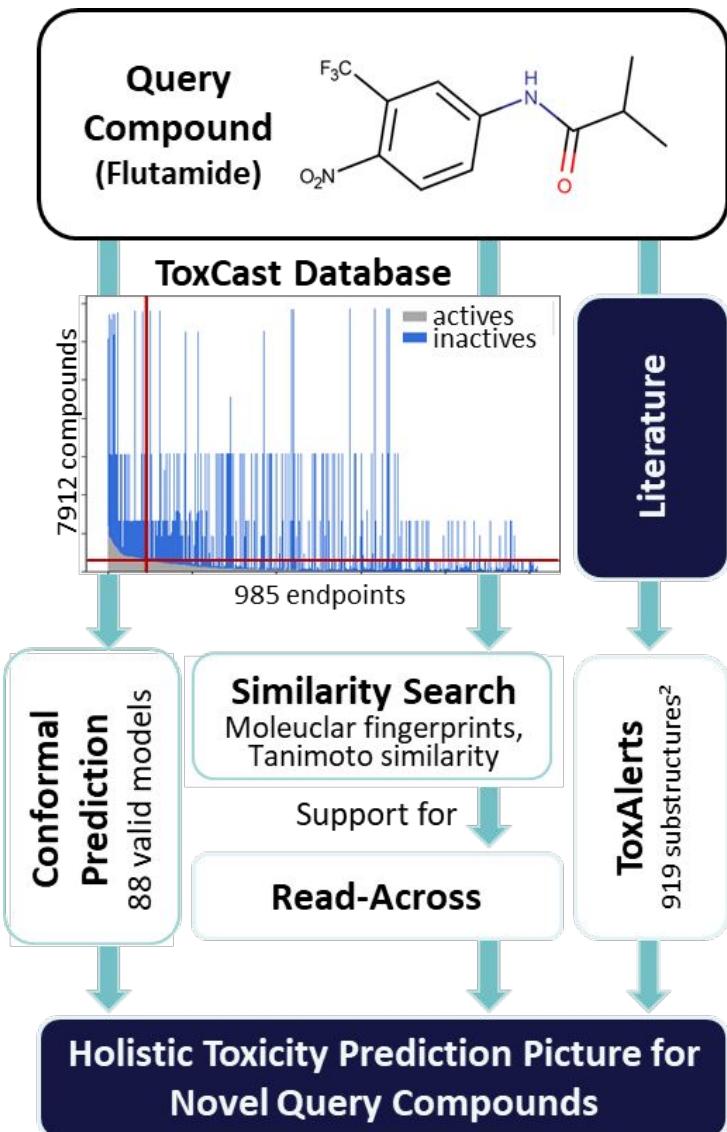
# Know-Tox - Predictions with Confidence

- ToxCast data set
  - ~ 8000 compounds \* 1000 endpoints (sparse)
  - Pharmaceutical, pesticides, environmental chemicals
  - Cell cycle, steroid receptors, cytotoxicity, ...

Can we apply the model to new (BASF) data?

- Conformal prediction
  - Case study: Antiandrogen activity (AA)
  - Validity, efficiency and accuracy like other studies

	Dataset	Efficiency	Accuracy (SCPs)			# toxic/non-toxic
			all	cl.1	cl.0	
cross-validation	ToxCast AA	0.87	0.78	0.80	0.78	868/5842
external	Norinder <sup>3</sup>	0.79	0.68	0.70	0.67	160/201
In-house	BASF (I)	0.94	<b>0.56</b>	<b>0.97</b>	<b>0.07</b>	280/254



