











Speaker: Raúl Fernández-Díaz (PhD Candidate UCD – IBM Research)

UCD: R. Cossio-Pérez, C. Agoni, D.C. Shields

IBM Research: T.L. Hoang, V. Lopez

Novo Nordisk: R. Ochoa







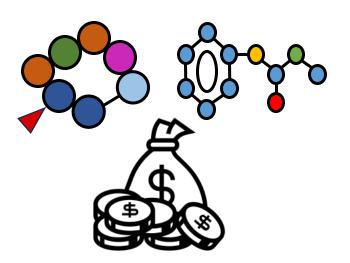
# A tale of two peptides

- Protein sequences (20 aa)
- Cheap
- Not great drugs



## Natural (or canonical) peptides Synthetic peptides (or non-canonical) and peptidomimetics

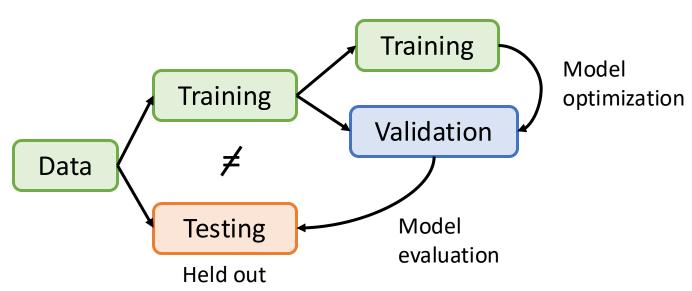
- Modified peptides + small molecules
- Expensive
- Better drug candidates



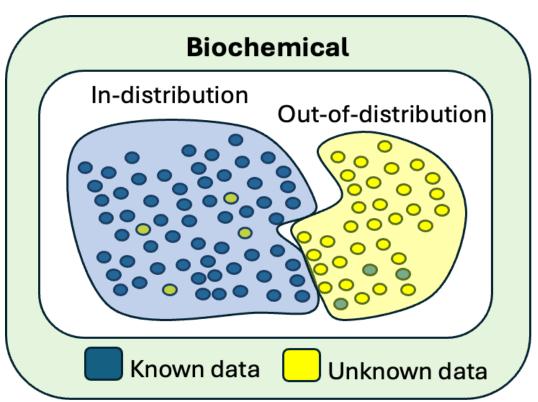


# **Dataset partitioning**

Central Assumption of ML: "Training data is representative of prediction data"



Training and testing need to have different molecules, otherwise we cannot evaluate generalization/extrapolation

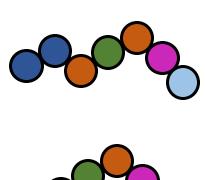


**Objective:** Estimate model performance in unseen data



# Peptide representation

We need to transform molecules into numerical vectors



#### **Methods**

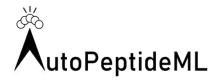


[1, 0.3, 0.4, 2.1, ...]

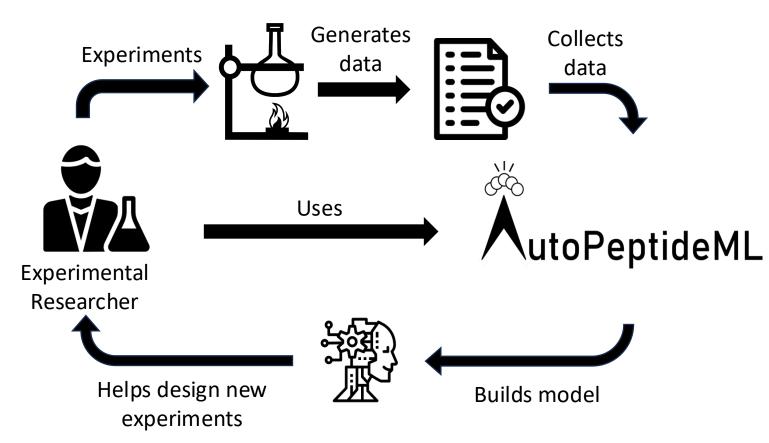
- **Heuristic**: Molecular FPs
  - Substructure-based: ECFP/Avalon
  - Monomer sequence: PepFuNN
- **Data-driven**: Pre-trained representation models
  - [Protein/Chemical/Peptide] Language Models
  - GNNs



