# AI-Accelerated Discovery of Novel Aqueous Amines for CO$_2$ Capture

**Team Aspirin**

Charlotte Breakwell

Luis Carvalho

Kieran Nehil-Puleo

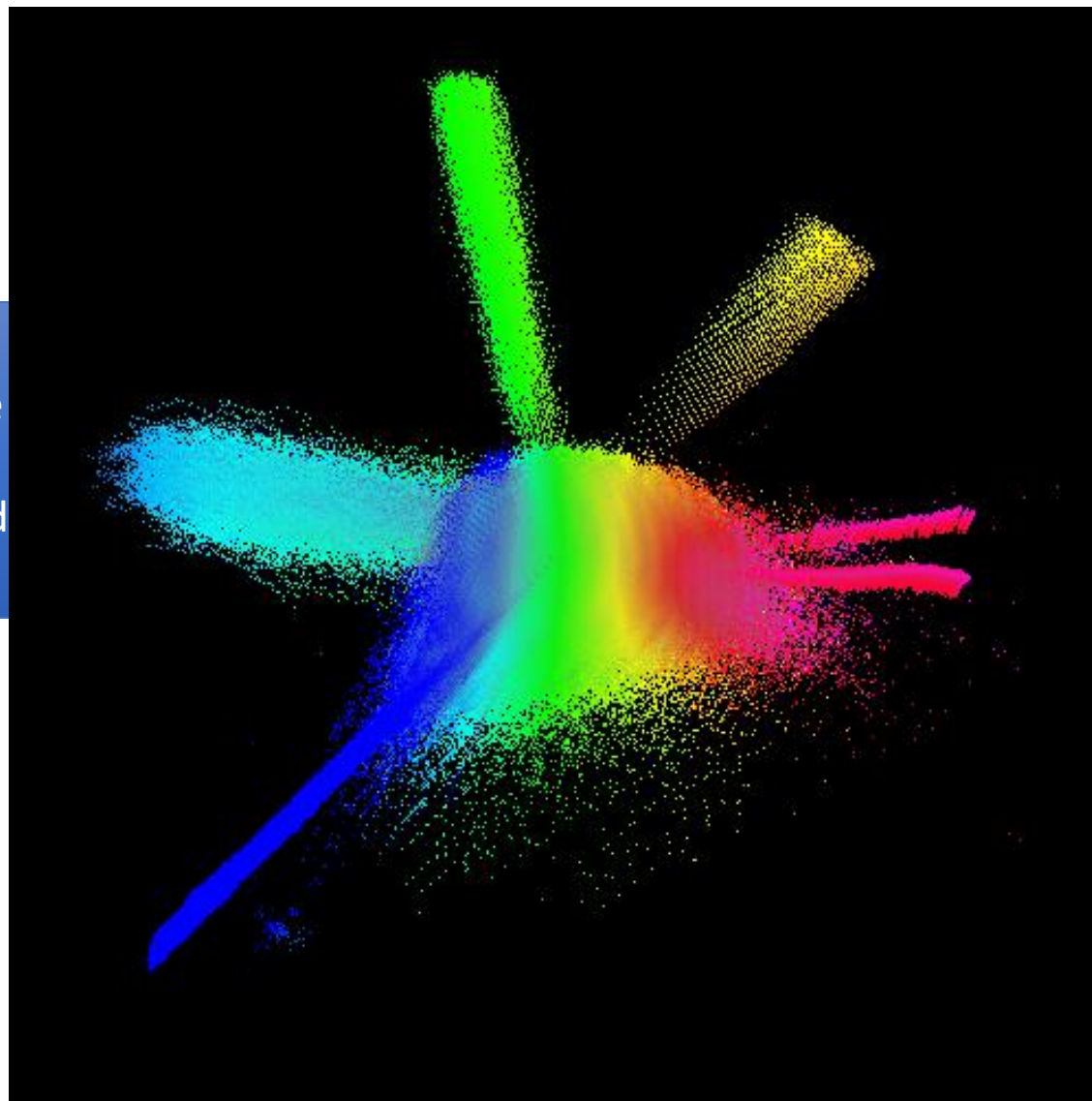Evangelos