

Enhancing Drug Discovery with Machine Learning: ADMET Property Modeling

Marcin Kowiel, PhD | April 06, 2024



Developing therapeutics
at the forefront of oncology

Agenda

01

About Ryvu

02

Drug discovery process

03

AI in drug discovery

04

Property prediction model training pipeline at Ryvu

05

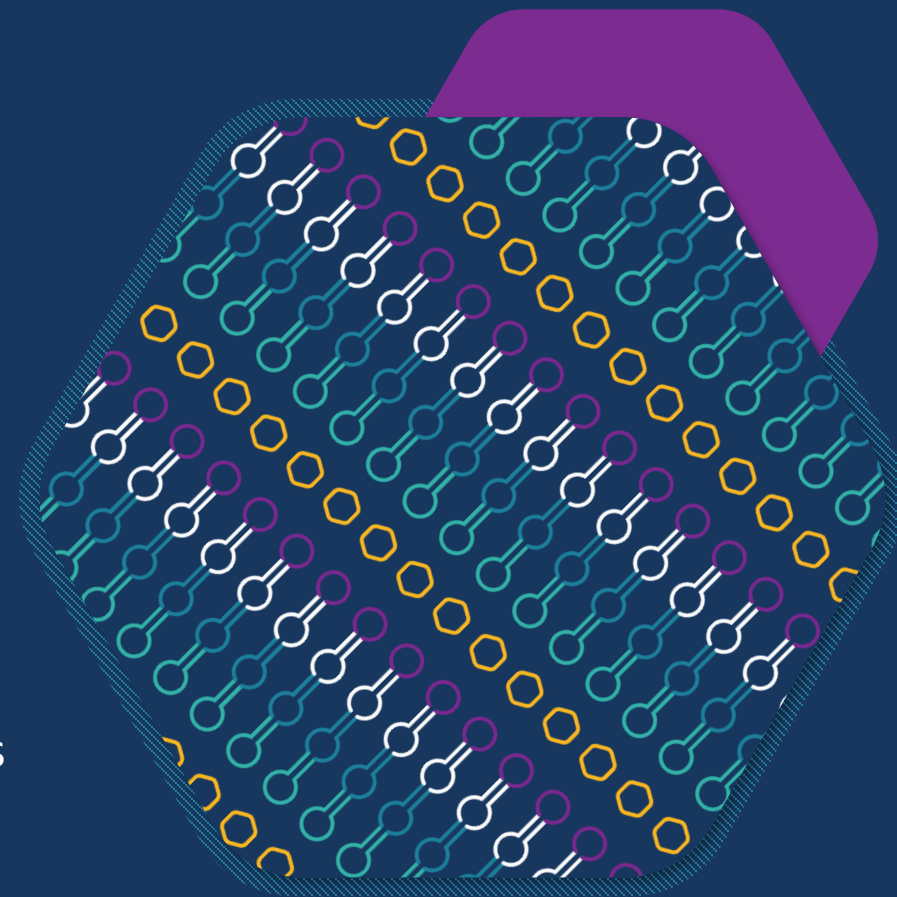
Pipeline stages

06

Summary

Takeaway points

- Ryvu Therapeutics is leading Polish drug discovery company
- AI can improve every stage of the drug discovery process
- Property prediction models can improve lead optimization process
- High quality models can be achieved through combination of data science, data engineering and understanding of the field
- Models' interpretability is important for medicinal chemists
- It is possible to automate many steps of the property prediction model training pipeline



Ryvu Therapeutics – Developing therapeutics at the forefront of oncology

Ryvu Therapeutics is a clinical stage **biopharmaceutical** company developing novel small molecule therapies addressing emerging targets in **precision oncology**.

The company was founded in **2007** and has its headquarters in **Kraków**. It was previously known as **Selvita** until it was renamed Ryvu Therapeutics after the spinning out of the services segment.