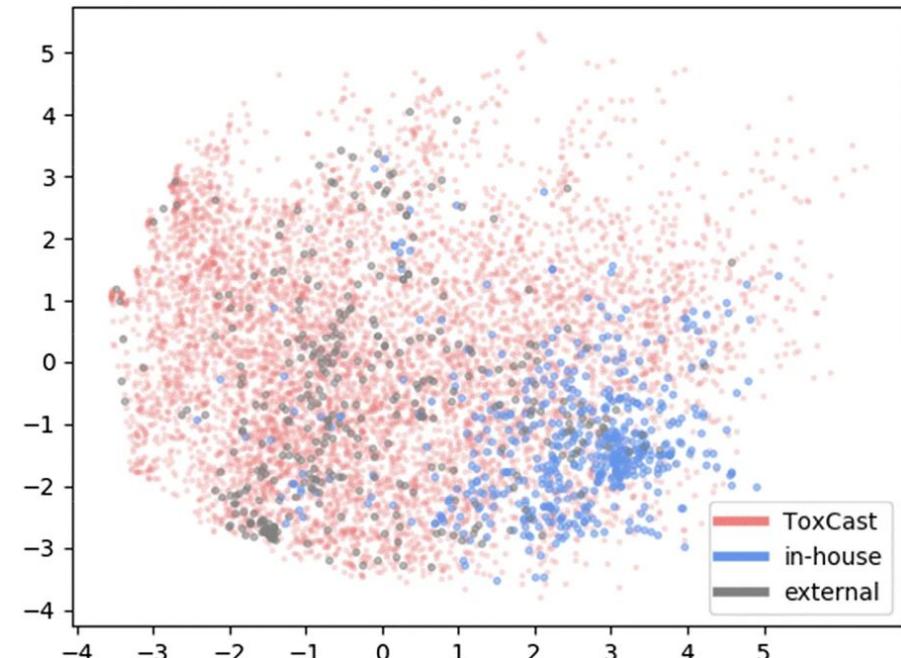
# Know-Tox – Predictions with Confidence

- ToxCast data set
  - ~ 8000 compounds \* 1000 endpoints (sparse)
  - Pharmaceutical, pesticides, environmental chemicals
  - Cell cycle, steroid receptors, cytotoxicity, ...

Can we apply the model to new (BASF) data?

- Conformal prediction
  - Case study: Antiandrogen activity (AA)
  - Validity, efficiency and accuracy like other studies
  - kNN normalization & balancing
  - Less efficient, but higher accuracy

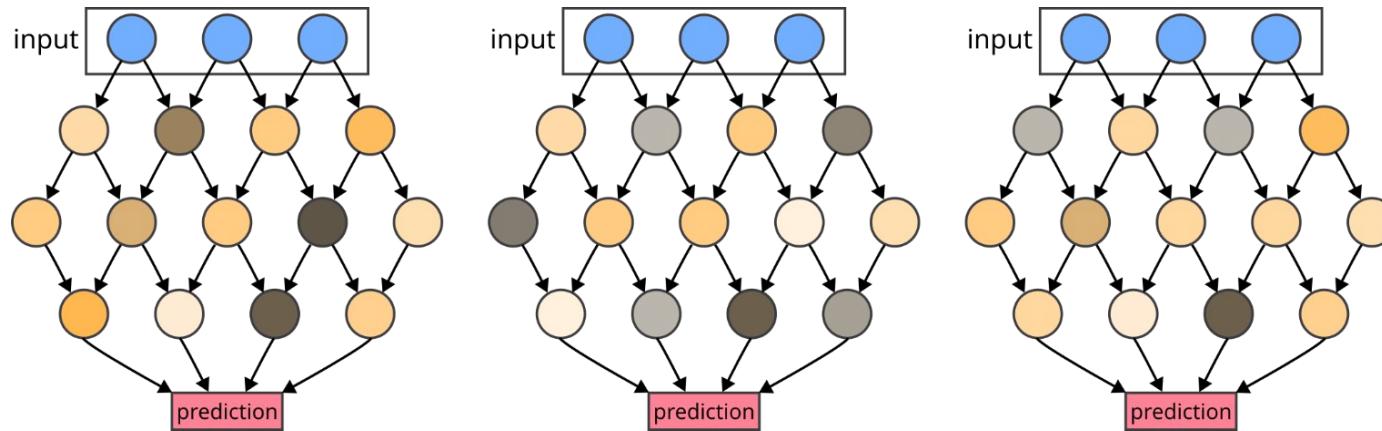
	<b>Dataset</b>	<b>Efficiency</b>	<b>all</b>	<b>Accuracy (SCPs)</b>	<b># toxic/non-toxic</b>	
				cl.1	cl.0	
cross-validation	ToxCast AA	0.87	0.78	0.80	0.78	
	Norinder <sup>3</sup>	0.79	0.68	0.70	0.67	
	BASF (I)	0.94	<b>0.56</b>	<b>0.97</b>	<b>0.07</b>	
In-house	BASF (II)	<b>0.20</b>	<b>0.75</b>	<b>0.80</b>	<b>0.71</b>	280/254