Tsochantaris

Eco-AI Hackathon

March 2024

# AI-Aided Discovery

Generate potential compounds → Select compounds with optimal properties ⇢ Synthesise candidates
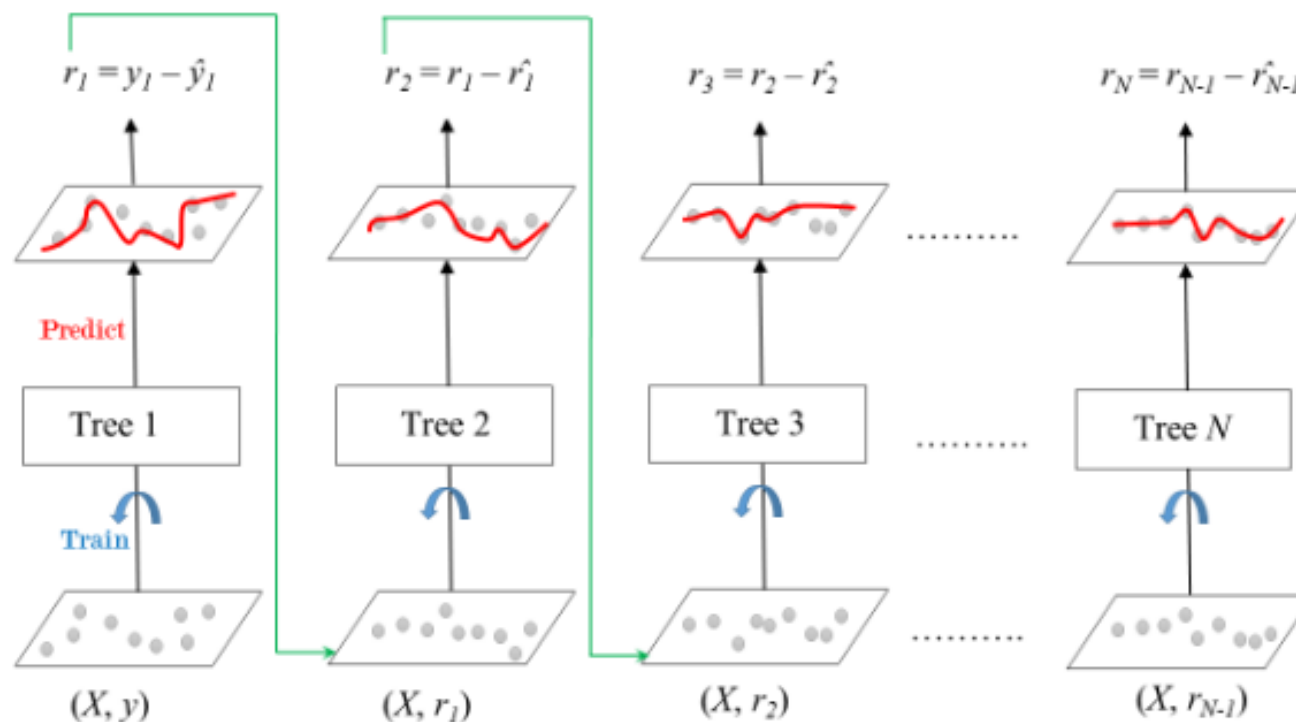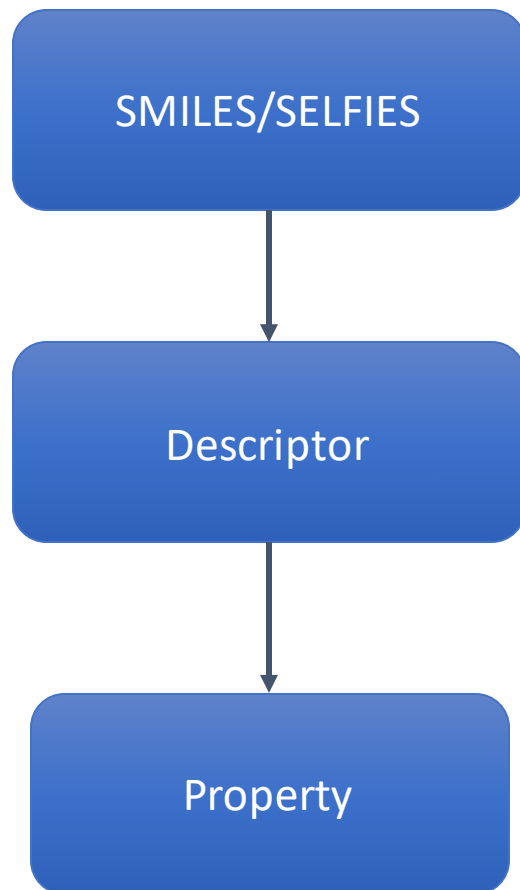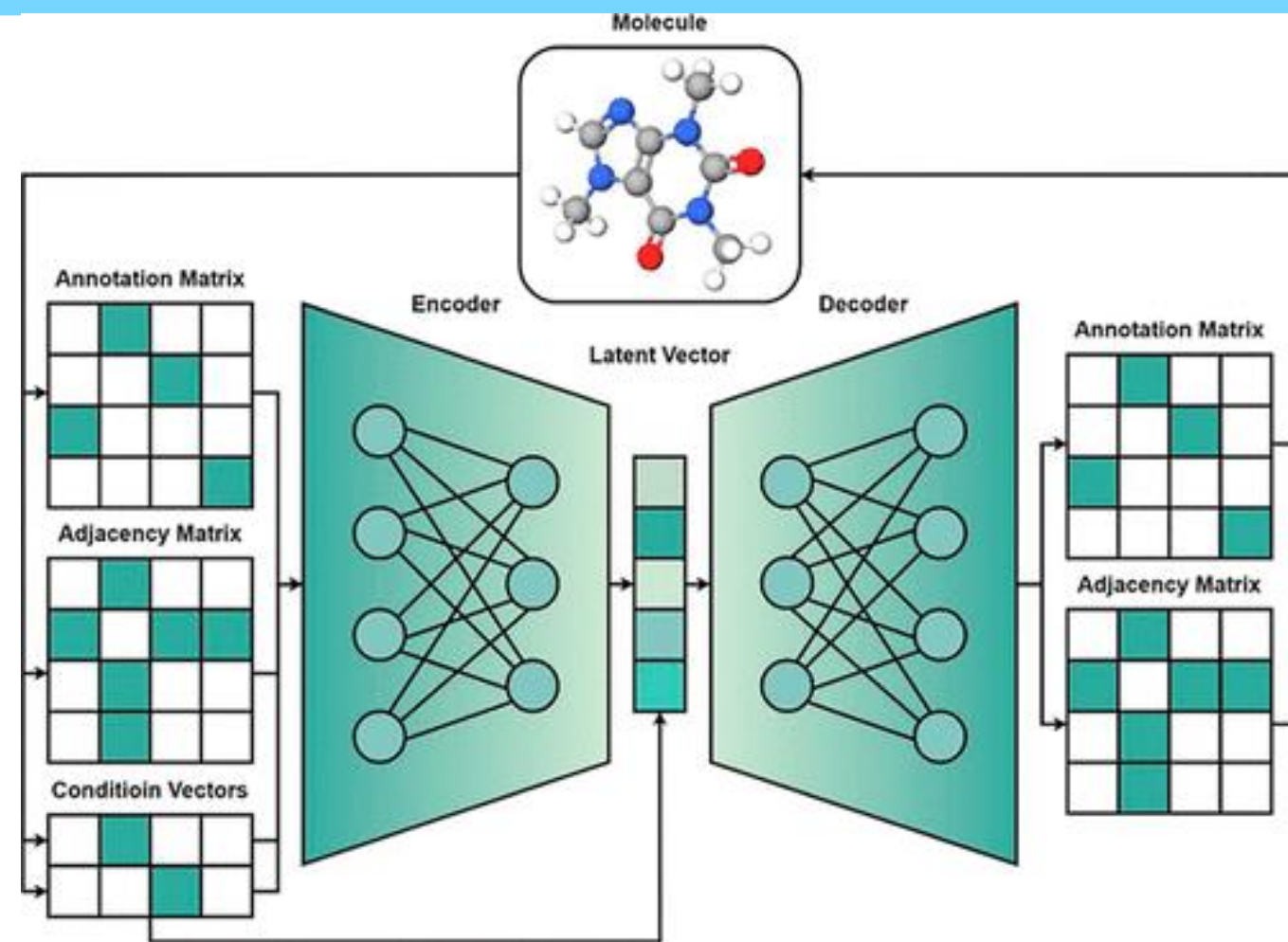
# AI-Aided Discovery
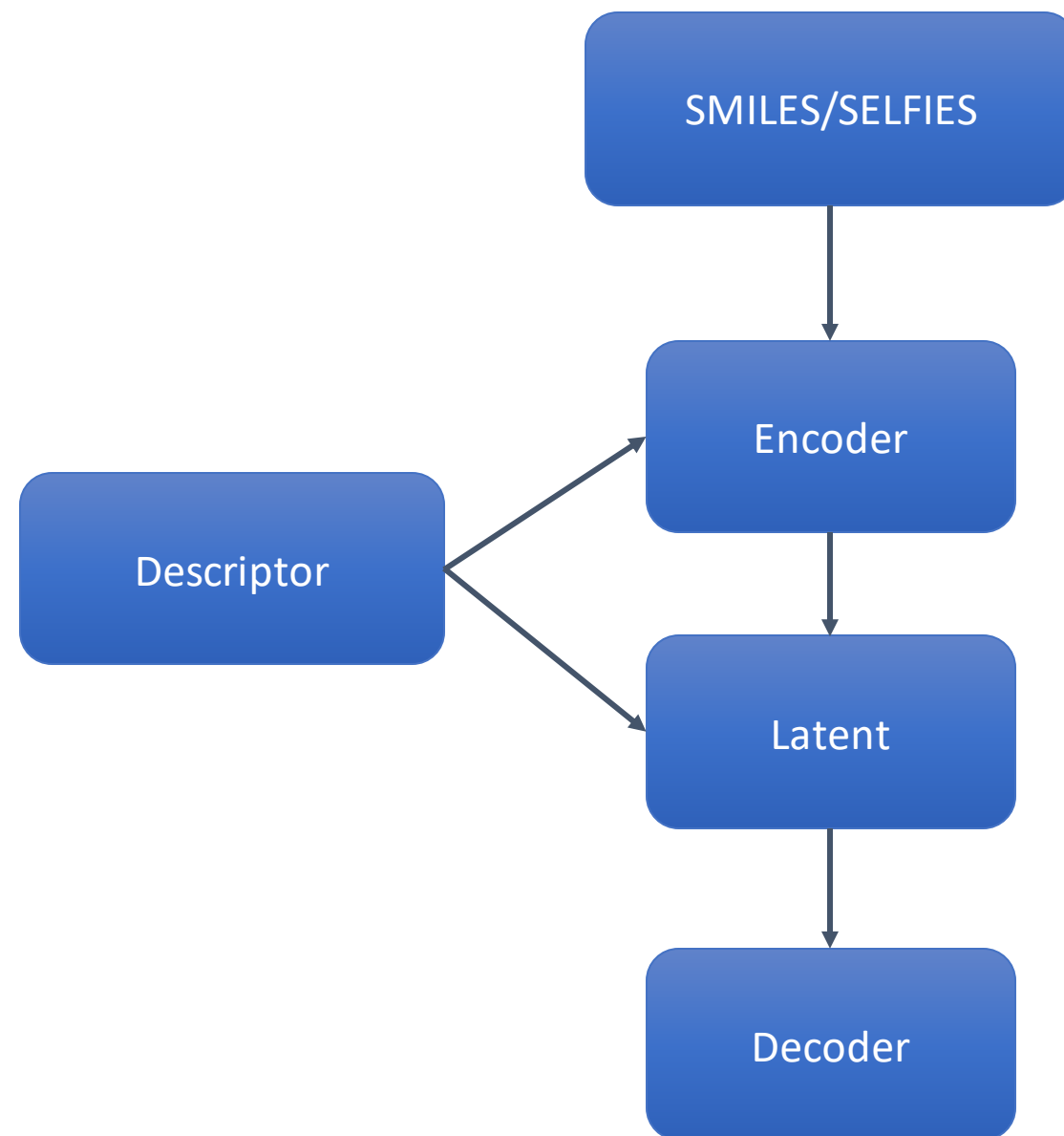


Generate potential compound

Synthesise candidates

# Proposed Workflow (Predictive Model)

SMILES/SELFIES

Descriptor

Property



$$r_1 = y_1 - \hat{y}_1$$

$$r_2 = r_1 - \hat{r_1}$$

$$r_3 = r_2 - \hat{r_2}$$

$$r_N = r_{N-1} - \hat{r_{N-1}}$$

Predict

Tree 1

Tree 2

Tree 3

Tree $N$

Train

$(X, y)$

$(X, r_1)$
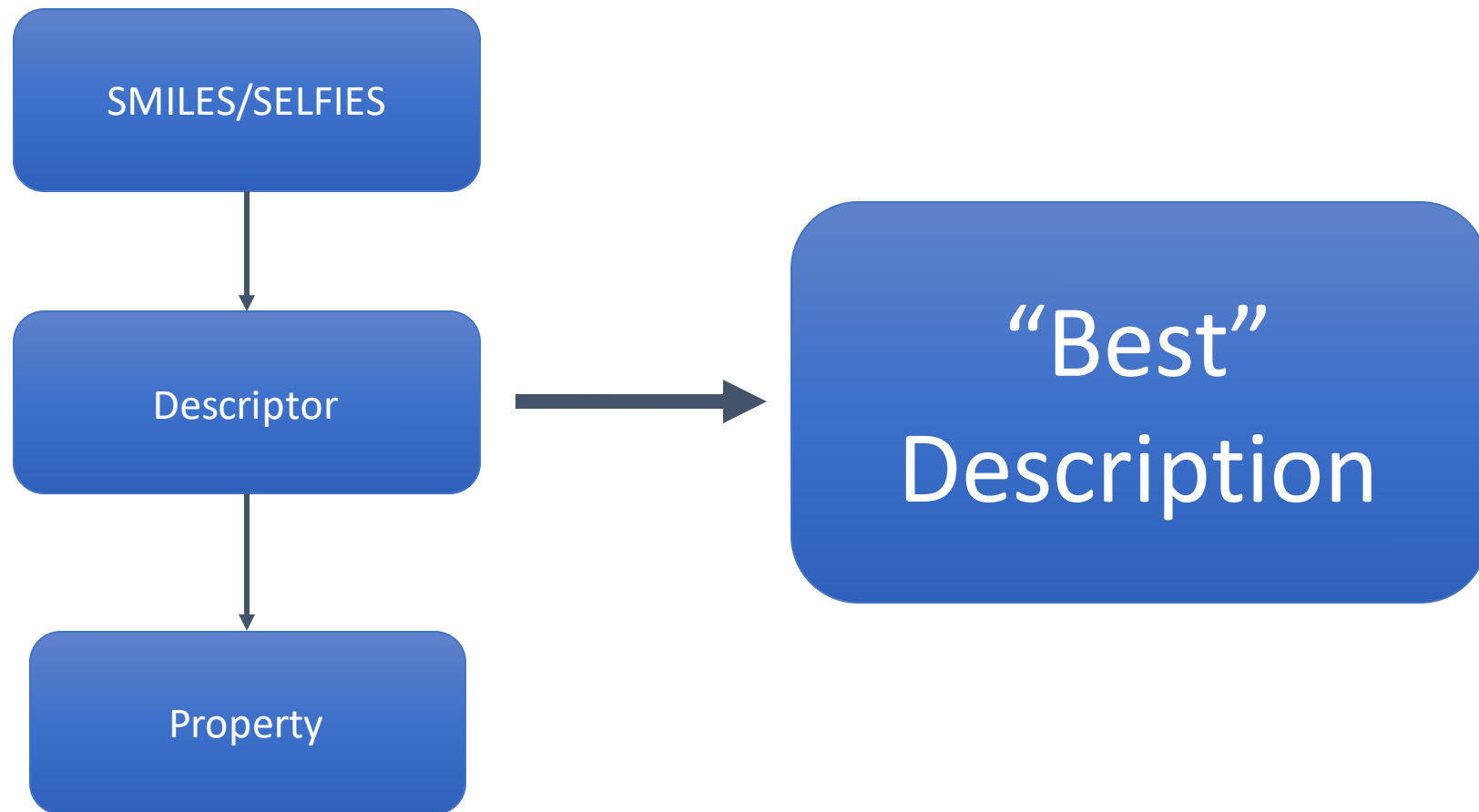
$(X, r_2)$

$(X, r_{N-1})$

https://www.geeksforgeeks.org/ml-gradient-boosting/

# Proposed Workflow (Conditional VAE)



Lee, Myeonghun, and Kyoungmin Min. "MGCVAE: multi-objective inverse design via molecular graph conditional variational autoencoder." *Journal of chemical information and modeling* 62.12 (2022): 2943-2950.

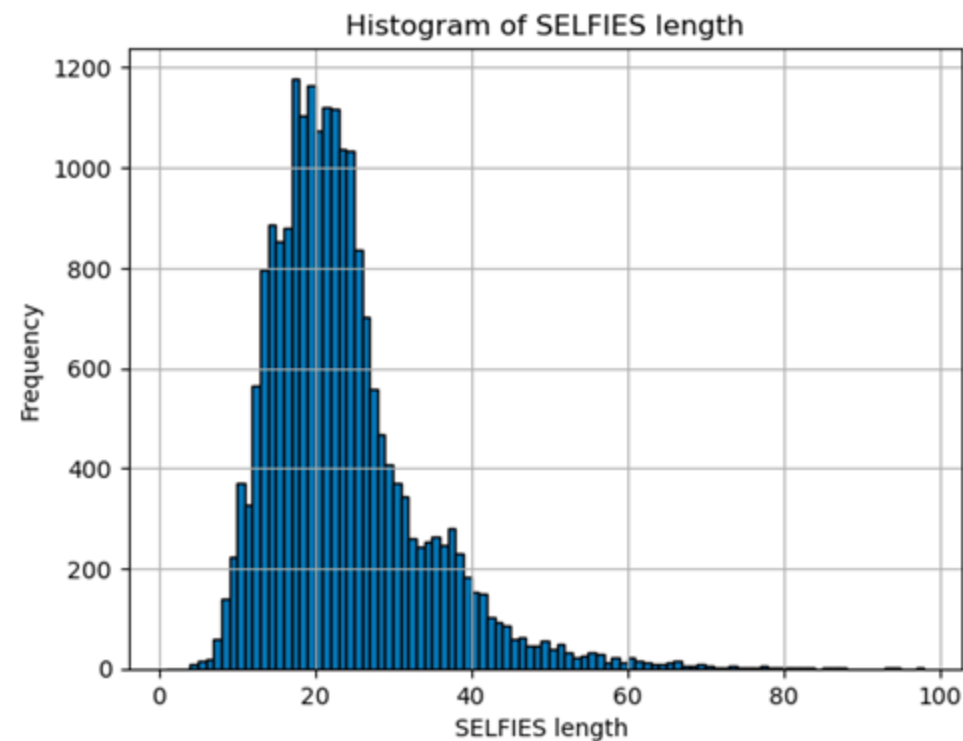# Proposed Workflow

# Smiles vs Selfies

SMILES: CC(=O)Oc1ccccc1C(=O)O

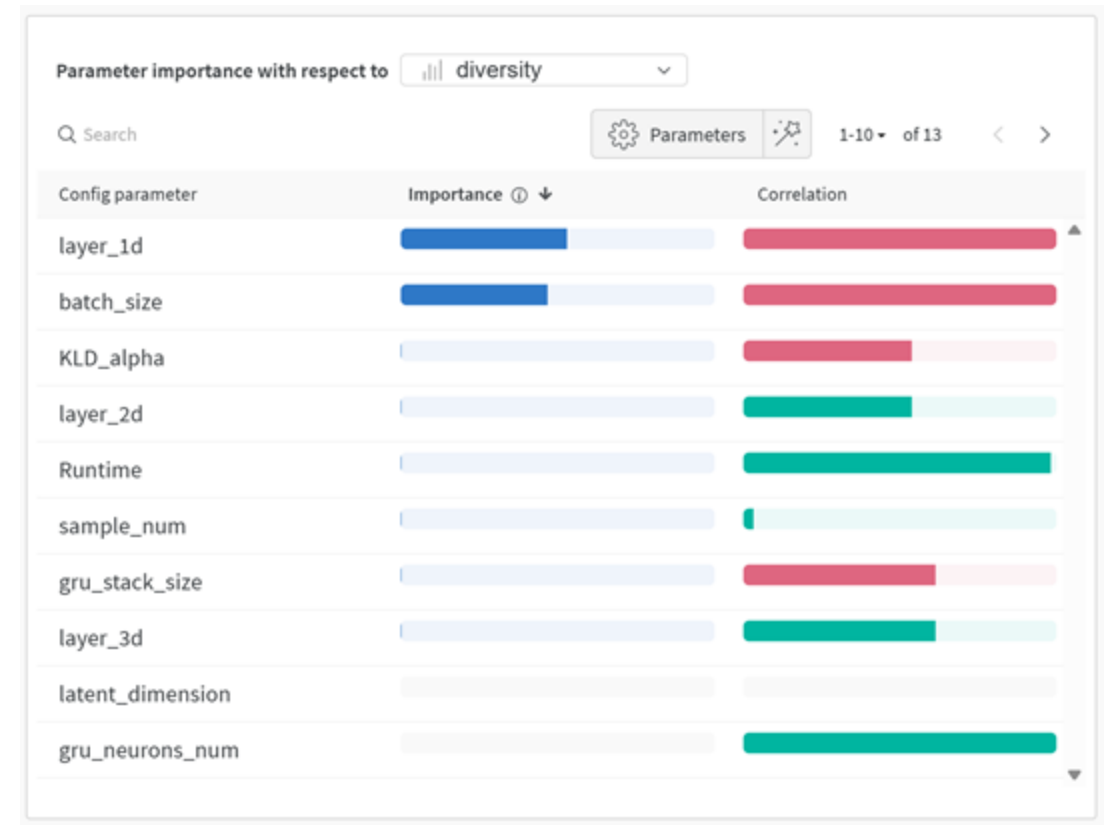SELFIES: [C][C](=[O])[O][C][c][c][c][c][c][c][C]([=[O])[O]