290+

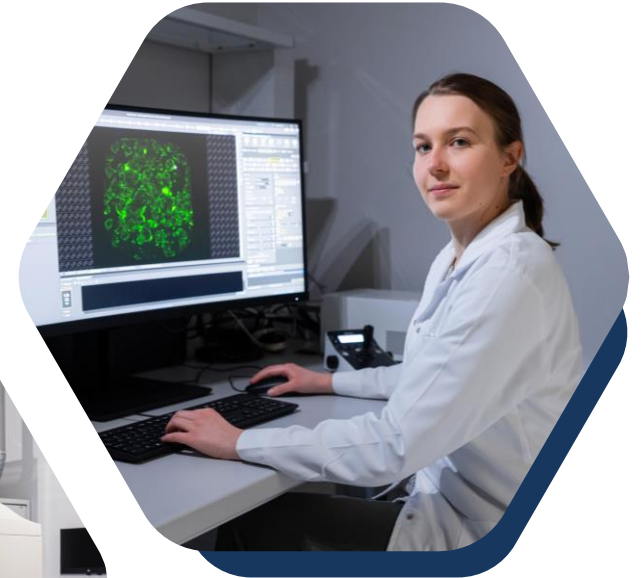
Employees

200

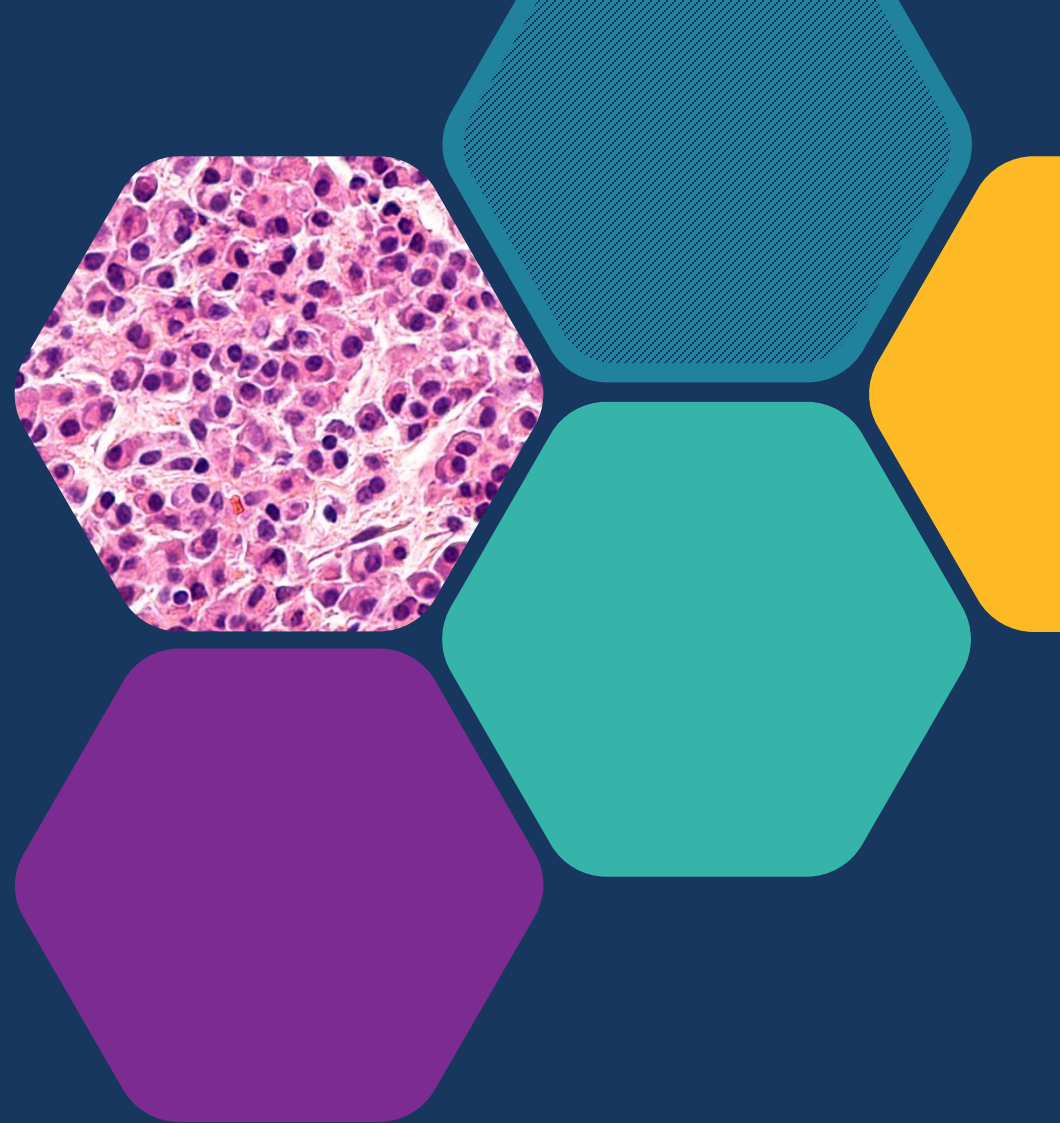
Scientists

20

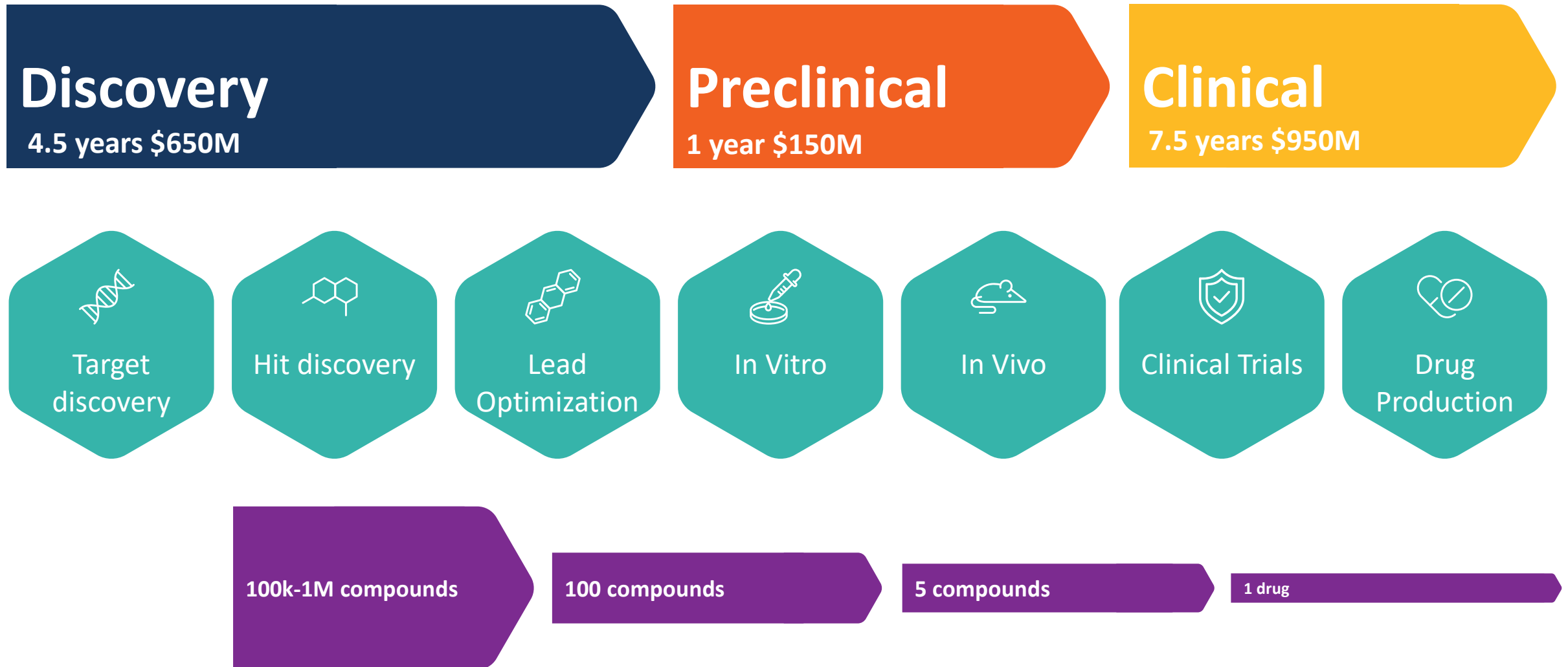
Machine Learning,
Cheminformatics and
Bioinformatics Engineers