# **Automation**



R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555



## **Design Requirements**

- 1. Robust evaluation
- 2. Reproducibility
- 3. Easier to integrate with experimental workflow



# **Objectives**

- 1. How to automatically build peptide property prediction models (and how to evaluate them)
- 2. How to extrapolate from natural to synthetic peptides or peptidomimetics





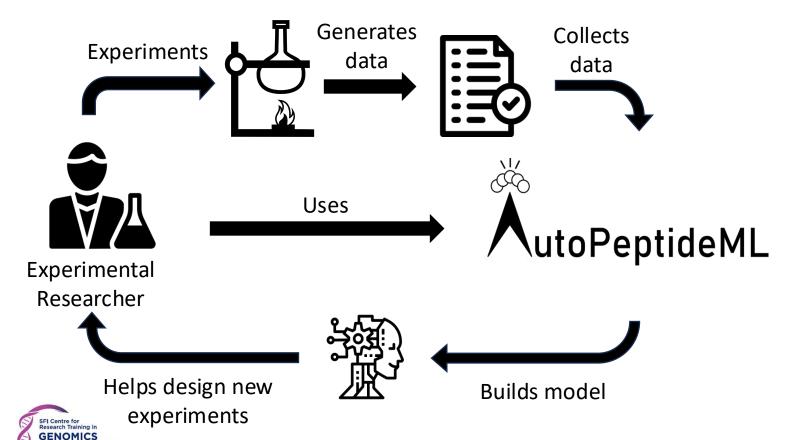




# **Automation**



R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555



## **Design Requirements**

- 1. Competitive performance
- 2. Reliable evaluation so that experimental scientist can trust the models
- 3. Easy to use



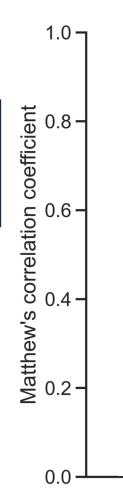


# How good can automation really be?

Collected 18
datasets used for
building different
peptide bioactivity
predictors



R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555



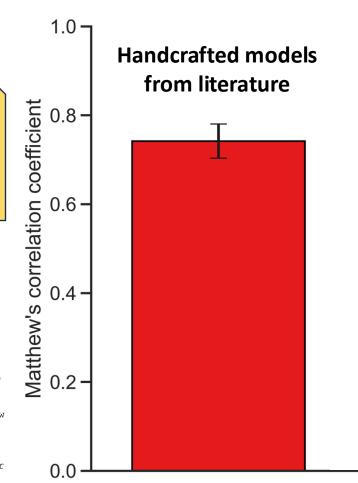


# How good can automation really be?

Collected 18
datasets used for
building different
peptide bioactivity
predictors



R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555



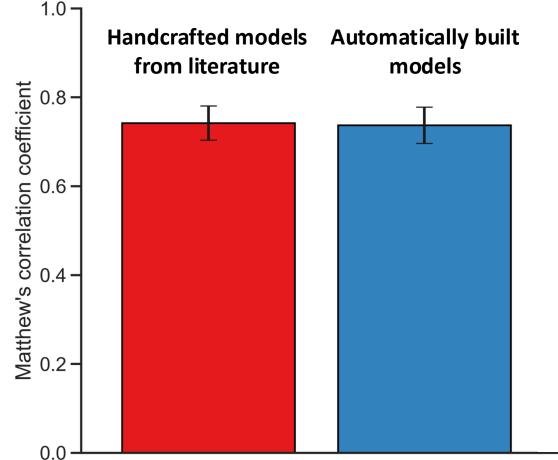


# Automation is as good as human developers

Collected 18
datasets used for
building different
peptide bioactivity
predictors

toPeptideML
Fernández-Díaz et al.,

R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555



learn

Low intensity computing



Bayesian Optimization for hyperparameter selection

Protein Language Models

General representation/ featurization



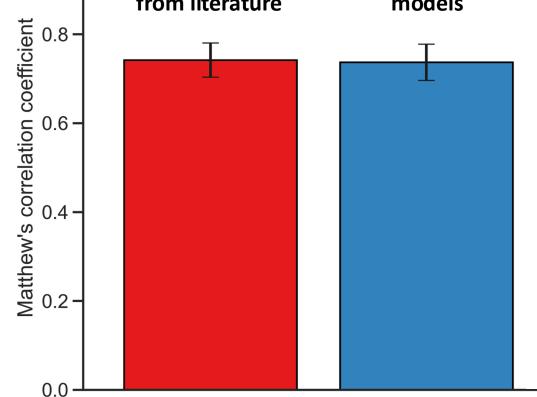
# How important is the choice of negative samples?

Collected 18
datasets used for
building different
peptide bioactivity
predictors

# **A**utoPeptideML

R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555





#### **Before:**

- Random peptides
- Peptides from Uniprot
- Protein fragments
- Scrambled sequences

**Problem:** They might differ from positives due to confounding biophysical characteristics (e.g. solubility, membrane crossing, etc.)

#### Now:

Other bioactive peptides

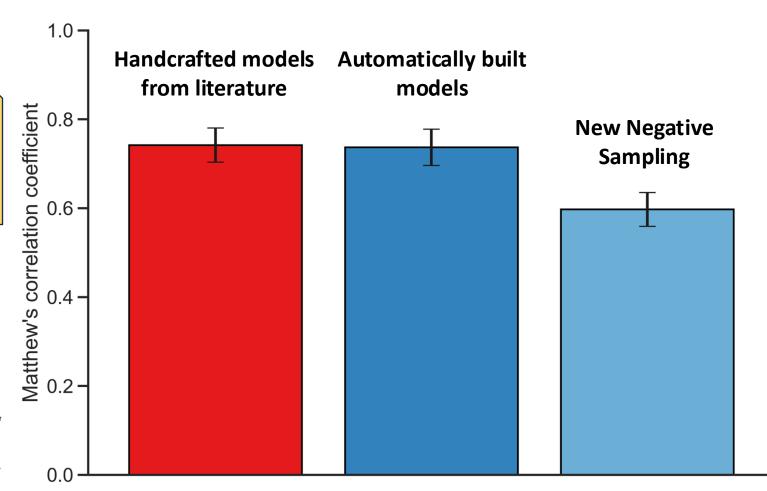
**Peptipedia** 



# Proper choice of negatives is quite important

Collected 18
datasets used for
building different
peptide bioactivity
predictors

R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555



#### **Before:**

- Random peptides
- Peptides from Uniprot
- Protein fragments
- Scrambled sequences

**Problem:** They might differ from positives due to confounding biophysical characteristics (e.g. solubility, membrane crossing, etc.)

#### Now:

Other bioactive peptides

**Peptipedia** 



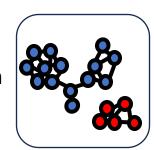
# How important is to have independent test sets?

Collected 18
datasets used for
building different
peptide bioactivity
predictors

1.0 -

Handcrafted models Automatically built from literature models coefficient 0.8 -**New Negative** Sampling 0.6 Matthew's correlation 0.0

Testing data should only include unseen molecules not similar to training (sequence alignment)



utoPeptideML

R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555

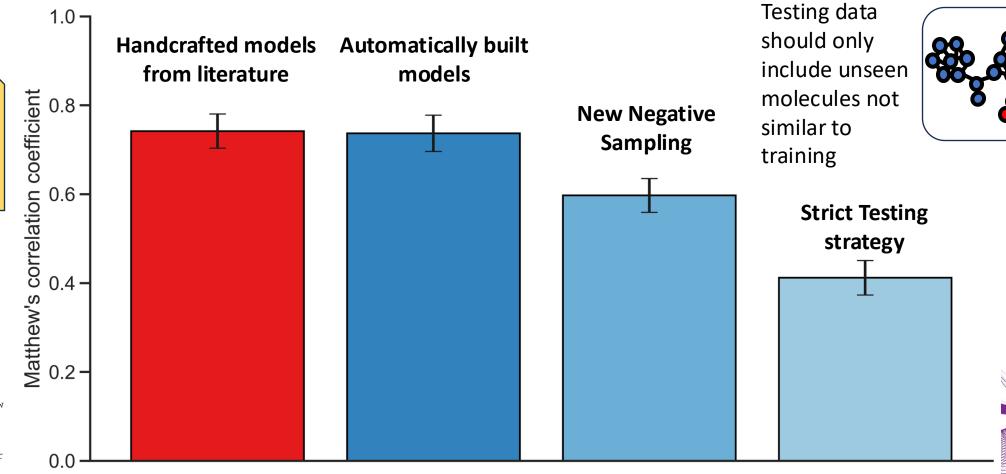


# Independence in test sets is crucial

Collected 18
datasets used for
building different
peptide bioactivity
predictors



R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555





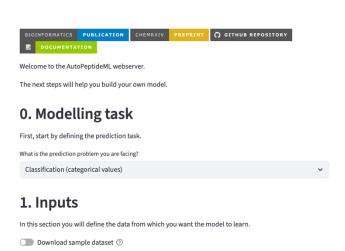
# Automating ML for natural peptides



R. Fernández-Díaz et al.,
AutoPeptideML: a study on how
to build more trustworthy
peptide bioactivity
predictors, Bioinformatics,
Volume 40, Issue 9, September
2024, btae555