Descriptor space of a 2-component PCA trained on ToxCast-AA data

# Using Uncertainty to Assess Model Confidence

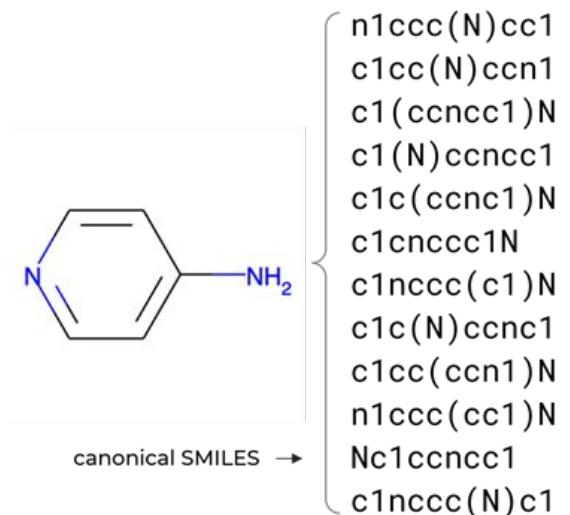
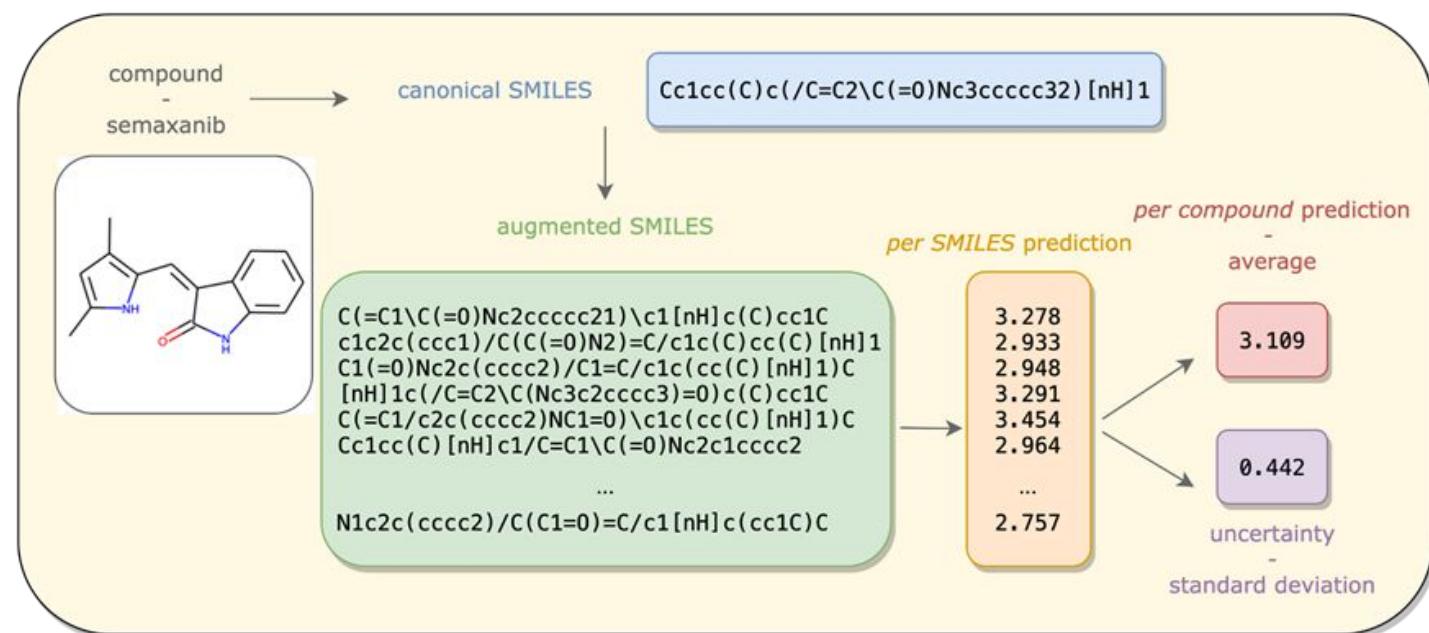
- Single deterministic methods
  - Classification setting: predicted class probabilities
  - Secondary model trained to predict uncertainty for an already trained model
- Ensemble methods
  - Ensemble of similar but different models -> Variance as uncertainty
  - Varying random seeds, resulting in different weights
  - Varying the training data (*bagging*)
- Training- and test-time data augmentation
  - Augmented set of points using stochastic noise or domain knowledge