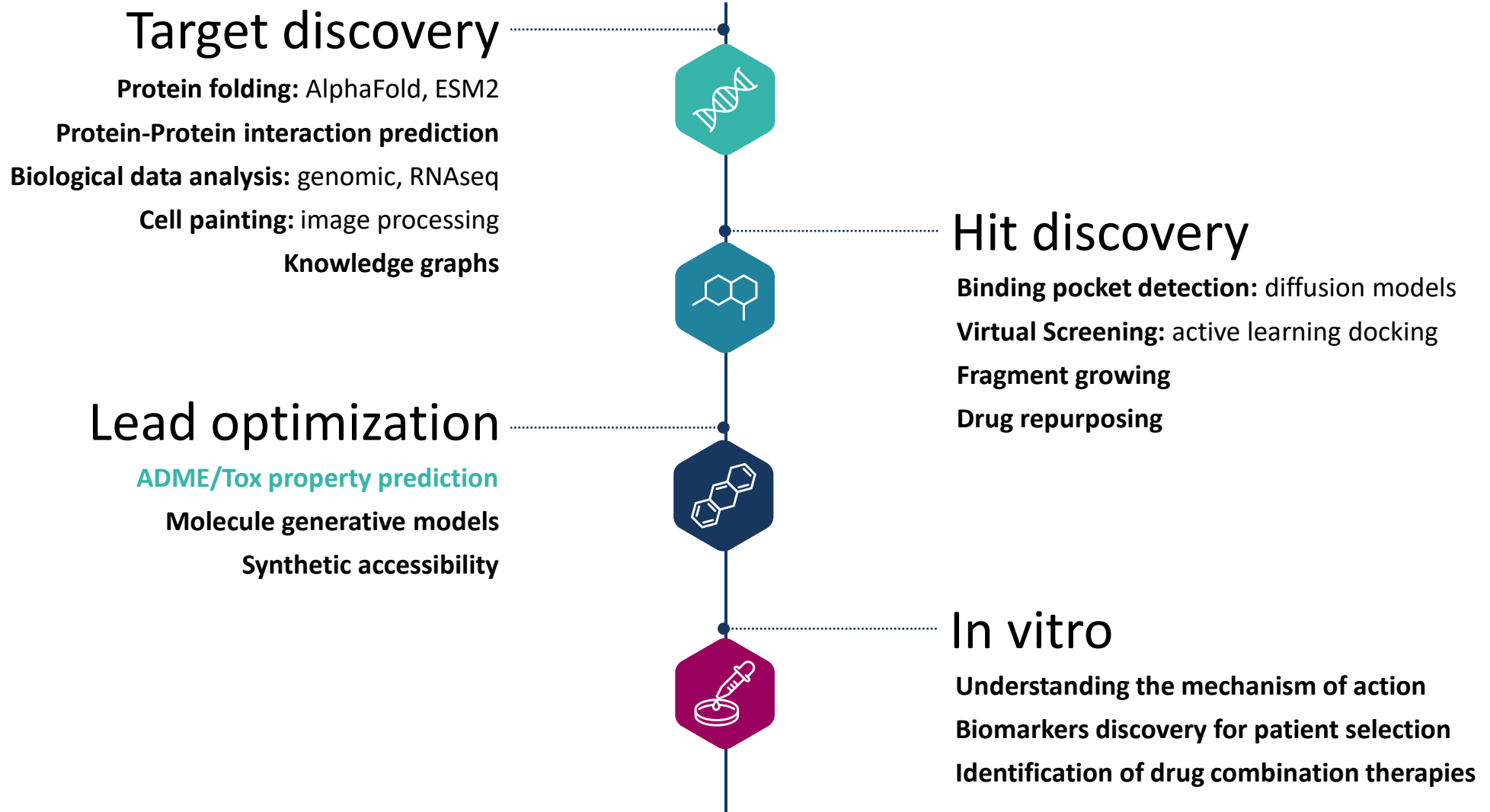
AI in drug discovery



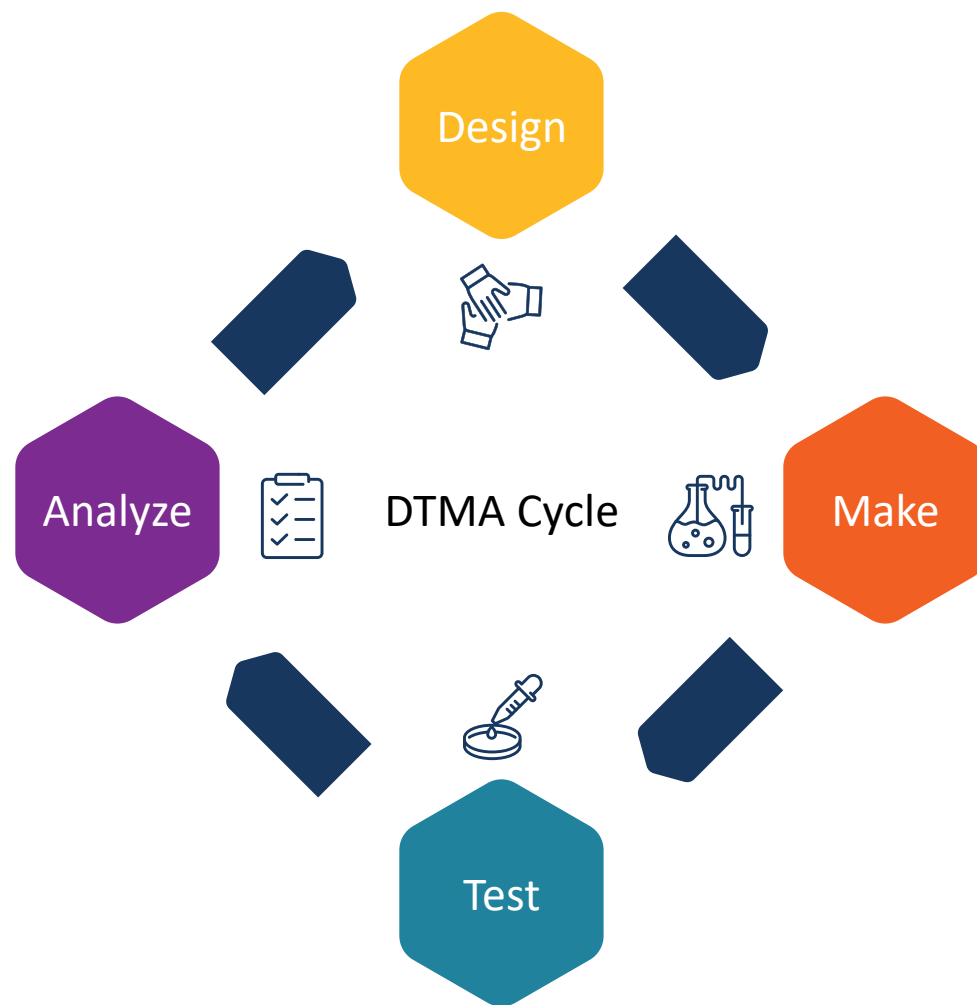
The drug discovery process is costly (\$1.7B) and time-consuming (10-15 years)



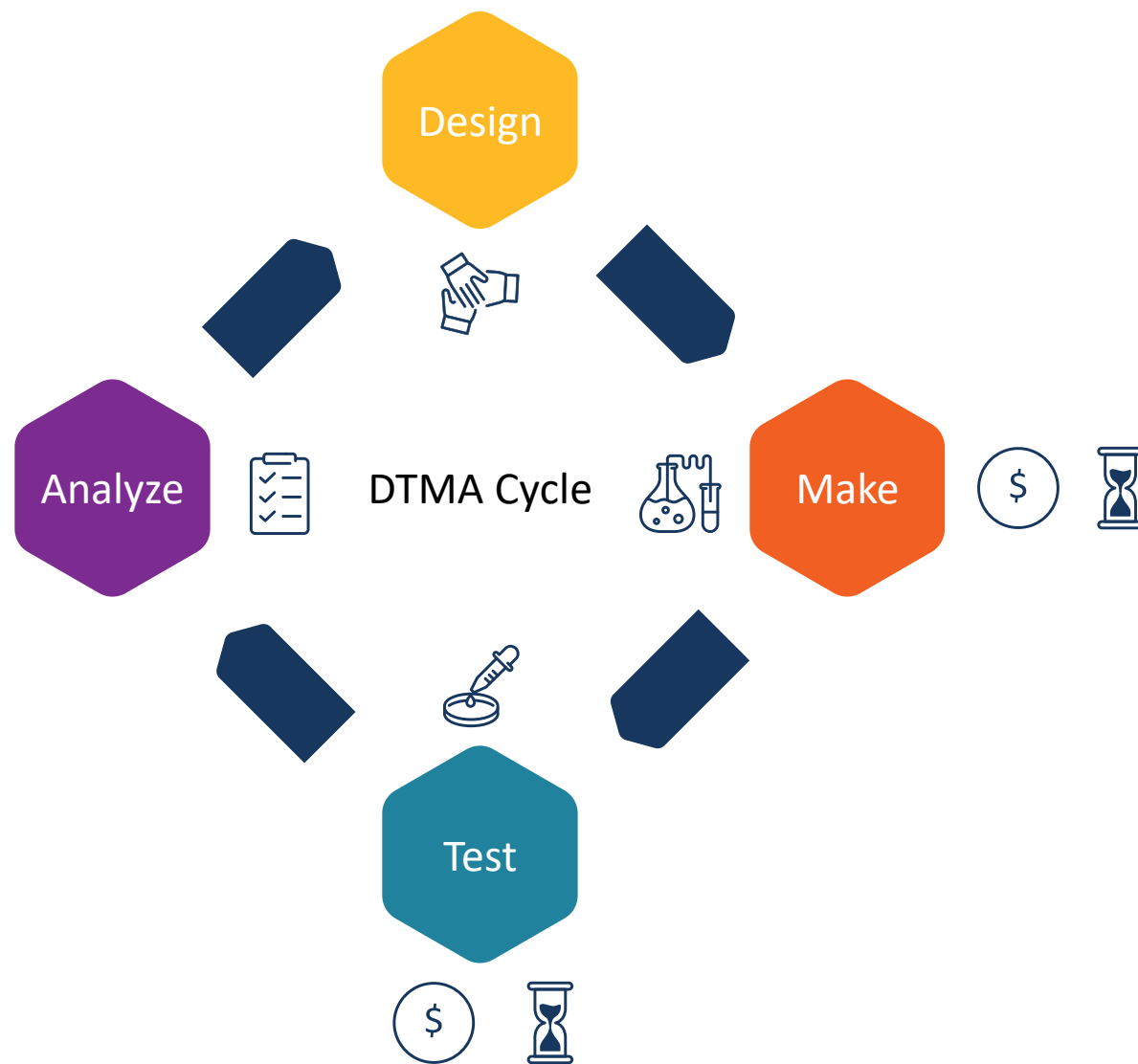
Broad application of Machine Learning in Drug Discovery



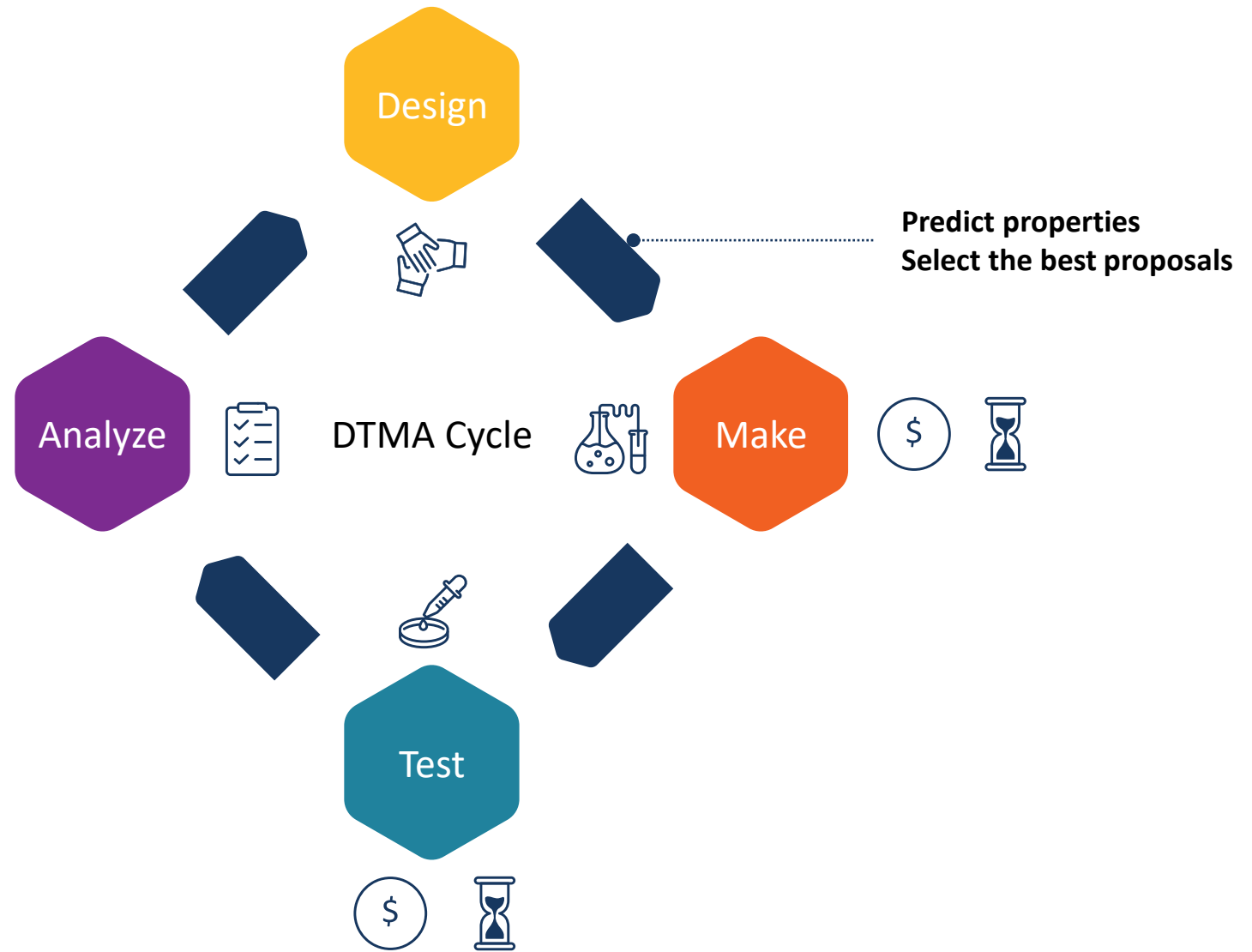
Problem: Rank molecule proposals before synthesis



