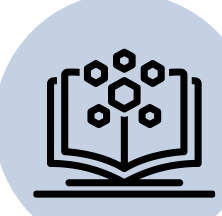# In silico generated reagents for detection of pesticides using mass spectrometry: An out-of-distribution task

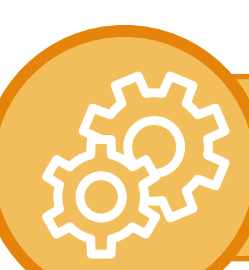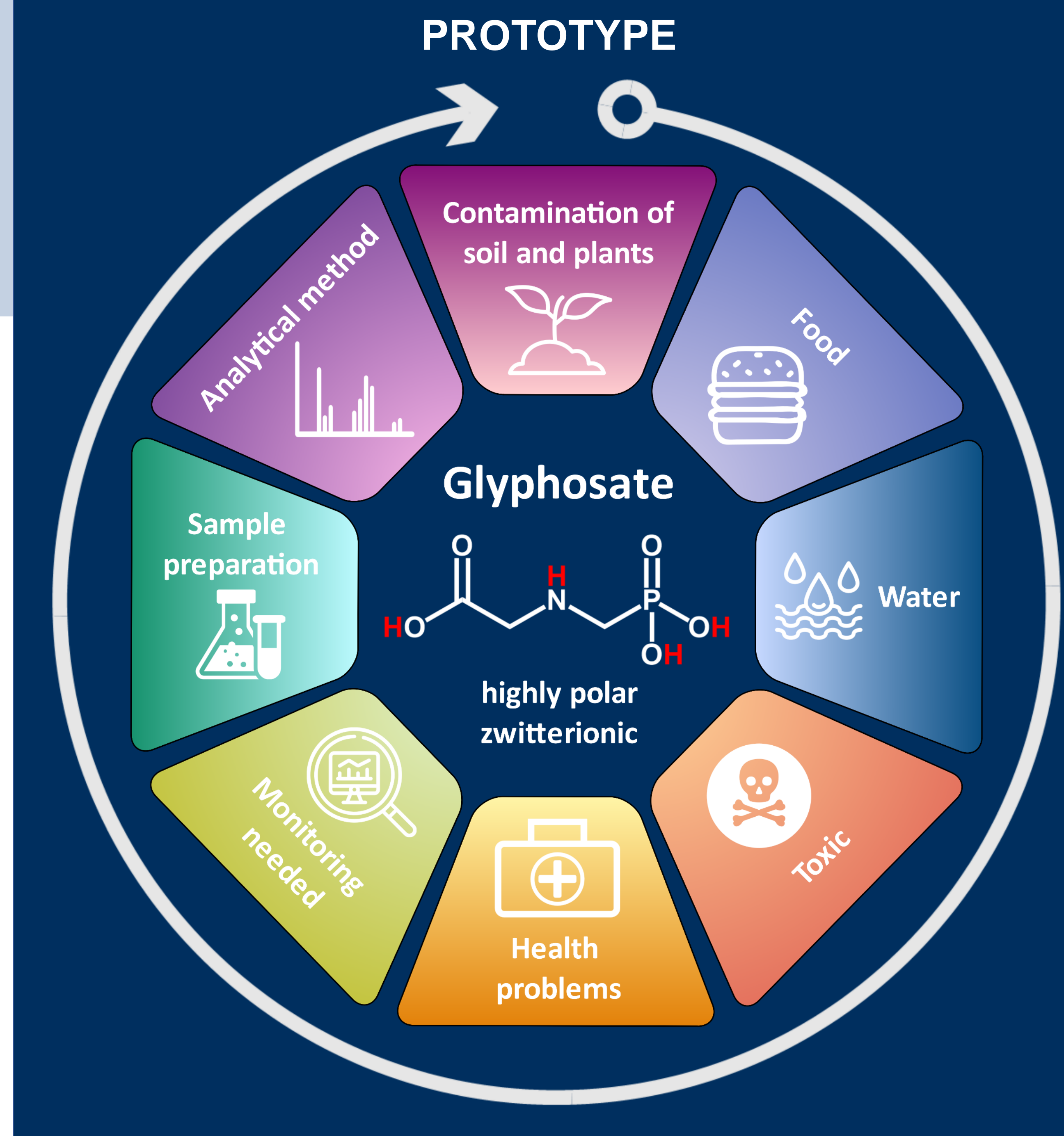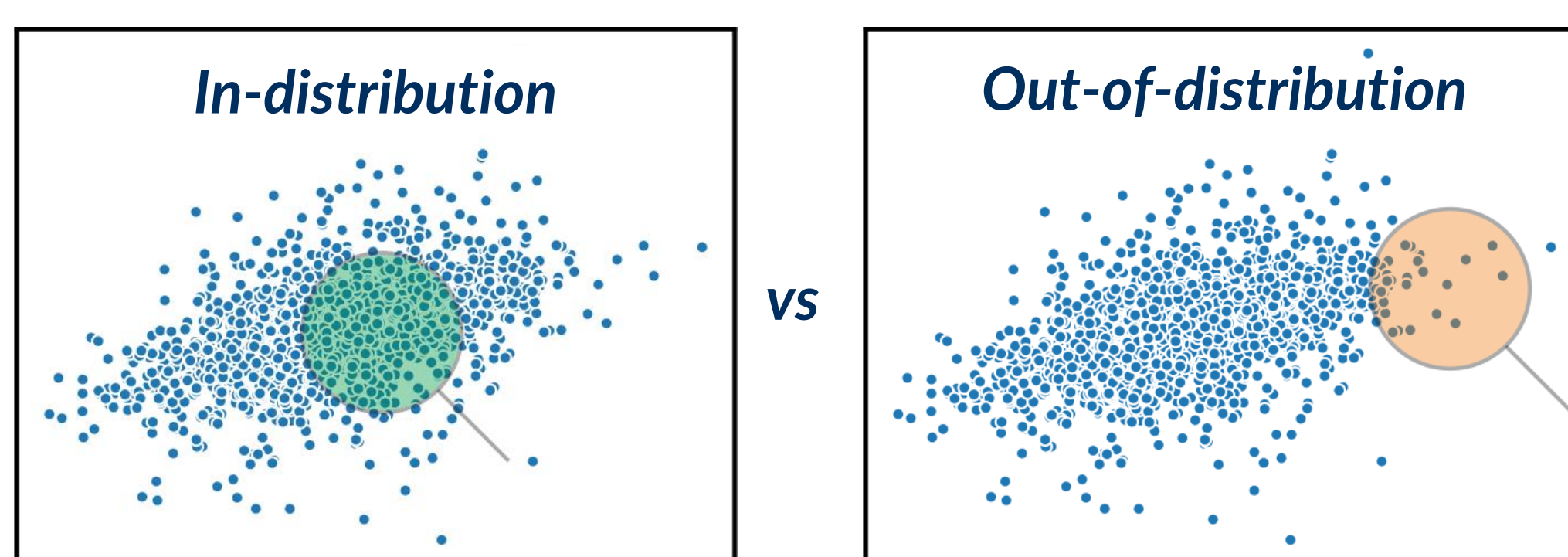Henrik Hupatz, Miguel Rivero Crespo, Berit Olofsson, Anneli Kruve

henrik.hupatz@mmk.su.se
Stockholm University Center of Circular and Sustainable Systems (SUCCeSS), Svante Arrhenius väg 16C, 114 16 Stockholm, Sweden

## BACKGROUND

- **Accumulation** of polar **pesticides** in the environment is **threatening** water quality.
- **Mass spectrometry** (MS) coupled with **chromatography** is a sensitive analytical chemistry method for a wide range of analytes.
- **Derivatization** facilitates the analysis of challenging **highly polar** small molecules.
- **Generative modeling** enables the **inverse molecular design** of easily **accessible** derivatizing reagents with **desired** chemical **properties**.
- **In-silico generated** reagents require **exploration** of **unknown chemical space** to predict chemical properties, such as ionization efficiency (*IE*).
- **IE prediction model** can be **optimized** for **out-of-distribution (OOD) data**.

**6 CLEAN WATER AND SANITATION**    **14 LIFE BELOW WATER**

*In-distribution*   vs   *Out-of-distribution*

**PROTOTYPE**

Analytical method · Contamination of soil and plants · Food · Water · Toxic · Health problems · Monitoring needed · Sample preparation

**Glyphosate**

**highly polar zwitterionic**

## COMPUTATIONAL METHODS

- 658 models were trained on the **1147 chemicals** of an *IE* dataset to predict log*IE*.

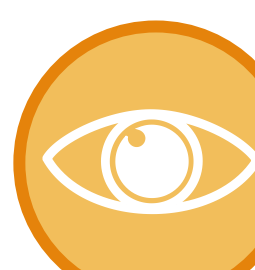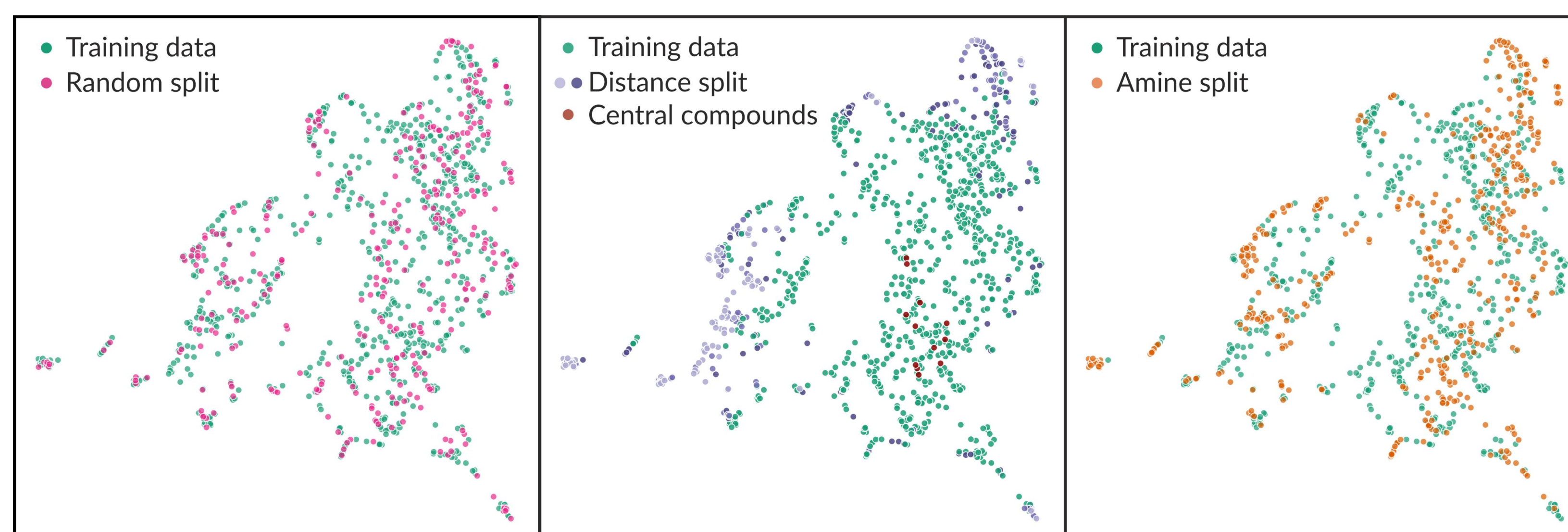**Input Features** [4]   **Train/Test Splitting** [3]   **Feature Cleaning** [9]   **Algorithms** [6]

## CHEMICAL SPACE VISUALIZATION AND DATA SPLITTING

- Training data · Random split
- Training data · Distance split · Central compounds
- Training data · Amine split
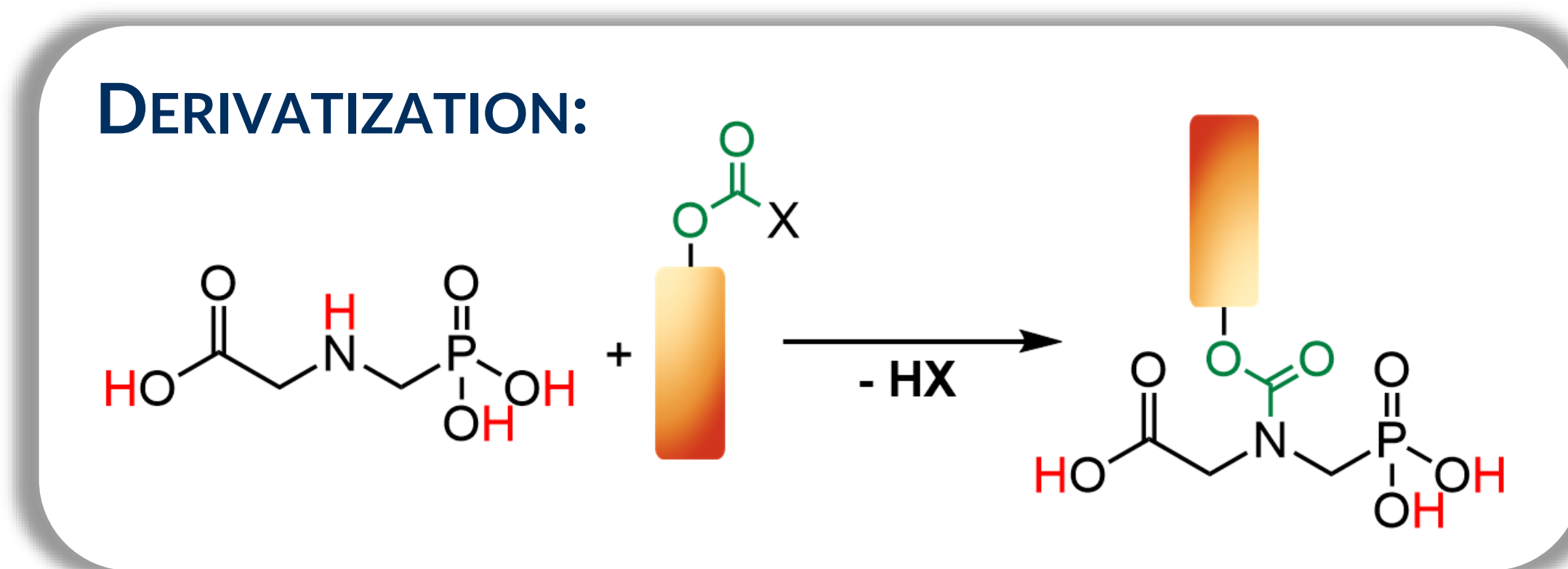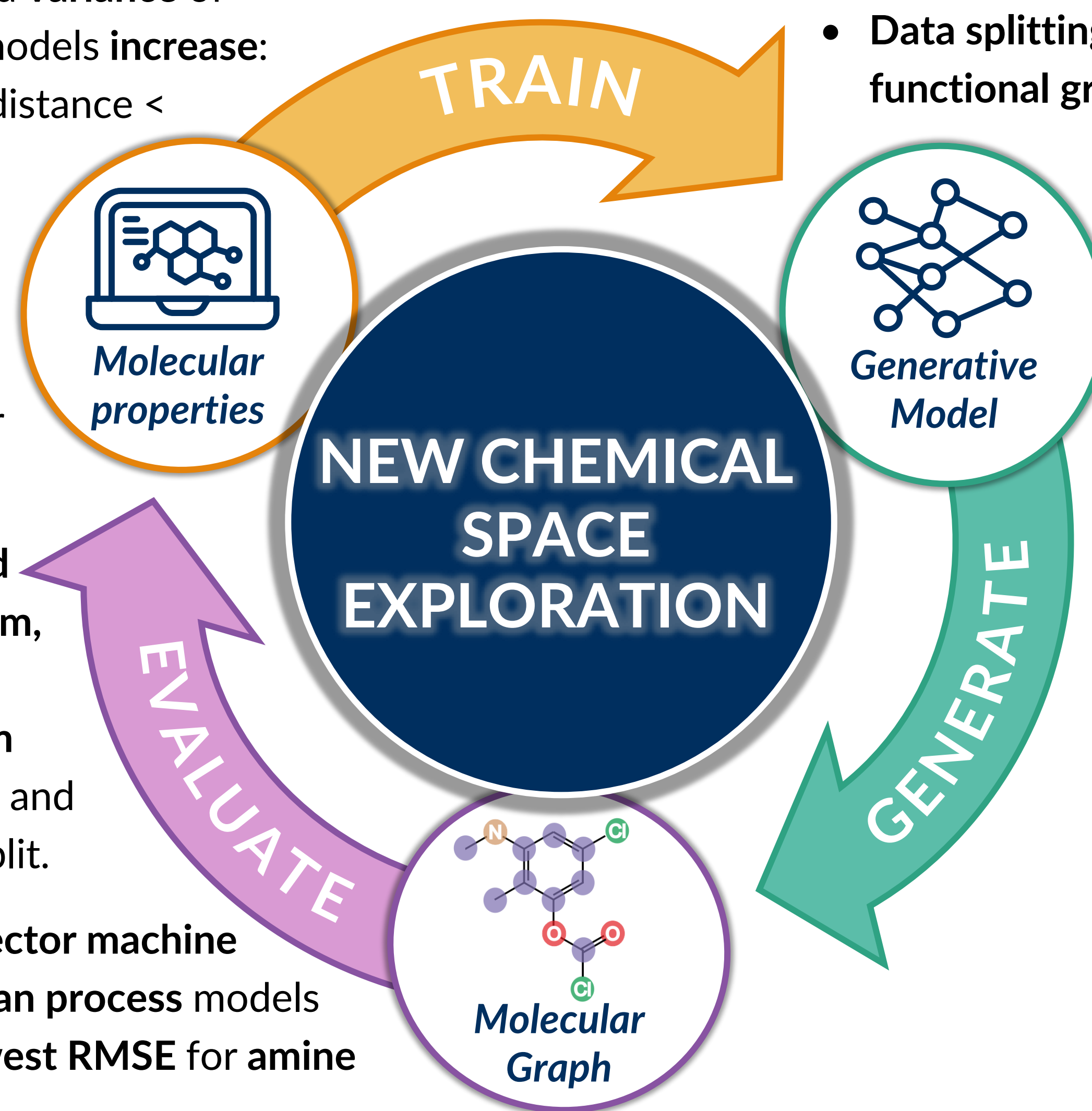