#### Webserver - GUI

# utoPeptideML



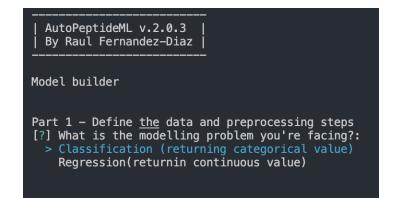
Browse files

Please upload dataset with your peptides and their labels if available

Drag and drop file here

Limit 200MB per file

#### **CLI tool**



#### **Python Package**

```
df = pd.read_csv(osp.join(PATH, 'original_data', f'c-{dataset}.csv'))
apml = AutoPeptideML(
    data=df,
    outputdir=f'apml-{dataset}',
    sequence_field='SMILES',
    label_field='labels'
)
apml.build_models(
    task='class',
    reps=['esm2-8m', 'peptideclm', 'chemberta-2', 'ecfp-16'],
    models=['svm', 'knn', 'rf', 'lightgbm', 'xgboost'],
    device='mps',
    n_trials=10
)
apml.create_report()
return apml
```



# Part 2 – Peptide representation



# Can we leverage data on cheaper experiments to prioritise more expensive experiments?



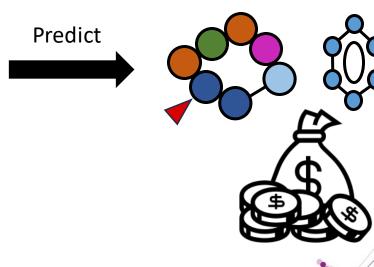
Fernández-Díaz R, et al. How to build machine learning models able to extrapolate from standard to modified peptides. ChemRxiv. 2025; doi:10.26434/chemrxiv-2025ggp8n-v3

Natural peptides (cheap)





Synthetic peptides and peptidomimetics (expensive but better drugs)





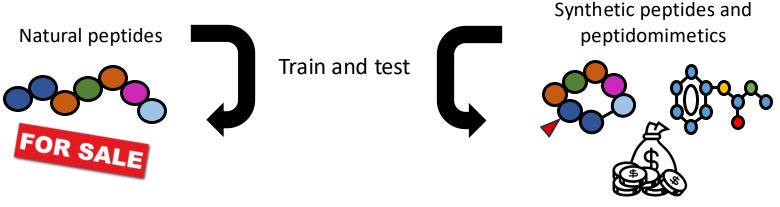


# **Computational experiments**

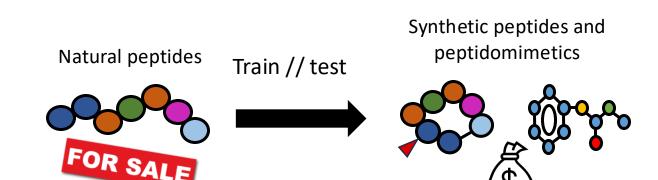
## **Interpolation**



Fernández-Díaz R, et al. How to build machine learning models able to extrapolate from standard to modified peptides. ChemRxiv. 2025; doi:10.26434/chemrxiv-2025-ggp8n-v3



## **Extrapolation**





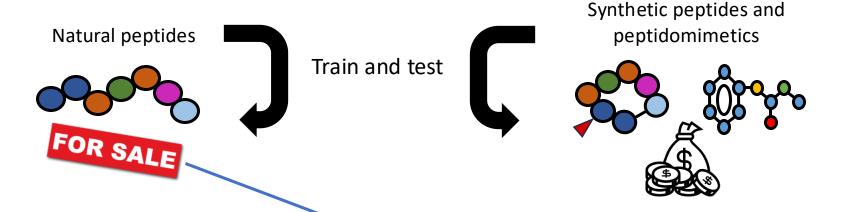


# **Computational experiments**

## **Interpolation**

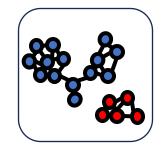


ternandez-Diaz R, et al. How to build machine learning models able to extrapolate from standard to modified peptides. ChemRxiv. 2025; doi:10.26434/chemrxiv-2025ggp8n-v3



How to build train/test splits: Should we use sequence alignment for measuring peptide similarity?