Problem: Rank molecule proposals before synthesis



Solution: ADMET property prediction models

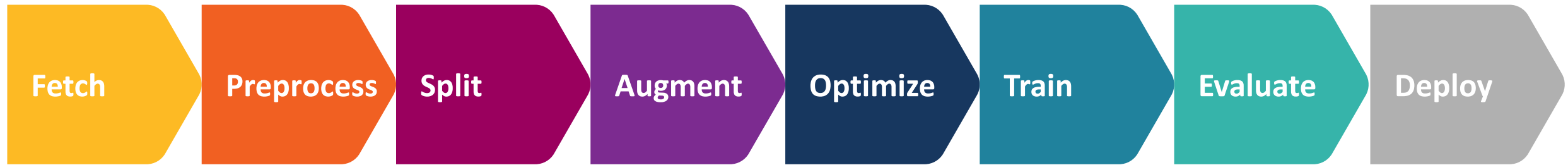


ADMET (absorption, distribution, metabolism, and excretion-toxicity)

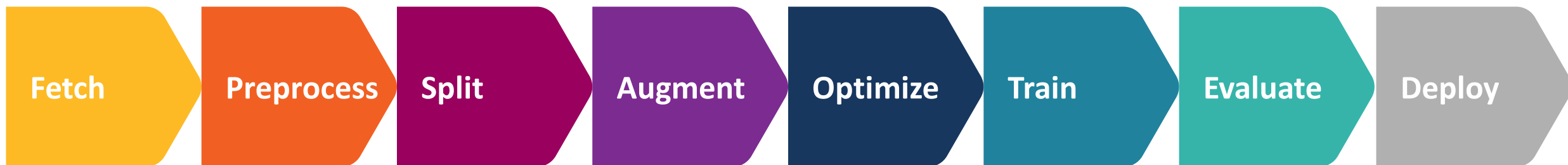
Property prediction model training pipeline at Ryvu



Typical machine learning model training pipeline



Multiple tools used in Ryvu property prediction model training pipeline



Data versioning is a key to experiment reproducibility

Fetch



Internal data
(train + test)

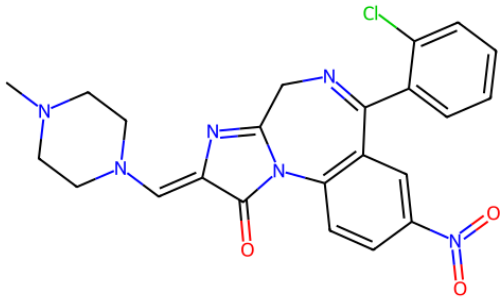


External data
(train only)

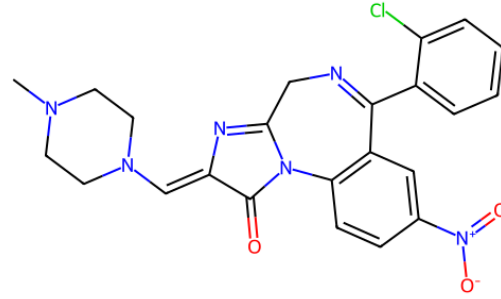


Different SMILES but the same compound: Importance of molecule standardization and data cleaning

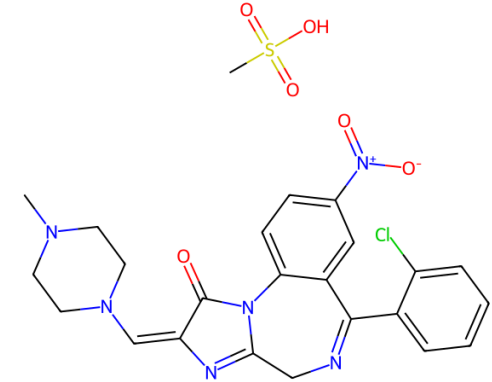
Preprocess



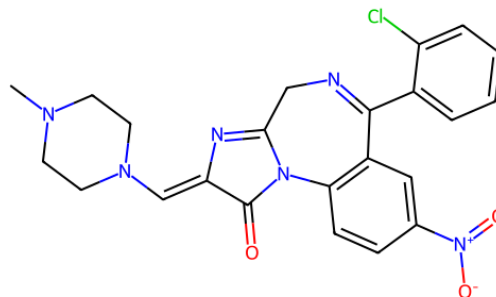
CN1CCN(CC1)\C=C1/N=C2CN=C(C3=CC=CC=C3Cl)C3=C(C=CC(=C3)[N](=O)=O)N2C1=O



CN1CCN(CC1)\C=C1/N=C2CN=C(C3=CC=CC=C3Cl)C3=C(C=CC(=C3)[N+](=[O-])=O)N2C1=O



CS(O)(=O)=O.CN1CCN(CC1)\C=C1/N=C2CN=C(C3=CC=CC=C3Cl)C3=CC(=CC=C3N2C1=O)[N+](=[O-])=O

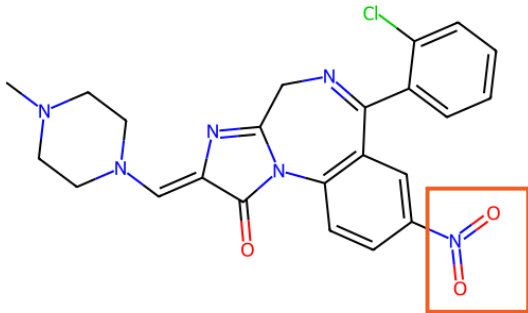


Loprazolam

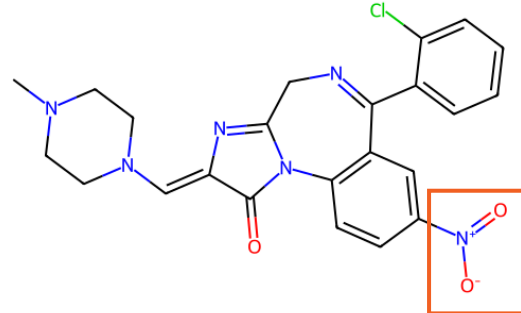
CN1CCN(C=C2N=C3CN=C(c4ccccc4Cl)c4cc([N+](=O)[O-])ccc4N3C2=O)CC1

Different SMILES but the same compound: Importance of molecule standardization and data cleaning

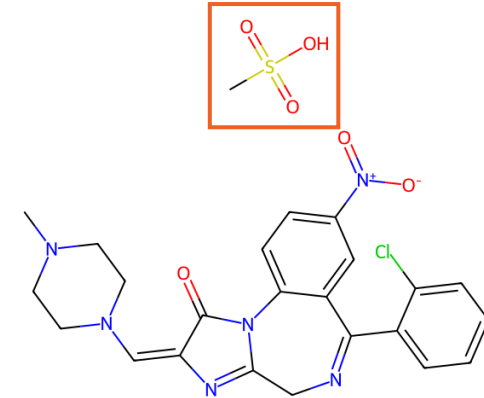
Preprocess



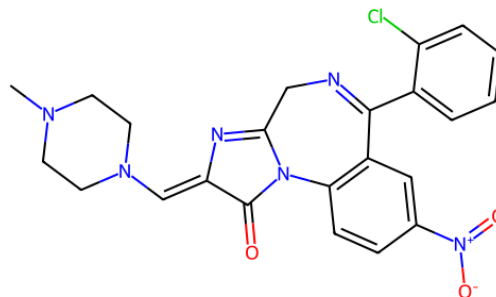
CN1CCN(CC1)\C=C1/N=C2CN=C(C3=CC=CC=C3Cl)C3=C(C=CC(=C3)[N](=O)=O)N2C1=O



CN1CCN(CC1)\C=C1/N=C2CN=C(C3=CC=CC=C3Cl)C3=C(C=CC(=C3)[N+](=[O-])=O)N2C1=O



CS(O)(=O)=O.CN1CCN(CC1)\C=C1/N=C2CN=C(C3=CC=CC=C3Cl)C3=CC(=CC=C3N2C1=O)[N+](=[O-])=O



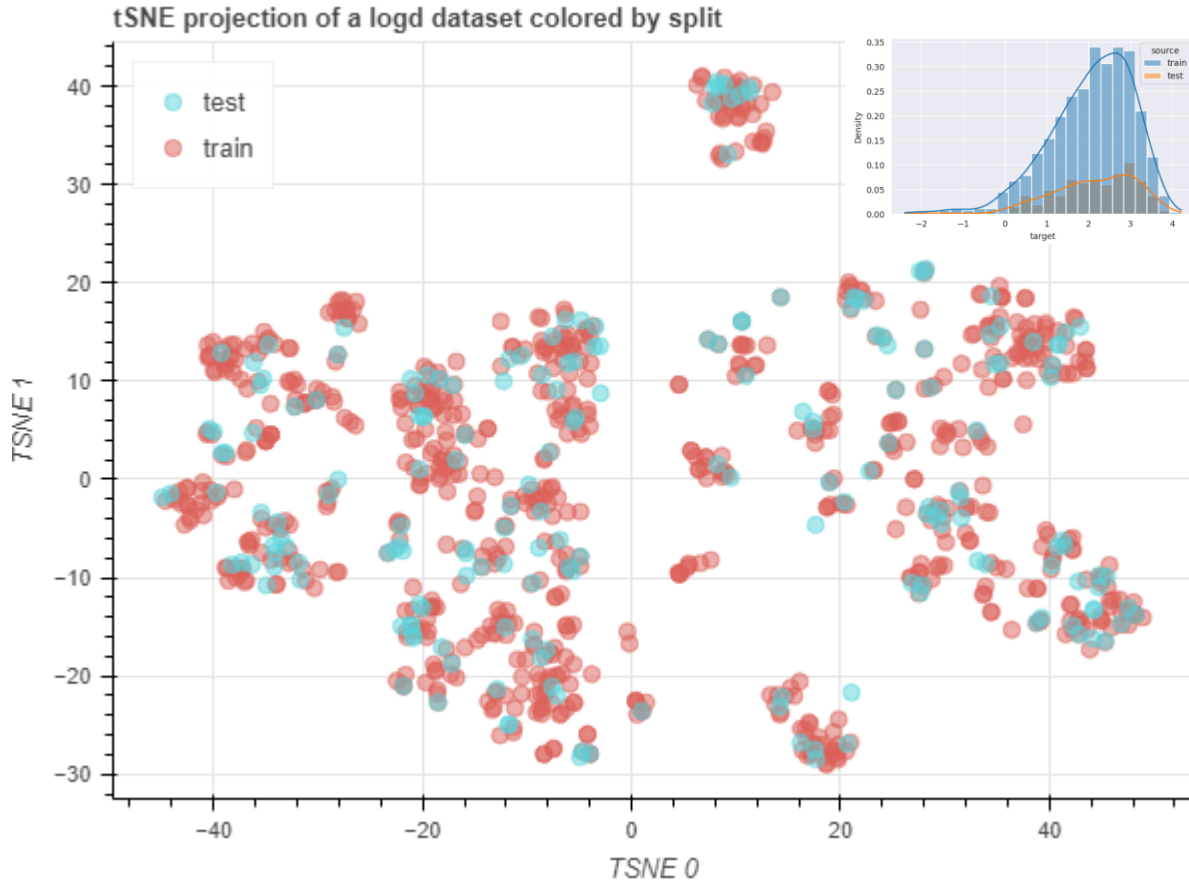
Loprazolam

CN1CCN(C=C2N=C3CN=C(c4ccccc4Cl)c4cc([N+](=O)[O-])ccc4N3C2=O)CC1

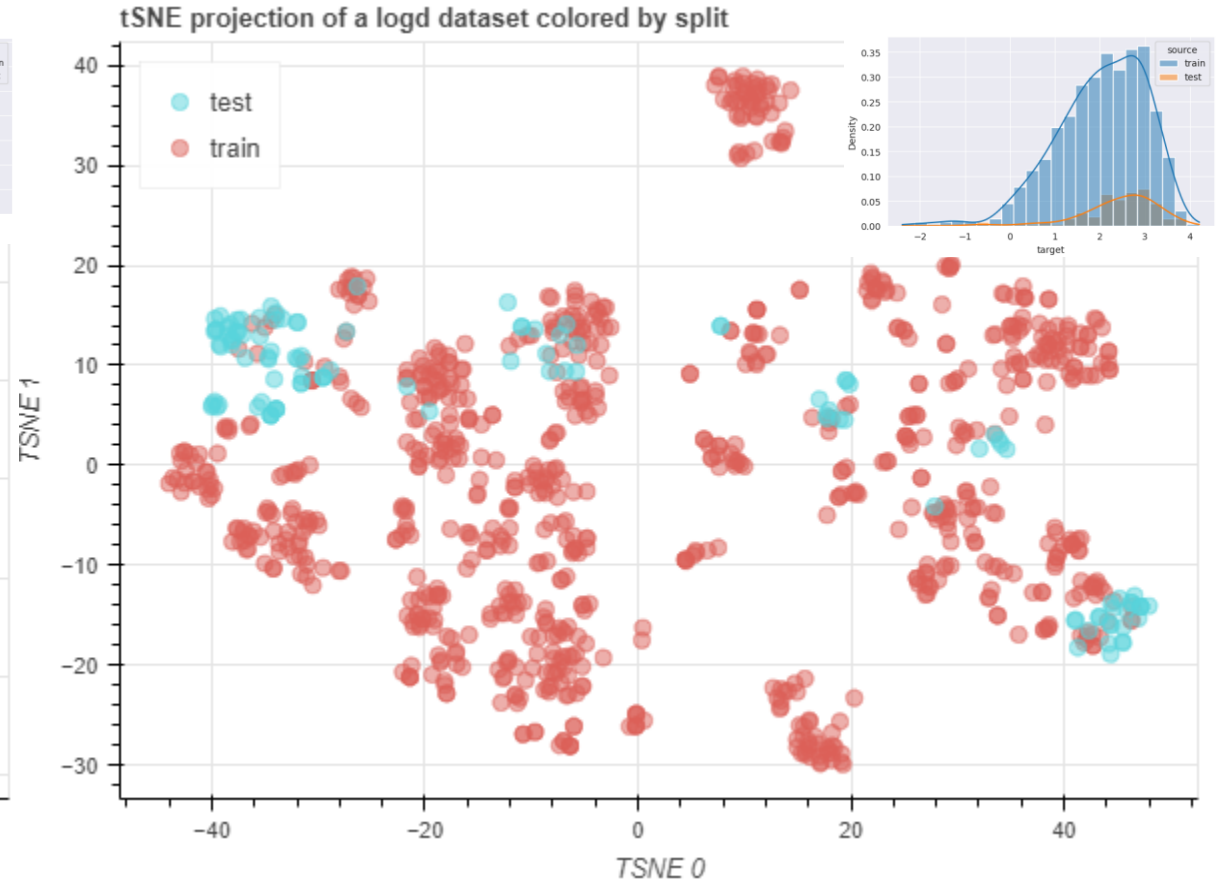
Use scaffold or time split to detect problems with generalization and have reliable estimation of model quality

Split

Random split



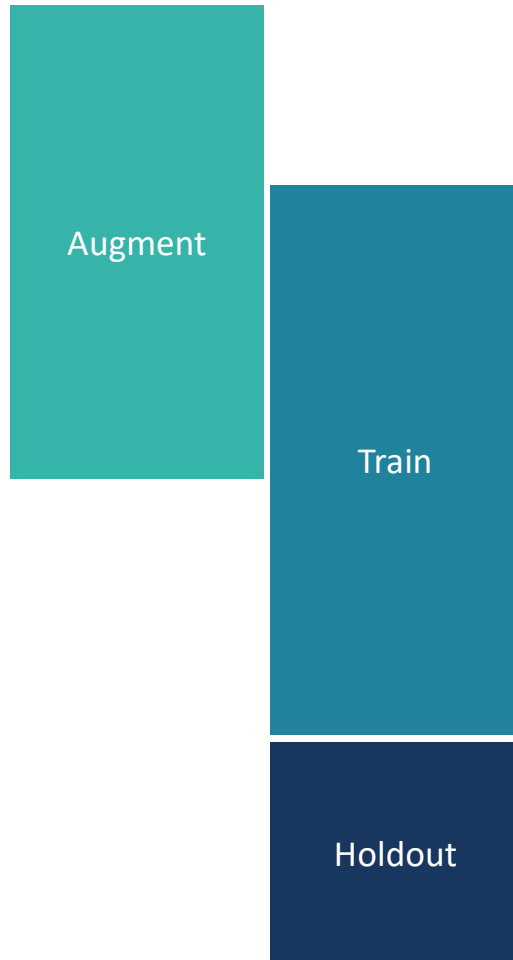
Time split (similar to scaffold split)



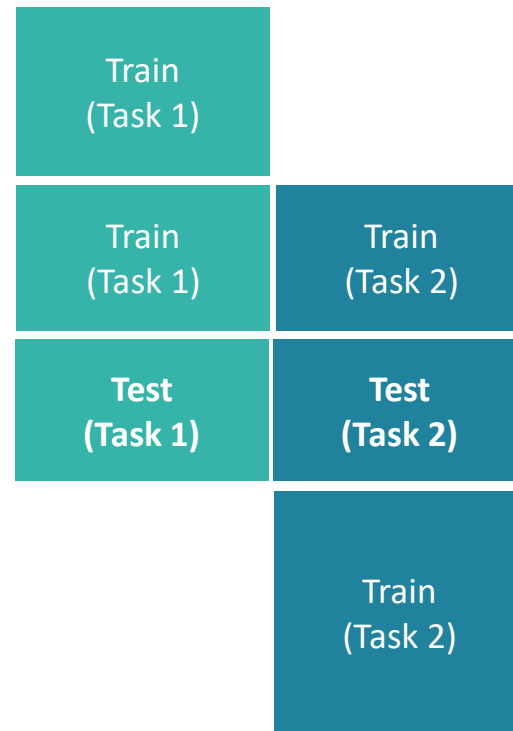
Multi-task learning or transfer learning can enhance model quality

Augment

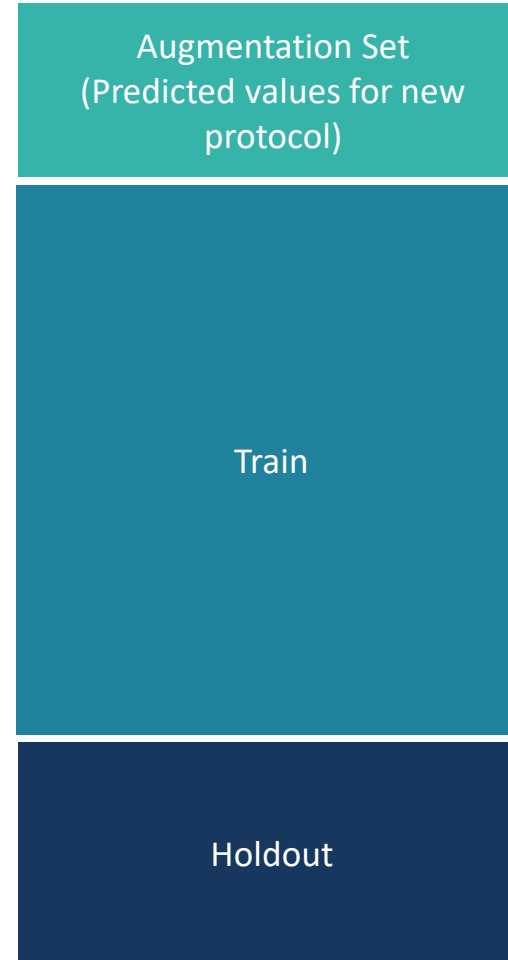
Old Protocol New Protocol



Multi-task learning



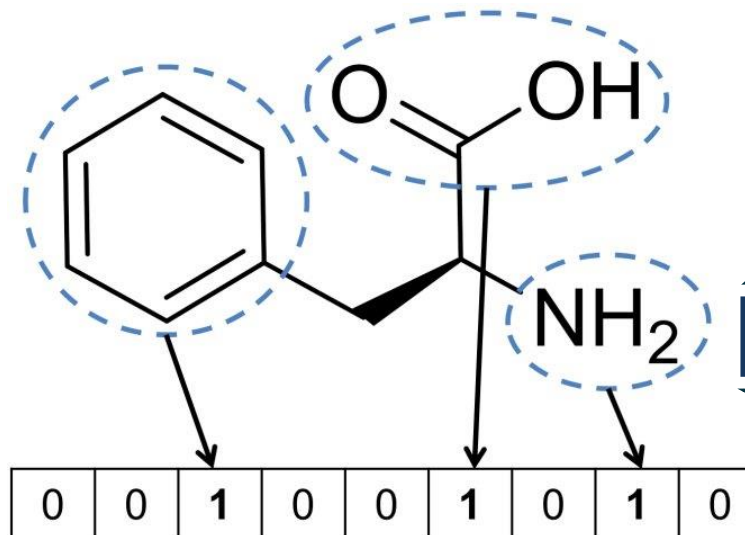
Final set (new protocol)



Needs special handling during CV

Find the optimal feature generation technique

Name	Description
MW	Molecular weight
nHDon	Number of donor atoms for hydrogen bond (HB)
nHAcc	Number of acceptor atoms for HB
SA	Total surface area
TPSA	Topological polar surface area
nSK	Number of non-H atoms
nsp3	Number of sp ³ hybridized carbon atoms
RBN	Number of rotatable bonds
ARR	Aromatic ratio
cLogP	Calculated partition coefficient between octanol and water
nAR	Number of aromatic rings
Fsp ³	Fraction of sp ³ carbon atoms

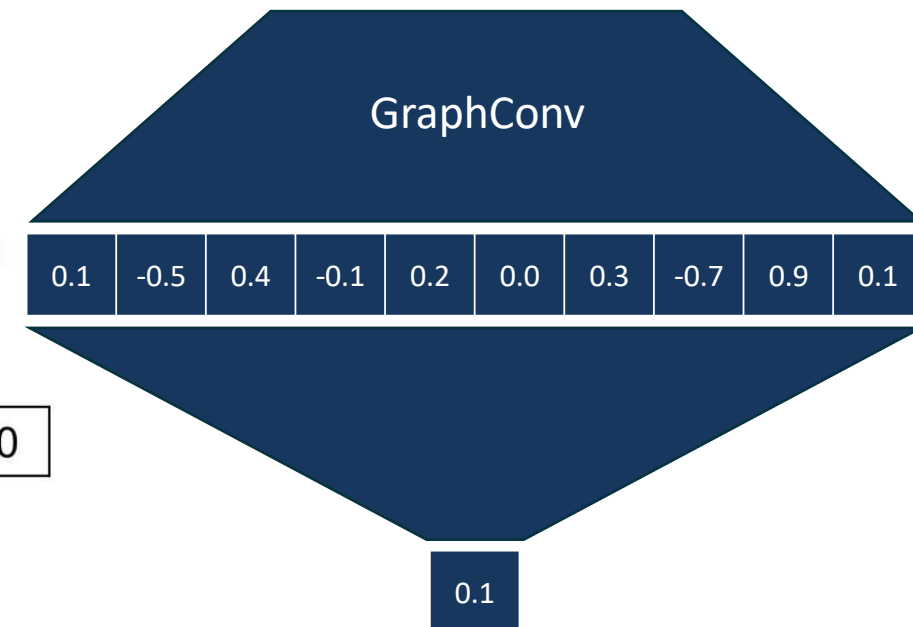
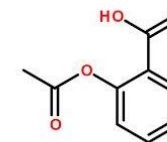


- **RDKit** descriptors
- ...
- **Mordred** descriptors

- **MACCS** keys
- **ECFP4** fingerprints
Extended-Connectivity Circular Fingerprints
- ...
- **PubChem** fingerprints

