

# AutoPeptideML 2: An open source library for democratizing machine learning for peptide bioactivity prediction

Raúl Fernández-Díaz\*, Thanh Lam Hoang, Vanessa Lopez, Denis Shields

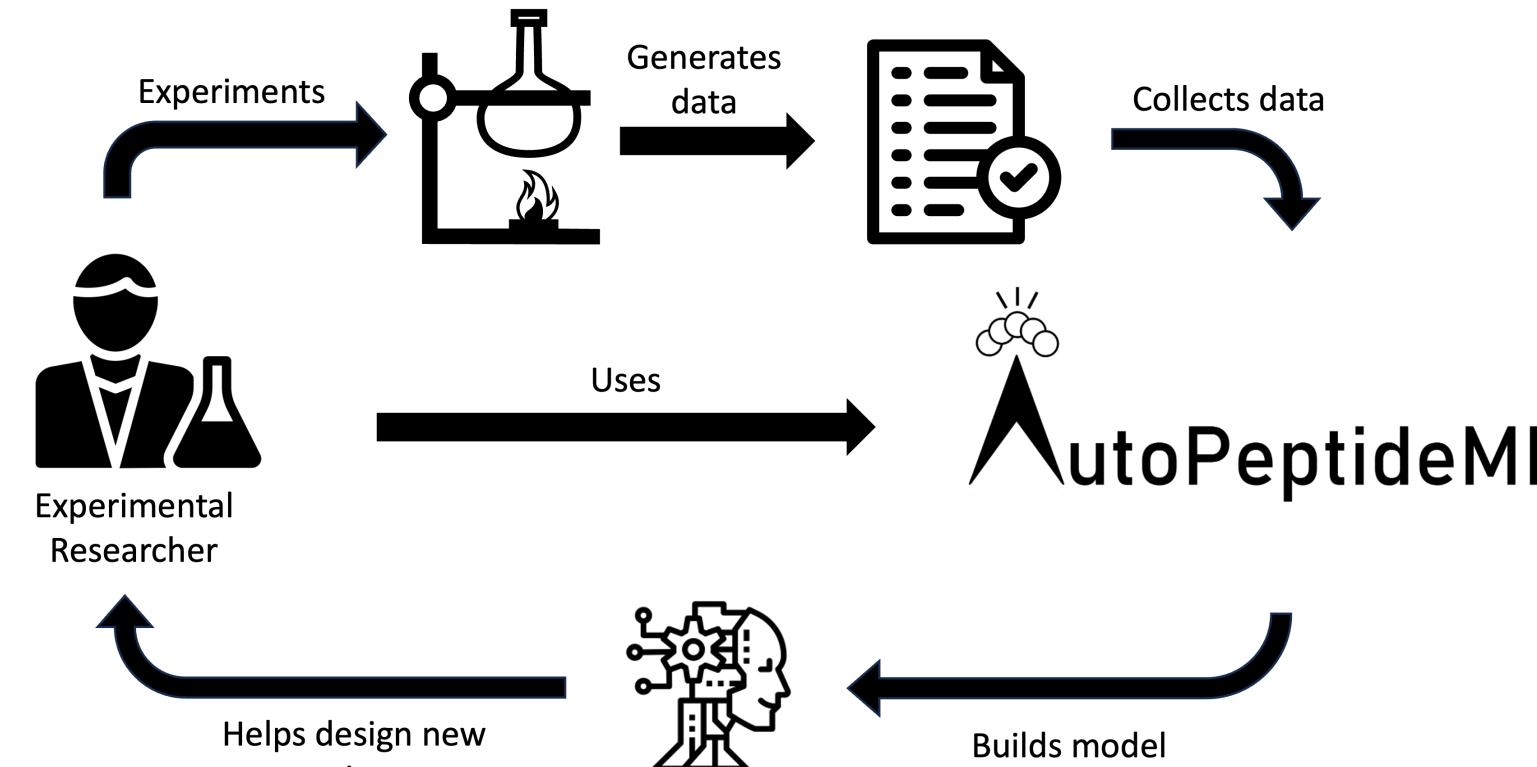
## Introduction

### Motivation

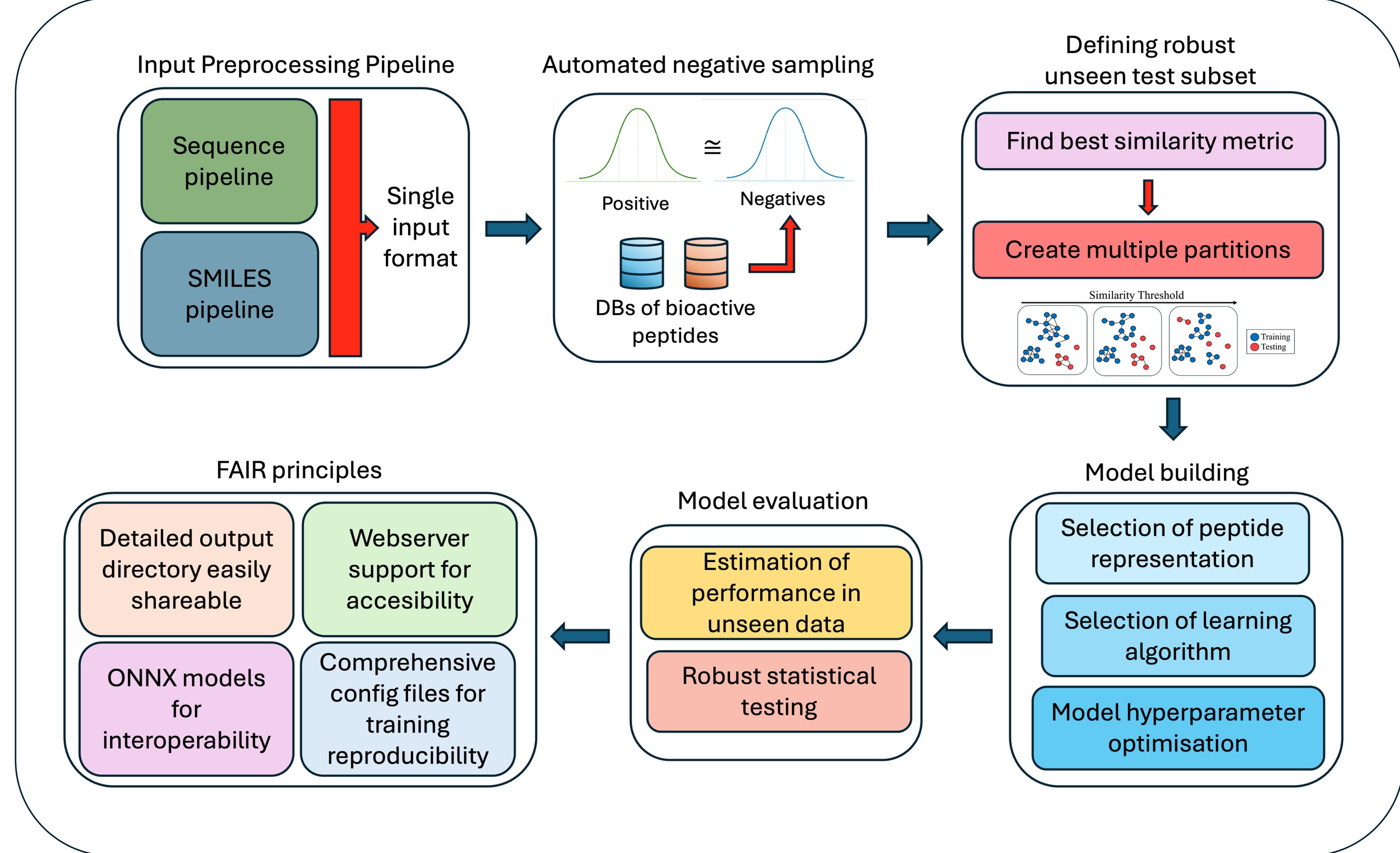
Integration of machine learning modelling to experimental workflows in peptide drug discovery/design can accelerate research.

### Objectives

1. Empower experimental researchers to build custom models.
2. Generate output that follows FAIR principles.
3. Handle canonical, non-canonical peptides, and peptidomimetic molecules.