<https://github.com/volkamerlab/teachopencadd>

# SMILES Augmentation Improves Model Performance

Maximizing molecular property prediction performance with confidence estimation using SMILES augmentation and DL

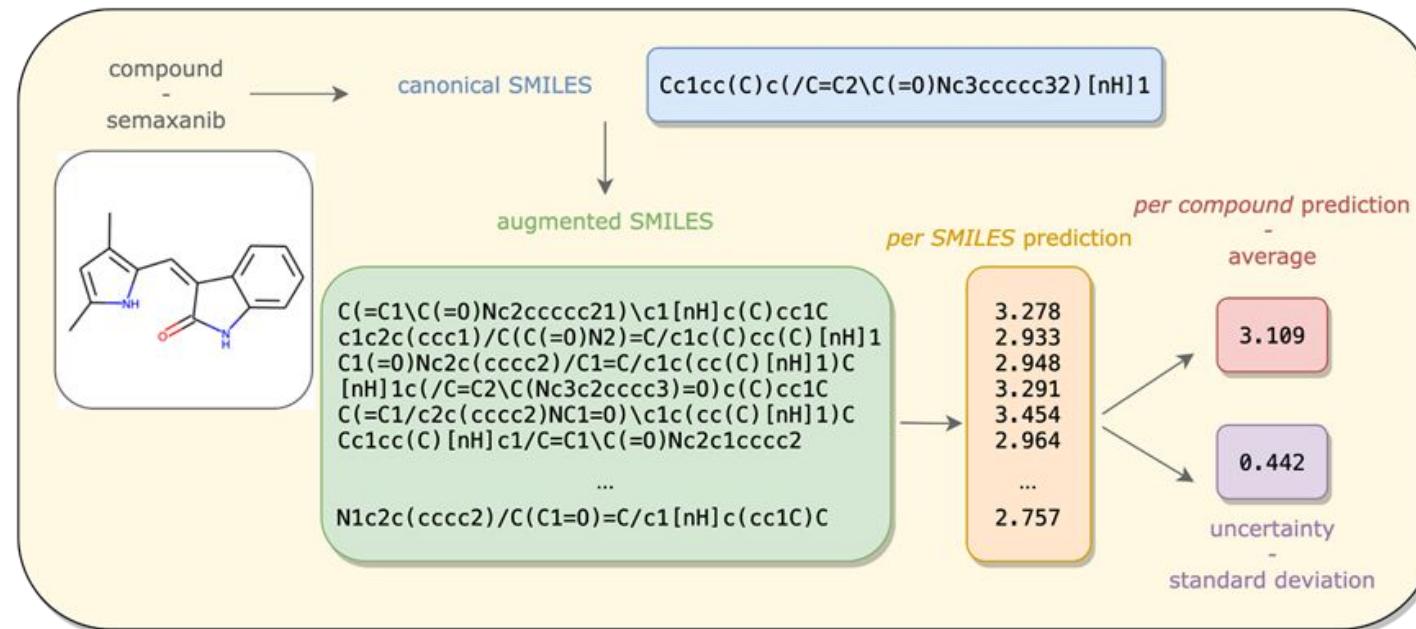


Kimber et al., AI in Life Sciences, 2021, 1, 100014

<https://github.com/volkamerlab/maxsmi>

# SMILES Augmentation Improves Model Performance

Maximizing molecular property prediction performance with confidence estimation using SMILES augmentation and DL

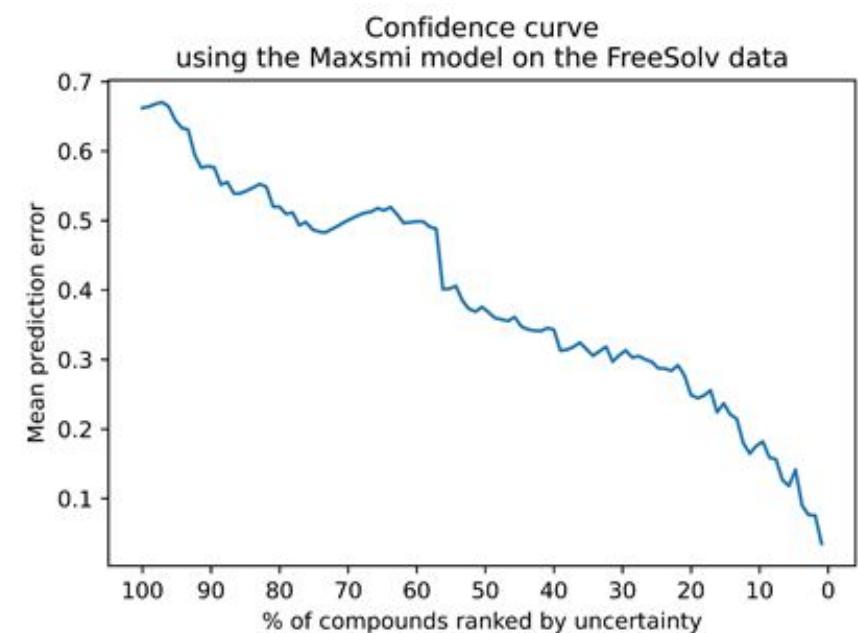


Kimber et al., AI in Life Sciences, 2021, 1, 100014

<https://github.com/volkamerlab/maxsmi>

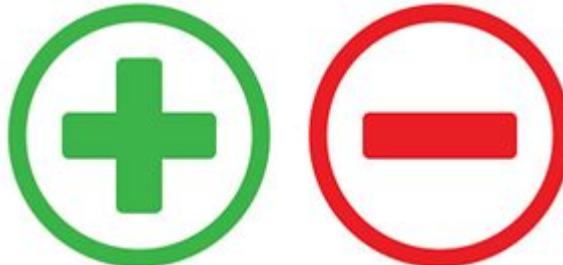
## Uncertainty prediction

More confident model implies smaller prediction error



# Take Home Messages

- More data available
- Increased computational power
- Advances in neural network algorithms
- Open-source libraries



- True predictive power
- Applicability and reliability
- Interpretability

- **Data quality, relevance and proper processing are key**
- Proper model evaluation, always consider baselines when applying ML/DL
- Check applicability of your model, ideally prospective studies
- Open source and reproducibility

# TeachOpenCADD in a Nutshell

<https://github.com/volkamerlab/teachopencadd>

# What? Pipelines for common CADD\* tasks using open resources

# How? Coding-based Jupyter notebooks<sup>1</sup> (Python) & GUI-based KNIME workflows<sup>2</sup>



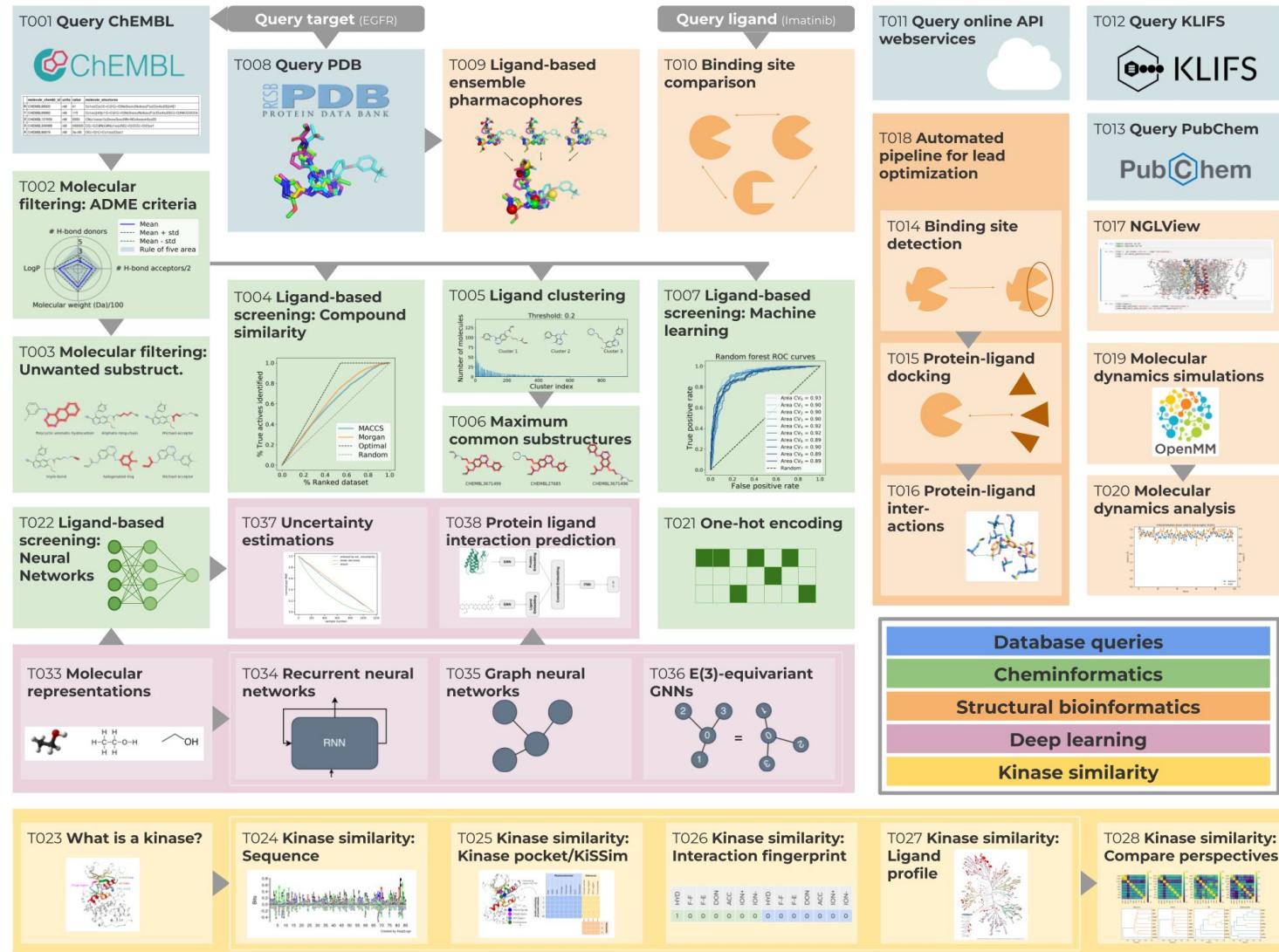
**For what?** Teaching, self studies & starting point for research projects