- **Myopic common edge subgraph distances** [1] (MCESD) were used to characterize **chemical similarity**.
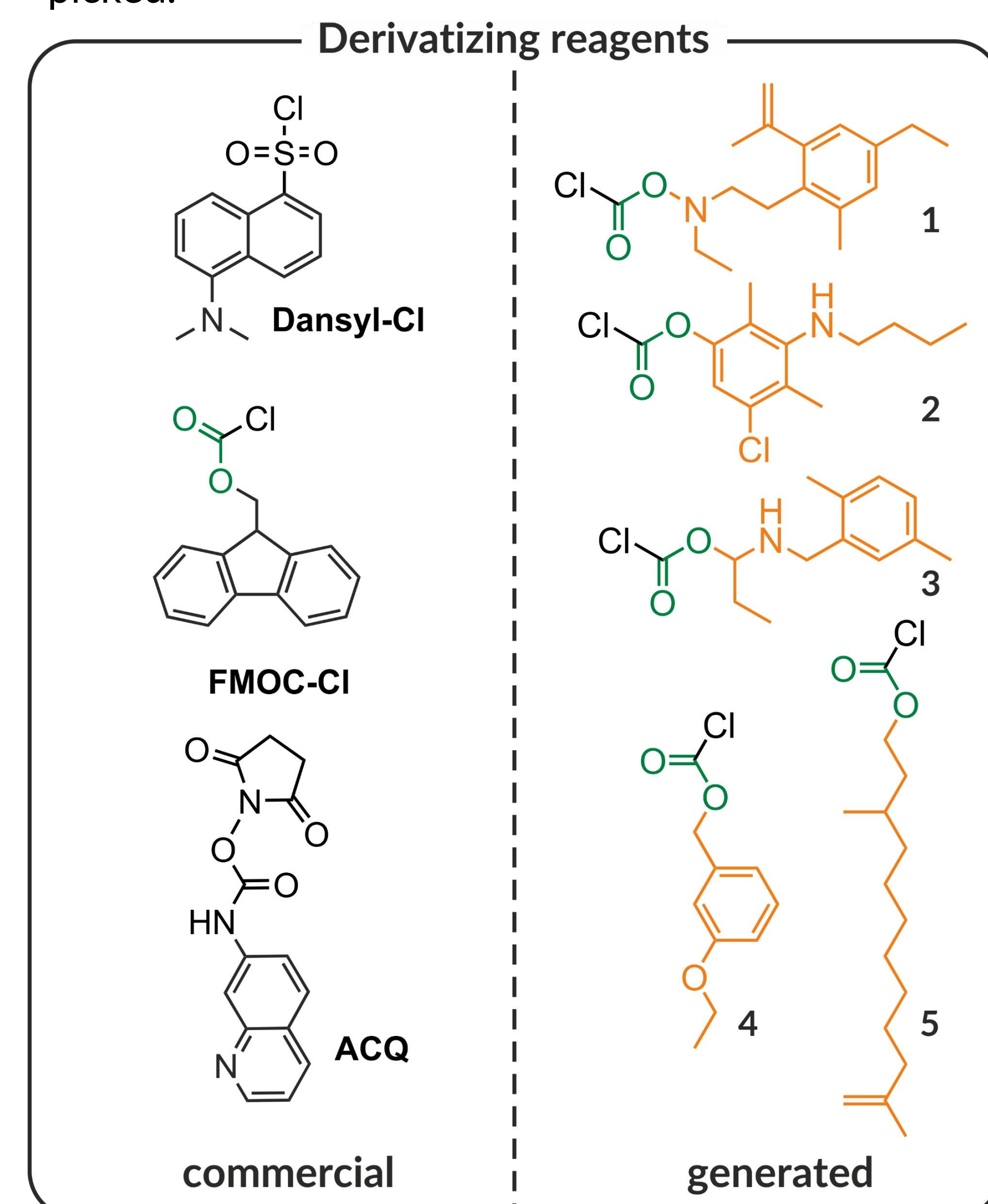- **UMAP reduction** to visualize chemical space of IE dataset.

- Most **characteristic** ("central") **compounds** were identified by **lowest average MCESD**.
- **Data splitting** based on **MCESD** to central compounds and most **abundant functional group** (amine) to investigate **OOD performance** of models.
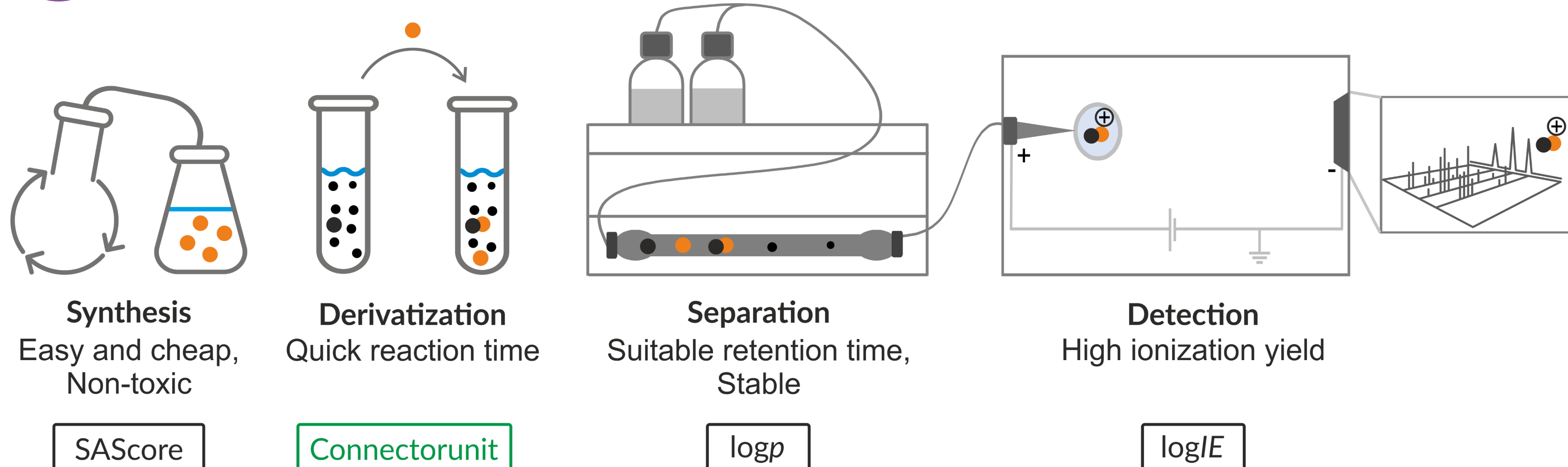
## ID VS OOD PERFORMANCE