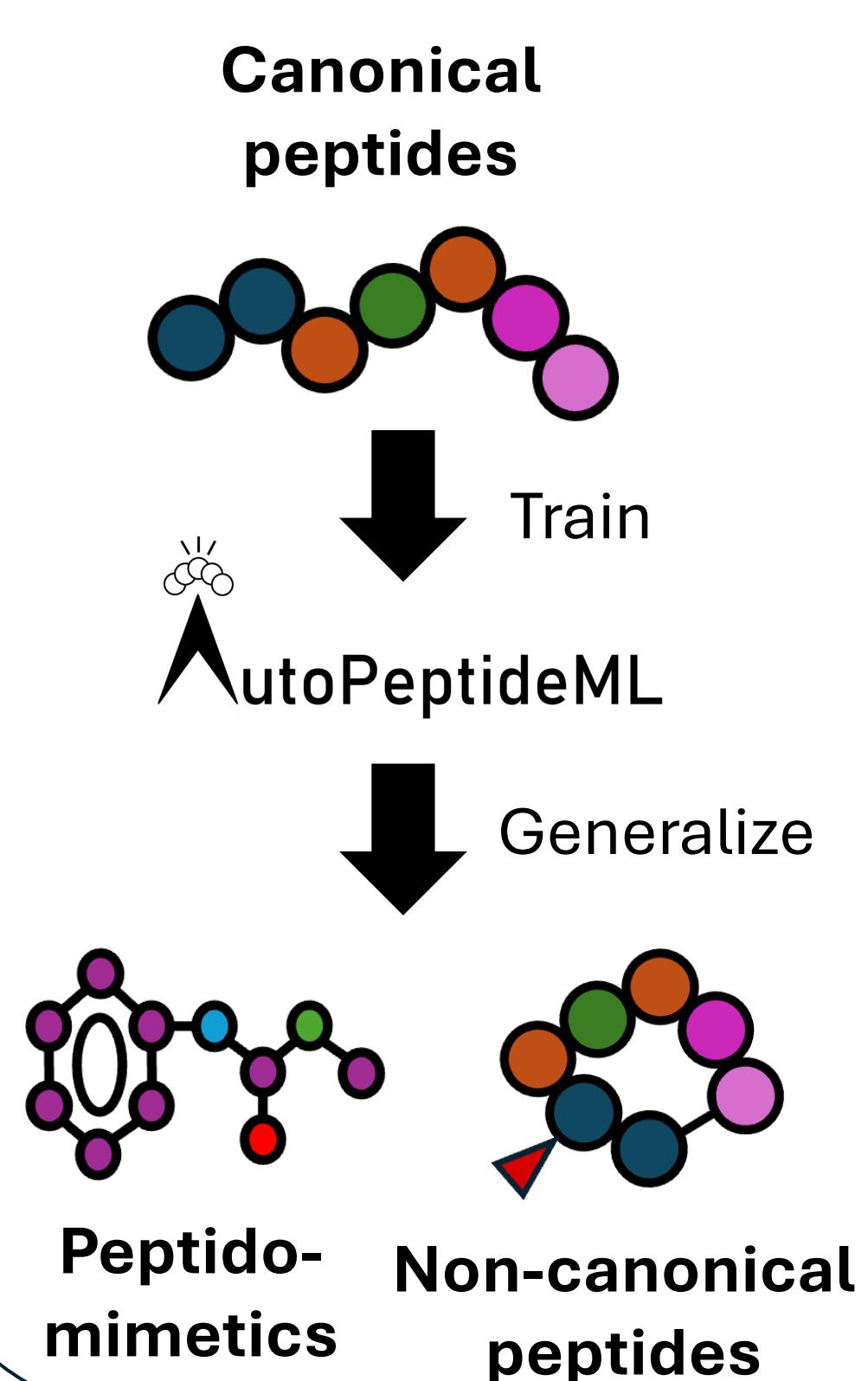


## Modules

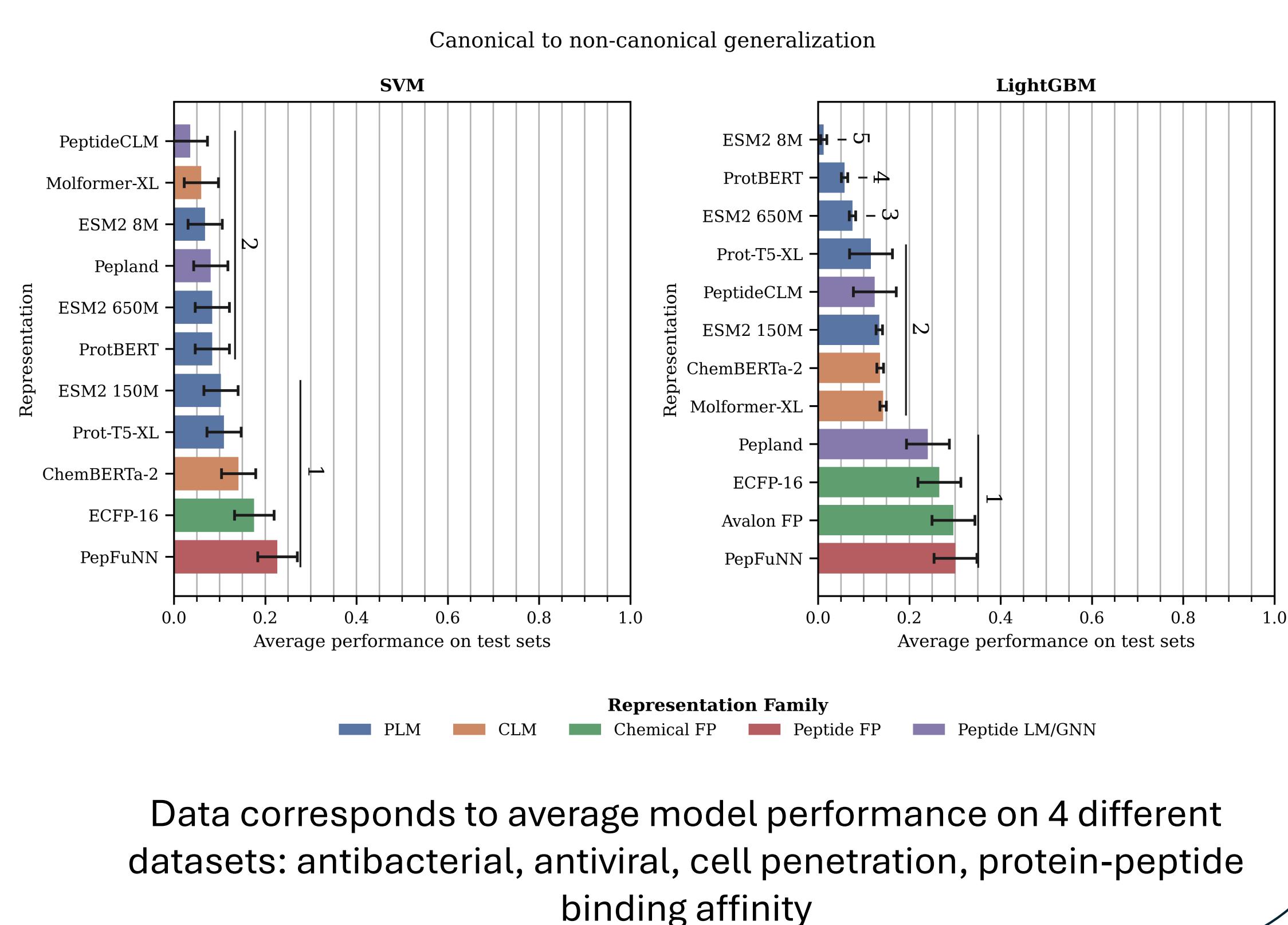


## Canonical to non-canonical generalization

### Overview



### Results



## Conclusions

A. AutoPeptideML is a flexible tool that makes model-building easier and handles:

1. Data preprocessing
2. Negative sampling
3. Dataset partitioning
4. Model and representation selection
5. Model hyperparameter optimization
6. Model evaluation
7. Training and evaluation PDF report

B. Its design facilitates compliance with FAIR principles and DOME guidelines.

- C. Handle canonical, non-canonical peptides, and peptidomimetic small molecules.
- D. Facilitates development of new methods of peptide representation by providing robust and reproducible workflows and evaluation

## Dependencies

### Data preparation



### Model training



### Evaluation + Report



### Webserver



### Webserver



### AutoPeptideML Model builder

Welcome to the AutoPeptideML webserver.