Encoder with
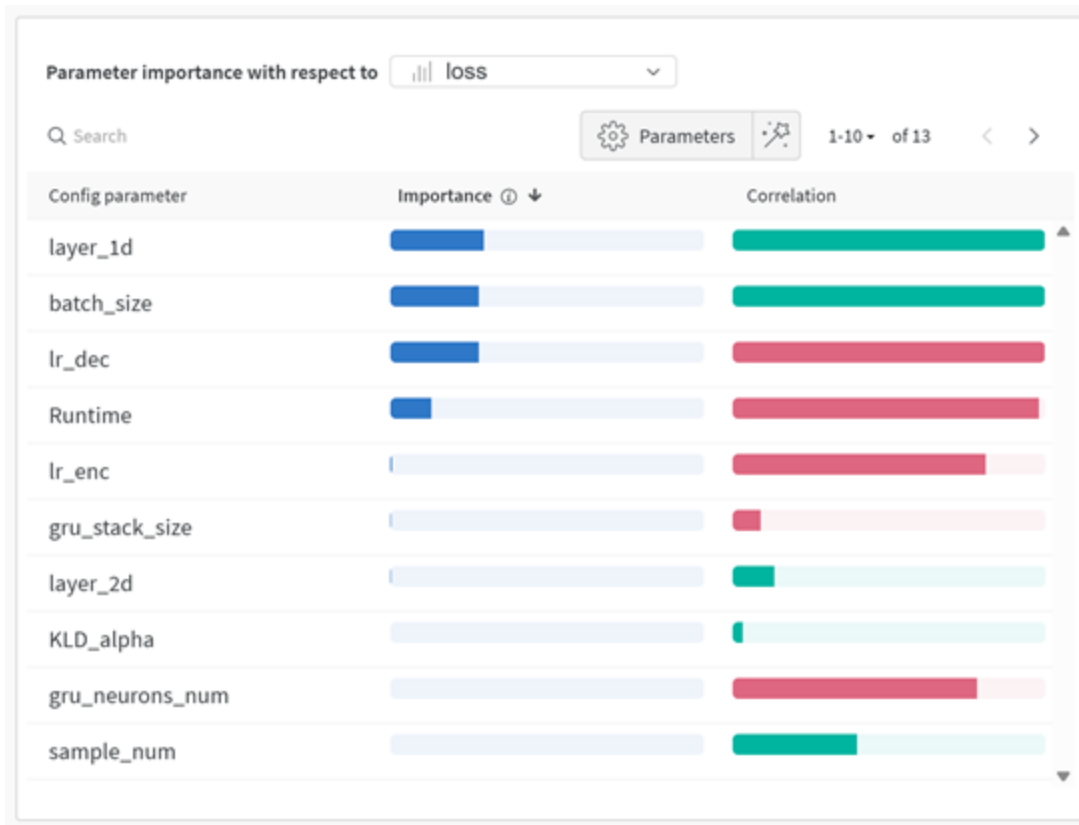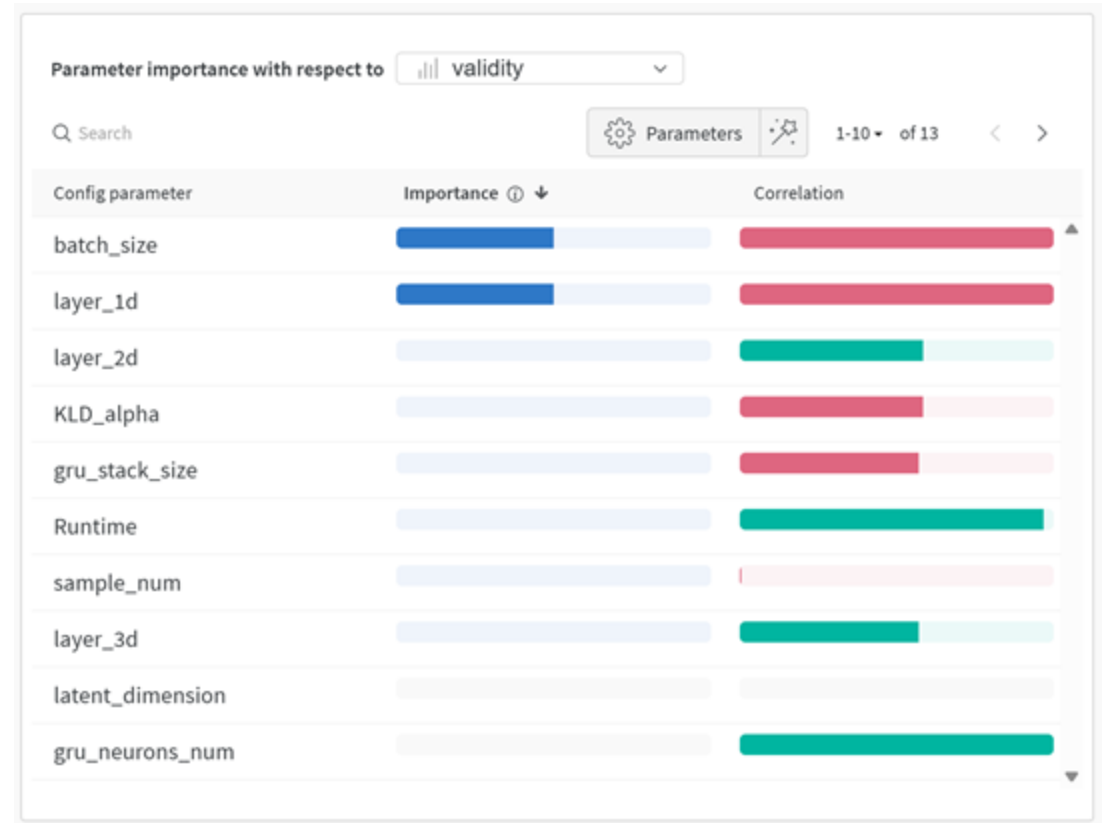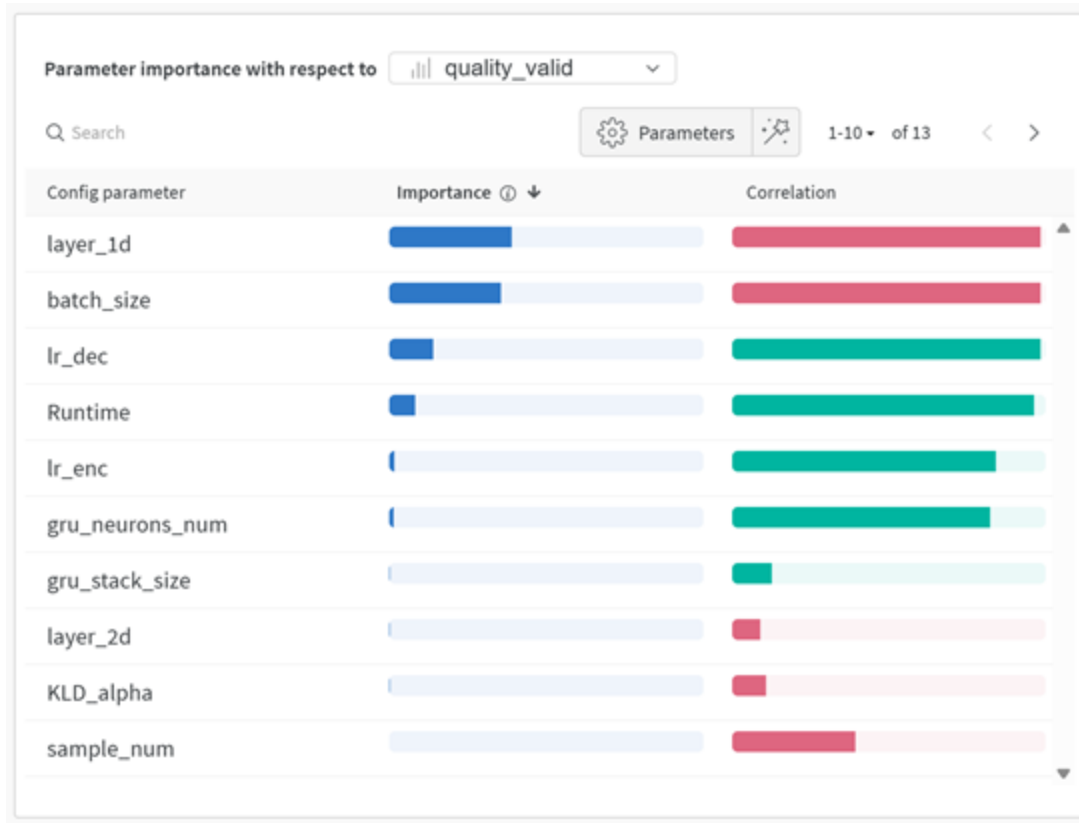
SMILES: ['C', 'O', '(', '=', ...]