- Molecule **Embeddings** (trainable)
- ...
- Sequence to sequence autoencoders

Optimize

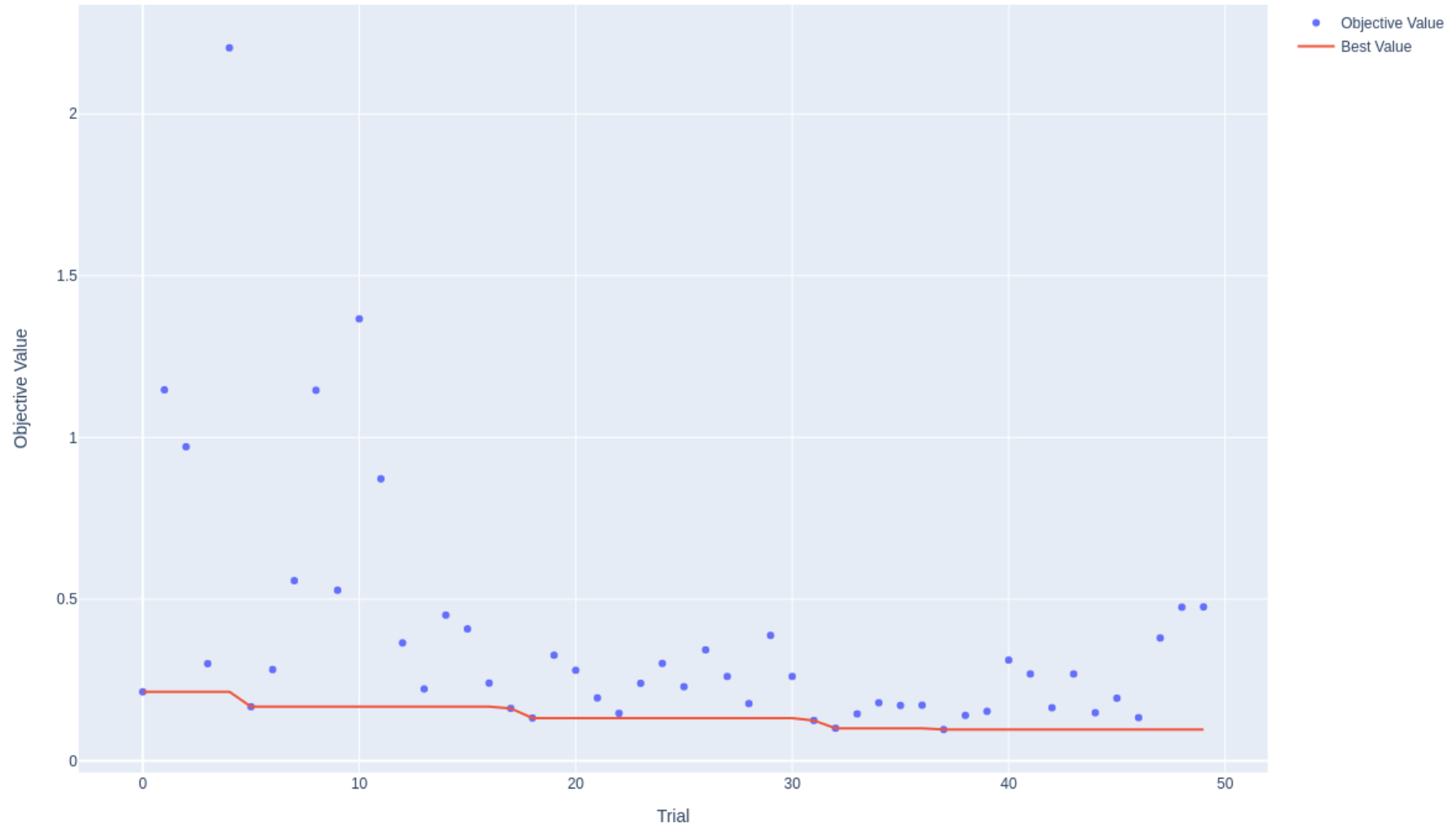


¹https://www.researchgate.net/figure/In-substructure-key-based-fingerprints-bits-are-set-according-to-the-substructures-that_fig10_315513438

Avoid grid search – use Optuna instead

Optimize

Optimization History Plot



Model zoo:

It is important to test various model architectures

Train

Random Forest

XGBoost
LightGBM

Graph Convolutional
Networks
(Message Passing)

Gaussian Process
Regressor

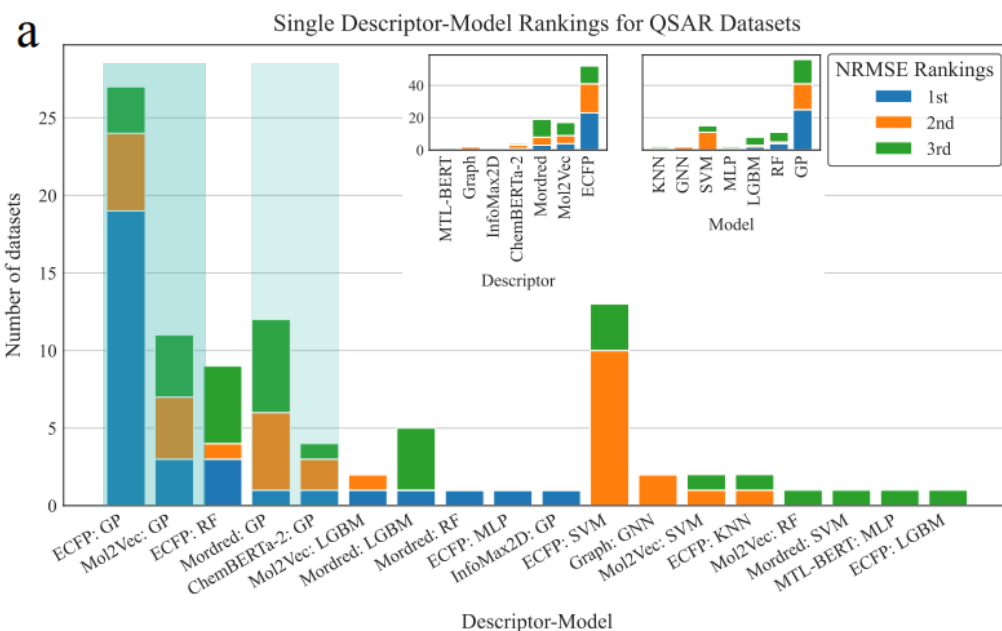


GPyTorch

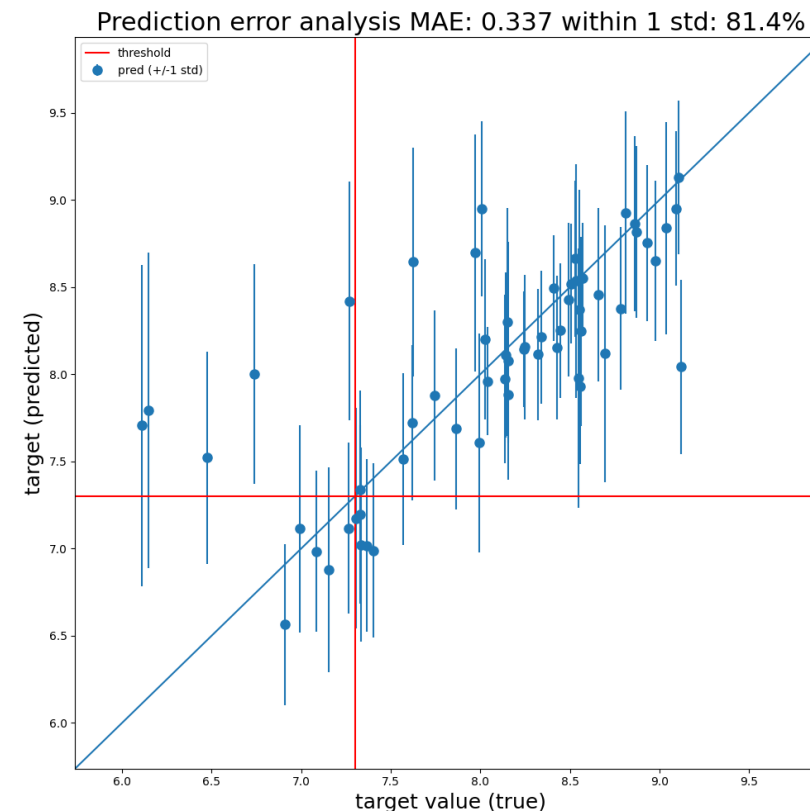
Gaussian process is effective for QSAR modeling

Train

- Gaussian Process Regressors are effective in Quantitative Structure-Activity Relationships (QSAR) modeling¹
- Gaussian processes have embedded prediction **uncertainty estimation**



Single descriptor-model rankings based on NRMSE¹



Gaussian processes have embedded prediction uncertainty estimation

¹ <https://browse.arxiv.org/pdf/2309.17161.pdf>

Pick the right metrics for the problem

Evaluate

Binary classification

Simple interpretation

ROC AUC

F1

Accuracy

Precision

Recall

Ranking

Prioritization

Spearman correlation

NDCG

Regression

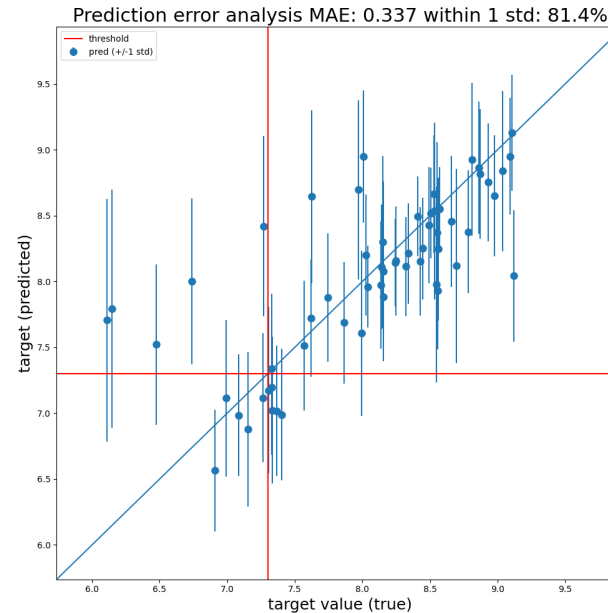
Exact value prediction

MSE

MAE

R2

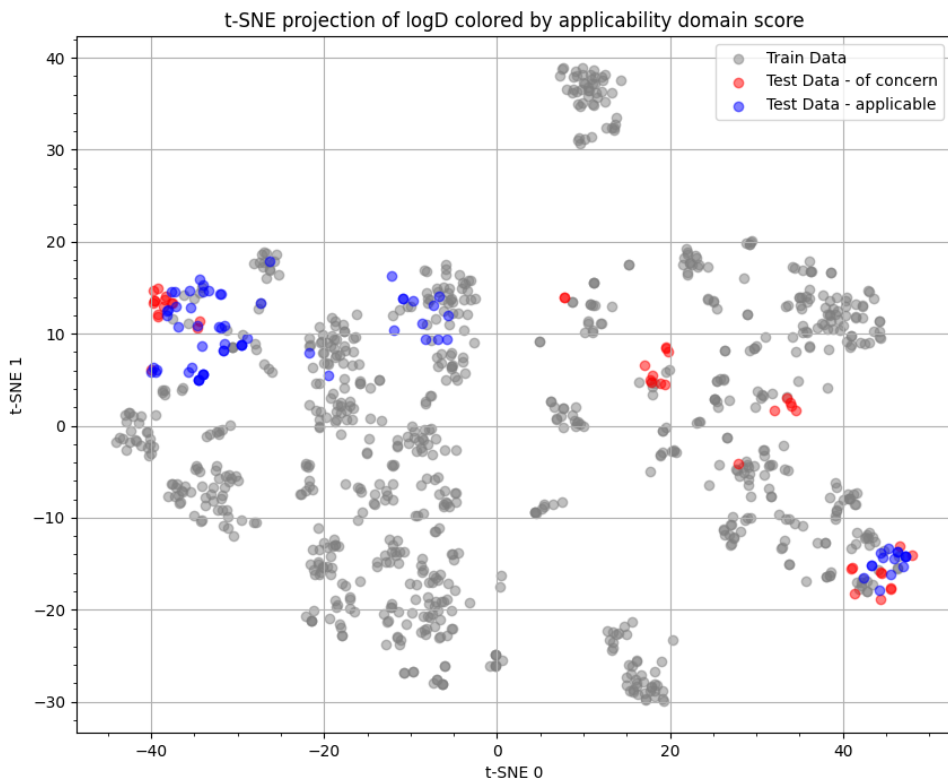
Pearson Correlation



Each user has different need – focus on regression

Estimation inference limits via applicability domain analysis

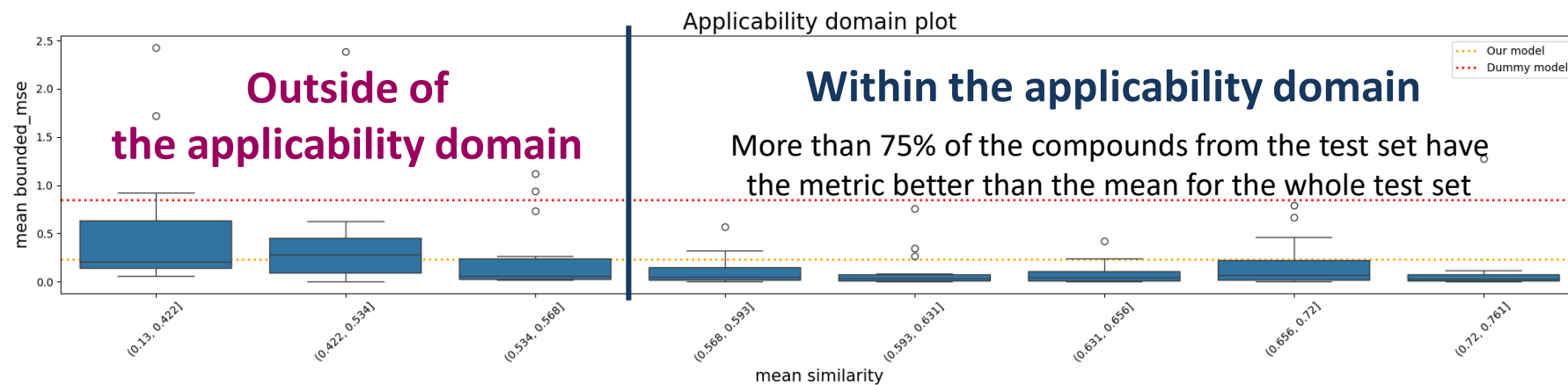
Evaluate



Problem

Chemists would like to know the model inference limitations.

Often, models can't generalize to different chemical series, further away from the training set.



Save time with deployment and model updates automation

Deploy

- **Regular model updates**
Ideally, once per DTMA cycle
- **Automation**
Deployment of multiple models can take time
- **Integration**
with software tools that chemists use for daily work



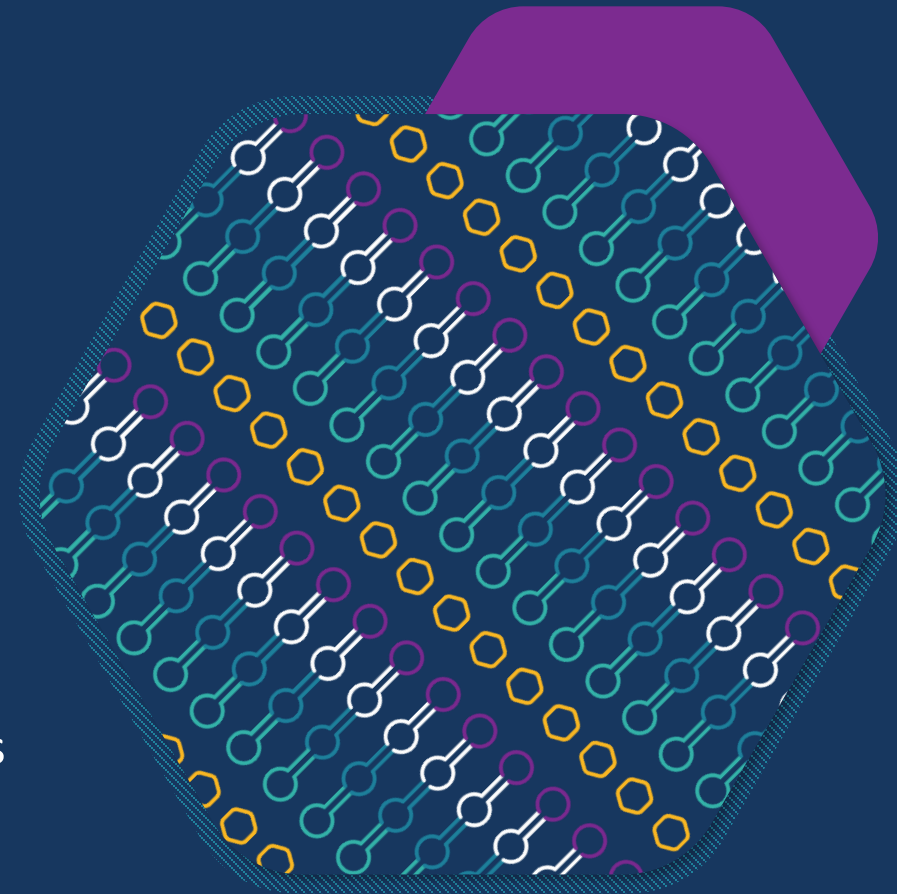
 FastAPI

Dotmatics

Regular model retraining is a key to keep high model quality

Summary

- Ryvu Therapeutics is leading Polish drug discovery company
- AI can improve every stage of the drug discovery process
- Property prediction models can improve lead optimization process
- High quality models can be achieved through combination of data science, data engineering and understanding of the field
- Models' interpretability is important for medicinal chemists
- It is possible to automate many steps of the property prediction model training pipeline



Thank you!

Marcin Kowiel

Team Lead Machine Learning Engineer

Data Science and AI Platforms

marcin.kowiel@ryvu.com



Open positions

<https://ryvu.com/careers/job-offers/>

Data Engineering Intern

Data Scientist/AI Research Intern

Computer Aided Drug Design Intern

Bioinformatics Intern

Computer Aided Drug Design Senior Scientist

Senior IT Administrator