Similarity-informed train/test split







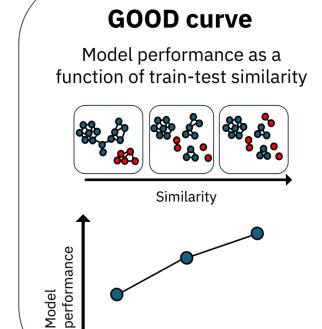


# Finding the best similarity metric



Fernandez-Diaz R, et al. A new framework for evaluating model out-of-distribution generalisation for the biochemical domain. InThe Thirteenth International Conference on Learning Representations 2025.

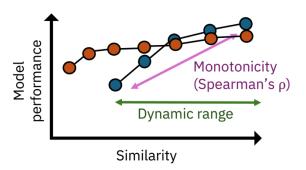
#### **Hestia-GOOD framework**



Similarity

#### Similarity metric selection

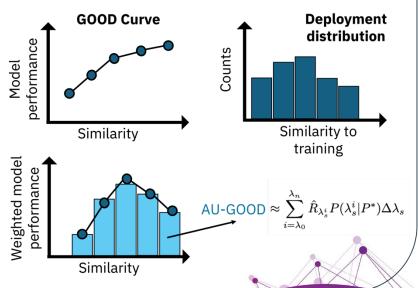
Quantitative analysis of best similarity function for a given task/dataset



- Monotonicity: Is model performance a function of train-test similarity?
- Dynamic range: What is the resolution of the similarity metric?

#### **AU-GOOD** metric

Estimation of model performance conditioned on a deployment distribution







#### Analysis of 8 datasets:

- 4 natural
- 4 synthetic
- Chemical FPs are better than sequence alignment for natural peptides

## Best metrics for each dataset

#### **Metrics explored:**

- ECFP various radii
- MAPc various radii
- MMSeqs2 (alignment) -
- Needleman-Wunsch (alignment)
- ESM2-8M embedding distance
  - Molformer-XL embedding distance





Fernandez-Diaz R, et al. A new framework for evaluating model out-of-distribution generalisation for the biochemical domain. In The Thirteenth International Conference on Learning Representations 2025.





## **Best metrics for each dataset**



Fernandez-Diaz R, et al. A new framework for evaluating model out-of-distribution generalisation for the biochemical domain. InThe Thirteenth International Conference on Learning Representations 2025.

#### **Metrics explored:**

- ECFP various radii
- Needleman-Wunsch (alignment)
- MAPc various radii
- ESM2-8M embedding distance
- MMSeqs2 (alignment) -
  - Molformer-XL embedding distance

- Analysis of 8 datasets:
  - 4 natural
  - 4 synthetic
- Chemical FPs are better than sequence alignment for natural peptides

| Dataset                             | Peptide type | Task           | Similarity Type    | Similarity | Dynamic range $(\uparrow)$ [a] | $ \begin{array}{c} \mathbf{Monotonicity} \\ (\uparrow) \ [\mathbf{b}] \end{array} $ |
|-------------------------------------|--------------|----------------|--------------------|------------|--------------------------------|---|
| Protein-peptide<br>binding affinity | Standard     | Regression     | Chemical FP        | MAPc-8     | 70 %                           | $0.8 \pm 0.1$   |
| Protein-peptide<br>binding affinity | Modified     | Regression     | Chemical FP        | MAPc-20    | 80 %                           | $0.95 \pm 0.03$   |
| Cell penetration                    | Standard     | Classification | Chemical FP        | MAPc-8     | 60 %                           | $0.98 \pm 0.04$   |
| Cell penetration                    | Modified     | Classification | Chemical FP        | MAPc-12    | 60 %                           | $0.5 \pm 0.2$   |
| Antibacterial                       | Standard     | Classification | Chemical FP        | MAPc-8     | 60 %                           | $0.97 \pm 0.02$   |
| Antibacterial                       | Modified     | Classification | Chemical FP        | ECFP-12    | 50 %                           | $0.9 \pm 0.1$   |
| Antiviral                           | Standard     | Classification | Sequence Alignment | MMSeqs2    | 80 %                           | $0.96 \pm 0.05$   |
| Antiviral                           | Modified     | Classification | Chemical FP        | MAPc-12    | 70 %                           | $0.6 \pm 0.2$   |





## **Best metrics for each dataset**



Fernandez-Diaz R, et al. A new framework for evaluating model out-of-distribution generalisation for the biochemical domain. InThe Thirteenth International Conference on Learning Representations 2025.

#### **Metrics explored:**

- ECFP various radii
- Needleman-Wunsch (alignment)
- MAPc various radii
- ESM2-8M embedding distance
- MMSeqs2 (alignment) -
  - Molformer-XL embedding distance

- Analysis of 8 datasets:
  - 4 natural
  - 4 synthetic
- Chemical FPs are better than sequence alignment for natural peptides

| Dataset                             | Peptide type | Task           | Similarity Type    | Similarity | Dynamic range $(\uparrow)$ [a] | $ \begin{array}{ll} \textbf{Monotonicity} \\ (\uparrow) \ [b] \end{array}$ |
|-------------------------------------|--------------|----------------|--------------------|------------|--------------------------------|--|
| Protein-peptide<br>binding affinity | Standard     | Regression     | Chemical FP        | MAPc-8     | 70 %                           | $0.8 \pm 0.1$  |
| Protein-peptide<br>binding affinity | Modified     | Regression     | Chemical FP        | MAPc-20    | 80 %                           | $0.95 \pm 0.03$  |
| Cell penetration                    | Standard     | Classification | Chemical FP        | MAPc-8     | 60 %                           | $0.98 \pm 0.04$  |
| Cell penetration                    | Modified     | Classification | Chemical FP        | MAPc-12    | 60 %                           | $0.5 \pm 0.2$  |
| Antibacterial                       | Standard     | Classification | Chemical FP        | MAPc-8     | 60 %                           | $0.97 \pm 0.02$  |
| Antibacterial                       | Modified     | Classification | Chemical FP        | ECFP-12    | 50 %                           | $0.9 \pm 0.1$  |
| Antiviral                           | Standard     | Classification | Sequence Alignment | MMSeqs2    | 80 %                           | $0.96 \pm 0.05$  |
| Antiviral                           | Modified     | Classification | Chemical FP        | MAPc-12    | 70 %                           | $0.6 \pm 0.2$  |





## **Best metrics for each dataset**



Fernandez-Diaz R, et al. A new framework for evaluating model out-of-distribution generalisation for the biochemical domain. InThe Thirteenth International Conference on Learning Representations 2025.

#### **Metrics explored:**

- ECFP various radii
- Needleman-Wunsch (alignment)
- MAPc various radii
- ESM2-8M embedding distance
- MMSeqs2 (alignment) -
  - Molformer-XL embedding distance

- Analysis of 8 datasets:
  - 4 natural
  - 4 synthetic
- Chemical FPs are better than sequence alignment for natural peptides

| Dataset                             | Peptide type | Task           | Similarity Type    | Similarity | $\begin{array}{c} \textbf{Dynamic} \\ \textbf{range} \ (\uparrow) \ [\textbf{a}] \end{array}$ | $ \begin{array}{c} \mathbf{Monotonicity} \\ (\uparrow) \ [\mathbf{b}] \end{array} $ |
|-------------------------------------|--------------|----------------|--------------------|------------|---|---|
| Protein-peptide<br>binding affinity | Standard     | Regression     | Chemical FP        | MAPc-8     | 70 %  | $0.8 \pm 0.1$   |
| Protein-peptide<br>binding affinity | Modified     | Regression     | Chemical FP        | MAPc-20    | 80 %  | $0.95 \pm 0.03$   |
| Cell penetration                    | Standard     | Classification | Chemical FP        | MAPc-8     | 60~%  | $0.98 \pm 0.04$   |
| Cell penetration                    | Modified     | Classification | Chemical FP        | MAPc-12    | 60~%  | $0.5 \pm 0.2$   |
| Antibacterial                       | Standard     | Classification | Chemical FP        | MAPc-8     | 60~%  | $0.97 \pm 0.02$   |
| Antibacterial                       | Modified     | Classification | Chemical FP        | ECFP-12    | 50 %  | $0.9 \pm 0.1$   |
| Antiviral                           | Standard     | Classification | Sequence Alignment | MMSeqs2    | 80 %  | $0.96 \pm 0.05$   |
| Antiviral                           | Modified     | Classification | Chemical FP        | MAPc-12    | 70 %  | $0.6 \pm 0.2$   |





Average on 4

Statistical analysis are

Kruskal-Wallis with

post-hoc Wilcoxon

Significance defined

with Bonferroni

correction.

datasets.

test.

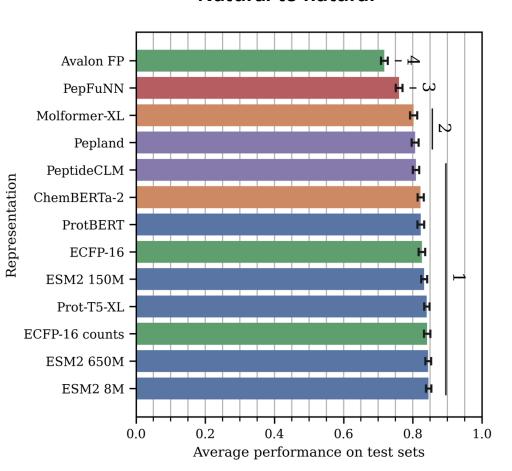
# Interpolation experiments



to build machine learning models able to extrapolate from standard to modified ChemRxiv. 2025;

Synthetic to synthetic 134/chemrxiv-2025-

#### Natural to natural





PLM

**CLM** 

**Representation Family** 

Chemical FP

Peptide FP





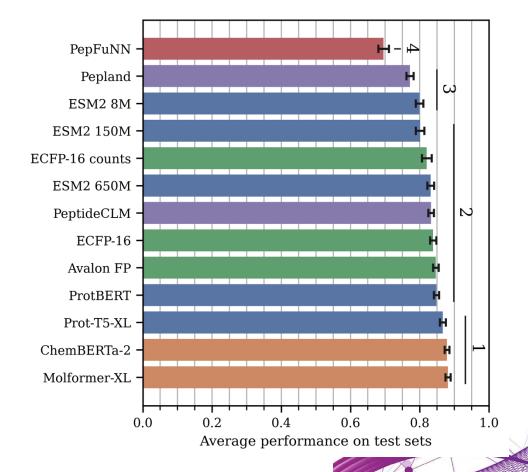
# Interpolation experiments



to build machine learning models able to extrapolate ChemRxiv. 2025;

#### Natural to natural Synthetic to synthetic 134/chemxxiv-2025-

- Average on 4 datasets.
- Statistical analysis are Kruskal-Wallis with post-hoc Wilcoxon test.
- Significance defined with Bonferroni correction.





Chemical FP

Peptide FP



# More information and contact info

# Interpolation experiments

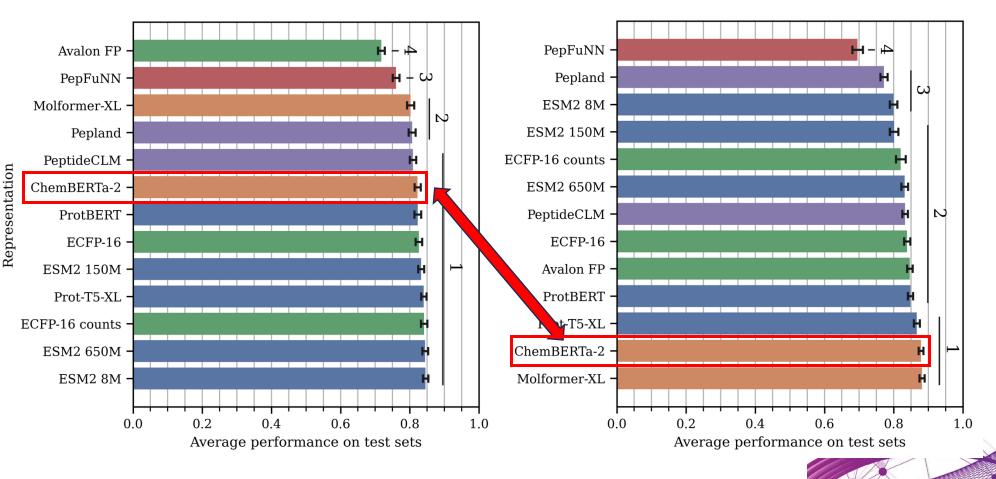
Natural to natural



to build machine learning models able to extrapolate

Synthetic to synthetic 134/chemrxiv-2025-

- Average on 4 datasets.
- Statistical analysis are Kruskal-Wallis with post-hoc Wilcoxon test.
- Significance defined with Bonferroni correction.





PLM

**CLM** 

**Representation Family** 

Chemical FP

Peptide FP

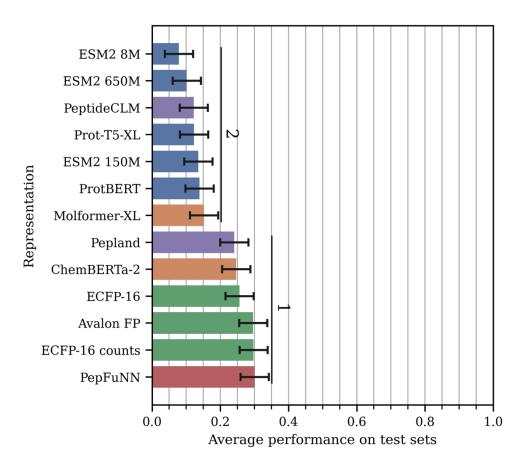


# Natural to synthetic extrapolation

## **Chem**Rxiv<sup>™</sup>

to build machine learning models able to extrapolate from standard to modified peptides. ChemRxiv. 2025; doi:10.26434/chemrxiv-2025-ggp8n-v3

- Average on 4 datasets.
- Statistical analysis are Kruskal-Wallis with post-hoc Wilcoxon test.
- Significance defined with Bonferroni correction.

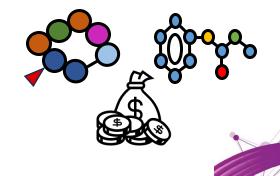




Natural peptides

Train // test

Synthetic peptides and peptidomimetics





**PLM** 

**CLM** 

**Representation Family** 

Chemical FP

Peptide FP

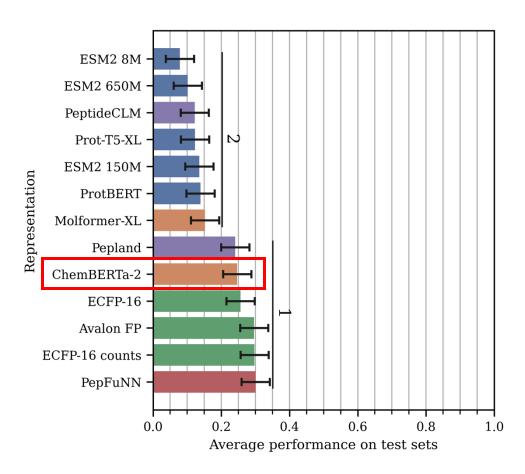


# Natural to synthetic extrapolation



fernandez-Diaz R, et al. How to build machine learning models able to extrapolate from standard to modified peptides. ChemRxiv. 2025; doi:10.26434/chemrxiv-2025ggp8n-v3

- Average on 4 datasets.
- Statistical analysis is ANOVA with post-hoc Tukey's HSD



#### **Conclusions**

- Natural to synthetic extrapolation is possible, but models are less reliable
- 2. ChemBERTa-2 appears to be the most versatile tool to work with peptides
- Chemical and Peptide
   Fingerprints are robust
   options as well



PLM

**CLM** 

#### **Representation Family**

Chemical FP

Peptide FP



## **Conclusions**

- Dataset building (negative definition) and partitioning (train/test split) are key for proper model evaluation
- 2. Chemical fingerprints are better for partitioning natural and synthetic datasets than sequence alignment, contrary to standard practice
- 3. Natural to synthetic extrapolation is possible, but there is room for improvement. We release the benchmark datasets to make it easier for the community to improve
- 4. ChemBERTa-2 appears to be the most versatile tool, closely followed by chemical fingerprints



Contact info, papers, and slides of the presentation

















# Partitioning, representation, and automation in canonical and non-canonical peptide modelling

Speaker: Raúl Fernández-Díaz (PhD Candidate UCD – IBM Research)

UCD: R. Cossio-Pérez, C. Agoni, D.C. Shields

IBM Research: T.L. Hoang, V. Lopez

Novo Nordisk: R. Ochoa