**Features**
- ECFP6
- Mordred
- PaDEL
- SIRIUS FP

RMSE Test

**Algorithms**
- gaussprPoly
- gbm
- ranger
- svmPoly
- xgbLinear
- xgbTree

RMSE Test

Amine split · Distance split · Random split

- **Median** and **variance** of **RMSE of models increase**: random < distance < amine
- **Mordred** descriptors perform **better** over all **splits**.
- **Tree-based** models (**gbm, xgbTree**) **outperform** on **random** and **distance** split.
- **Support vector machine** and **gaussian process** models exhibit **lowest RMSE for amine** split.

## NEW CHEMICAL SPACE EXPLORATION

**TRAIN** → *Generative Model* → **GENERATE** → *Molecular Graph* → **EVALUATE** → *Molecular properties* →

## MONTE-CARLO TREE SEARCH

- **Molecular graphs** generated by MCTS from **carbon seed** to molecular weight of **400** with MCTS **depth of 2** and **width of 12**.
- **Evaluation function:** $Ev = \log IE - SAscore$.
- For each MCTS the structure with highest *Ev* was picked.

**Derivatizing reagents**

Dansyl-Cl · FMOC-Cl · ACQ

**commercial**   |   **generated** (1, 2, 3, 4, 5)

## EXPERIMENTAL PROPERTIES

**DERIVATIZATION:**

$- HX$

**Synthesis**
Easy and cheap, Non-toxic
SAScore

**Derivatization**
Quick reaction time
Connectorunit

**Separation**
Suitable retention time, Stable
log*p*

**Detection**
High ionization yield
log*IE*
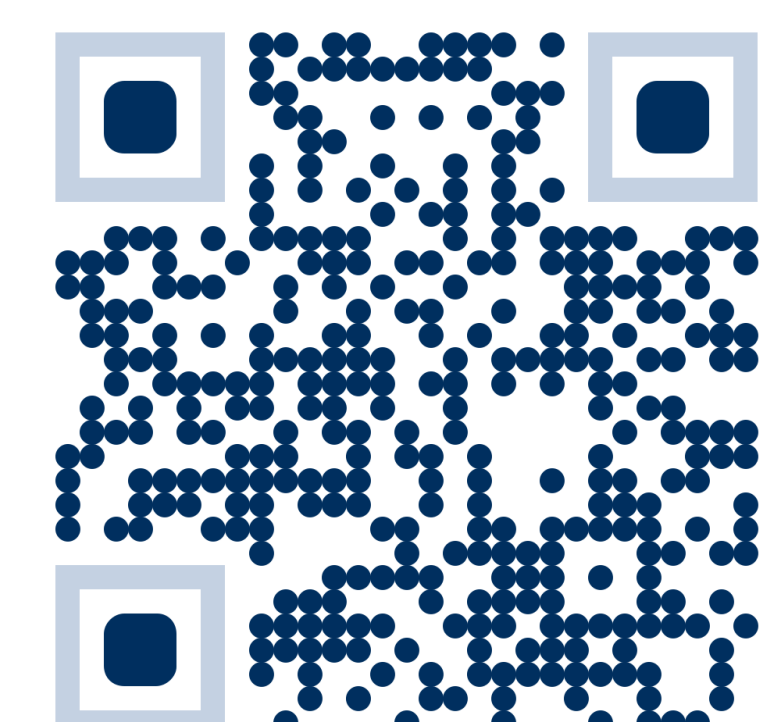
[1]For references, additional information and a poster copy:

*Kruve lab*

Stockholm University · SUCCeSS