**From and for whom?** Students, group, community → Beginners & advanced users

\*CADD = computer-aided drug design



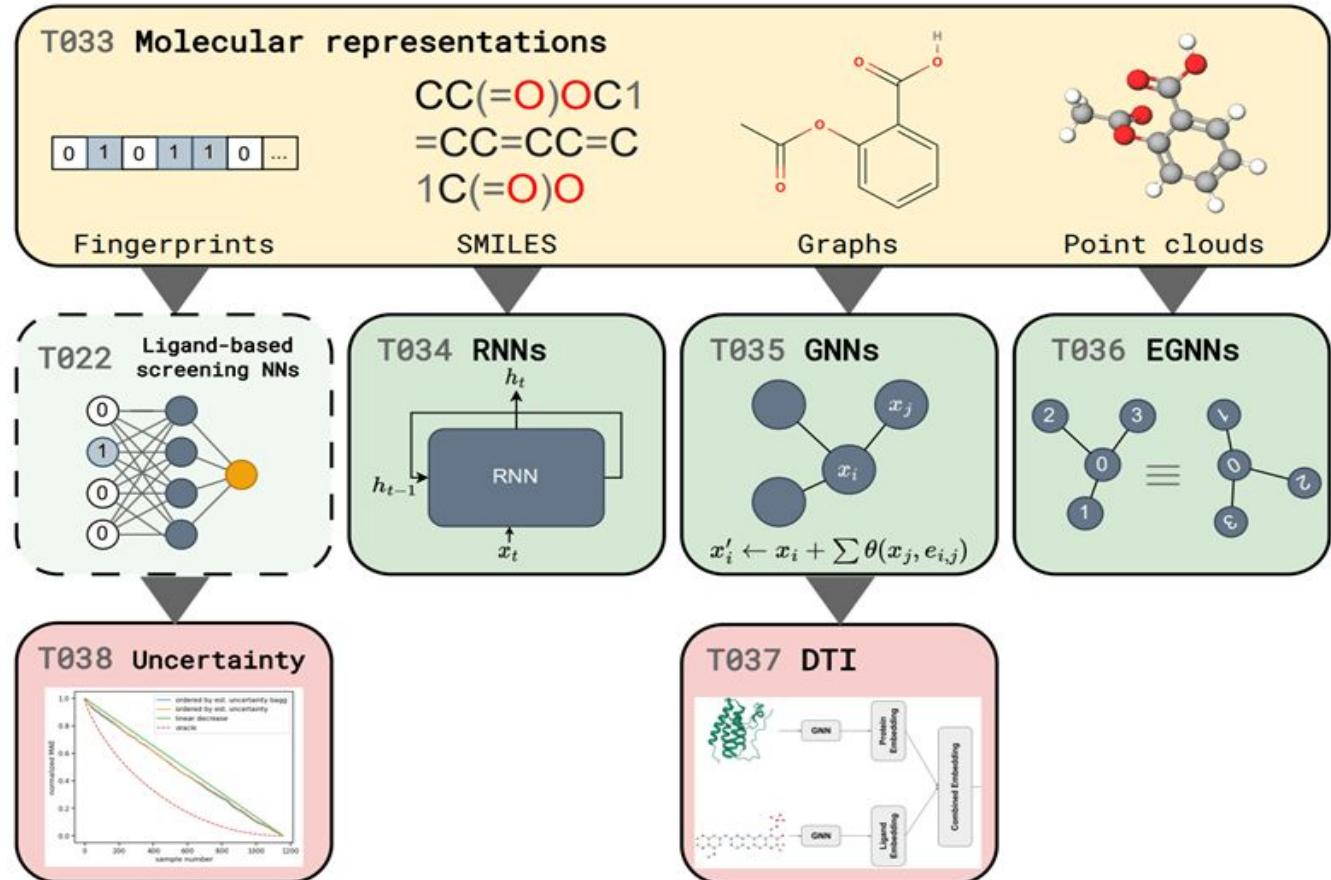
# TeachOpenCADD – Open-Source Python Pipelines



# TeachOpenCADD goes Deep Learning (DL)

Introduction to molecular representations and architectures

- Molecular representation
  - Fingerprints
  - Text (SMILES)
  - Molecular graphs
  - Point clouds
  - Embeddings (learned)
- Different architectures
  - Multi-layer neural network
  - Recurrent neural network (RNN)
  - Graph neural network (GNN)
  - E(3)-invariant GNN





UNIVERSITÄT  
DES  
SAARLANDES

Thank you for your attention!