SELFIES: ['[C]', '([=[O])', '[O]', ...]

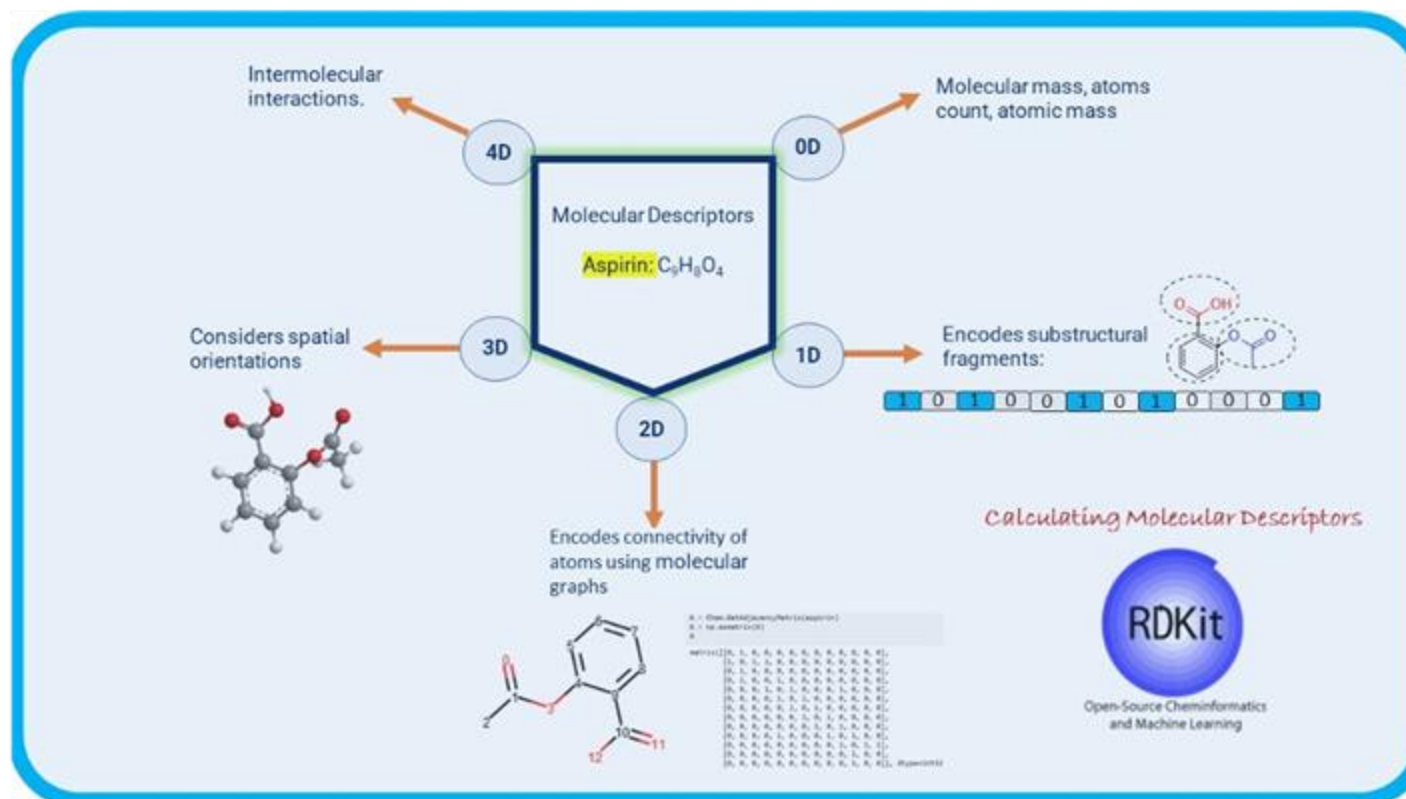# Hyperparameter Optimisation

# Hyperparameter Optimisation

# Finding the Right Molecule (Cheminformatics)

**Cheminformatics depend on representations of molecules by descriptors that capture their structural characteristics and properties.**

- Experimental measurements

- Theoretical measurements
  - RDKit descriptors
  - Simulations (e.g. DFT, MD)
  - Predictions from existing ML models

# Finding the Right Molecule (Cheminformatics)

```python
from rdkit import Chem as Chem
from rdkit.Chem import Descriptors

# Featurize SMILES strings
def featurize_smiles(smiles):
    mol = Chem.MolFromSmiles(smiles)
    features =
Descriptors.CalcMolDescriptors(mol)
    return features

# Apply featurization to SMILES strings
in the dataset
features =
data['smiles'].apply(featurize_smiles)
```

**Raw Dataset**

| 1 | SMILES |
|---|--------|
| 2 | $CO_2$ adsorption capacity |
| 3 | n_nitrogen |
| 4 | Molecular Mass |

**Dataset with RDKit Descriptors**

| 1 | SMILES |
|---|--------|
| 2 | $CO_2$ adsorption capacity |
| 3 | n_nitrogen |
| 4 | Molecular Mass |
| 5 | NumValenceElectrons |
| 6 | MolLogP |
|   | ... |
| 210 | BertzCT |
| 211 | MaxAbsPartialCharge |
| 212 | MinAbsEStateIndex |
| 213 | BalabanJ |
| 214 | FpDensityMorgan1 |
| 215 | VSA_EState7 |

# Finding the Right Molecule (Cheminformatics)



**Raw Dataset**
Random Forest
R-squared: -0.19

MSE = 0.297

**+ RDKit Descriptors**
Random Forest
R-squared: 0.30

MSE = 0.216

**+ QM Descriptors**
Random Forest
R-squared: 0.38

MSE = 0.191

* QM descriptors generated using xtb package by performing single point energy calculations in implicit water with PM6 functional (Semi-empirical DFT)

# Finding the Right Molecule (Cheminformatics)

**Feature selection**: By identifying the most important features, we can select a subset of relevant features for our target property.
- Reduced dimensionality and noise in the data

**\*Top Performing Model = Gradient Boosted Regression**
**$R^2$ = 0.49, MSE = 0.18 mol $CO_2$ / mol amine**



## All Descriptors

## Top 50 Descriptors

# Finding the Right Molecule (Cheminformatics)

No Properties

QM Properties

# Exploring The Chemical Space - Further Work

MCTS used in applications with large possibility spaces.

RL methods require suitable reward function.