The next steps will help you build your own model.

#### 0. Modelling task

First, start by defining the modelling task.

What is the modeling problem you are facing?

Classification (returning categorical values)

How do you want to define your peptide?

Macromolecules - allows for canonical, non-canonical, and peptidomimetics

Next step

<http://peptide.ucd.ie/autopeptideml>

### Python code

```

1 import pandas as pd
2 from autopeptideml import AutoPeptideML
3
4 # Load data
5 df = pd.read_csv('antibacterial_data_canonical.csv')
6 df2 = pd.read_csv('antibacterial_data_noncanonical.csv')
7 all_inputs = df['sequence'].tolist() + df['SMILES'].tolist()
8
9 # Initialise AutoPeptideML
10 apml = AutoPeptideML(
11     data=all_inputs,
12     outputdir='demo'
13 )
14
15 # Preprocess
16 apml.preprocess_data(
17     pipeline='to-smiles',
18     n_jobs=5,
19     verbose=True
20 )
21
22 # Build models
23 apml.build_models(split_strategy='min',
24     task='class',
25     reps=['chemberta-2', 'ecfp',
26           'esm2-8m',
27           'peptideclm'],
28     device='mps',
29     n_trials=200)
  
```

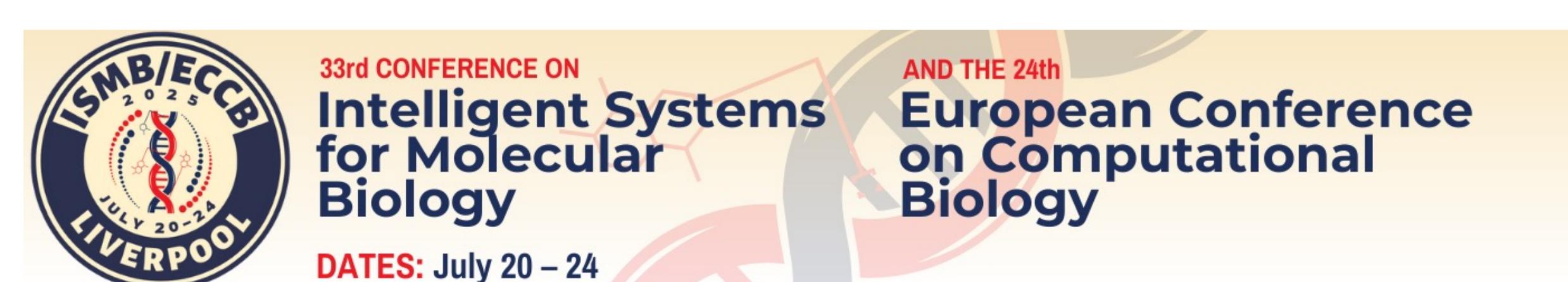
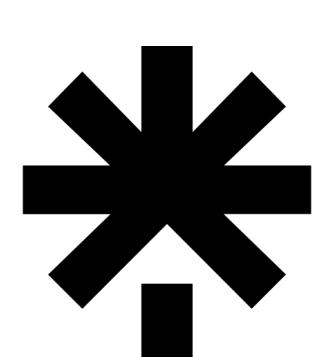
### PDF report

#### AutoPeptideML - Model Development report

##### Table of contents

1	Introduction	1
2	Guidance on the interpretation of model performance metrics	2
2.1	Classification	2
2.2	Regression	2
3	Metadata	3
4	Model and representation selection	4
4.1	Optimization history	4
4.2	Model vs peptide representation	5
4.3	Model vs model	6
5	Evaluation metrics - independent hold-out test set	7
5.1	Metrics	7
5.2	Calibration curve	7
5.3	ROC curve	8
5.4	Precision-recall curve	9

## Contact info:



## Github Repository:

