

How to generalize machine learning models to both canonical and non-canonical peptides

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

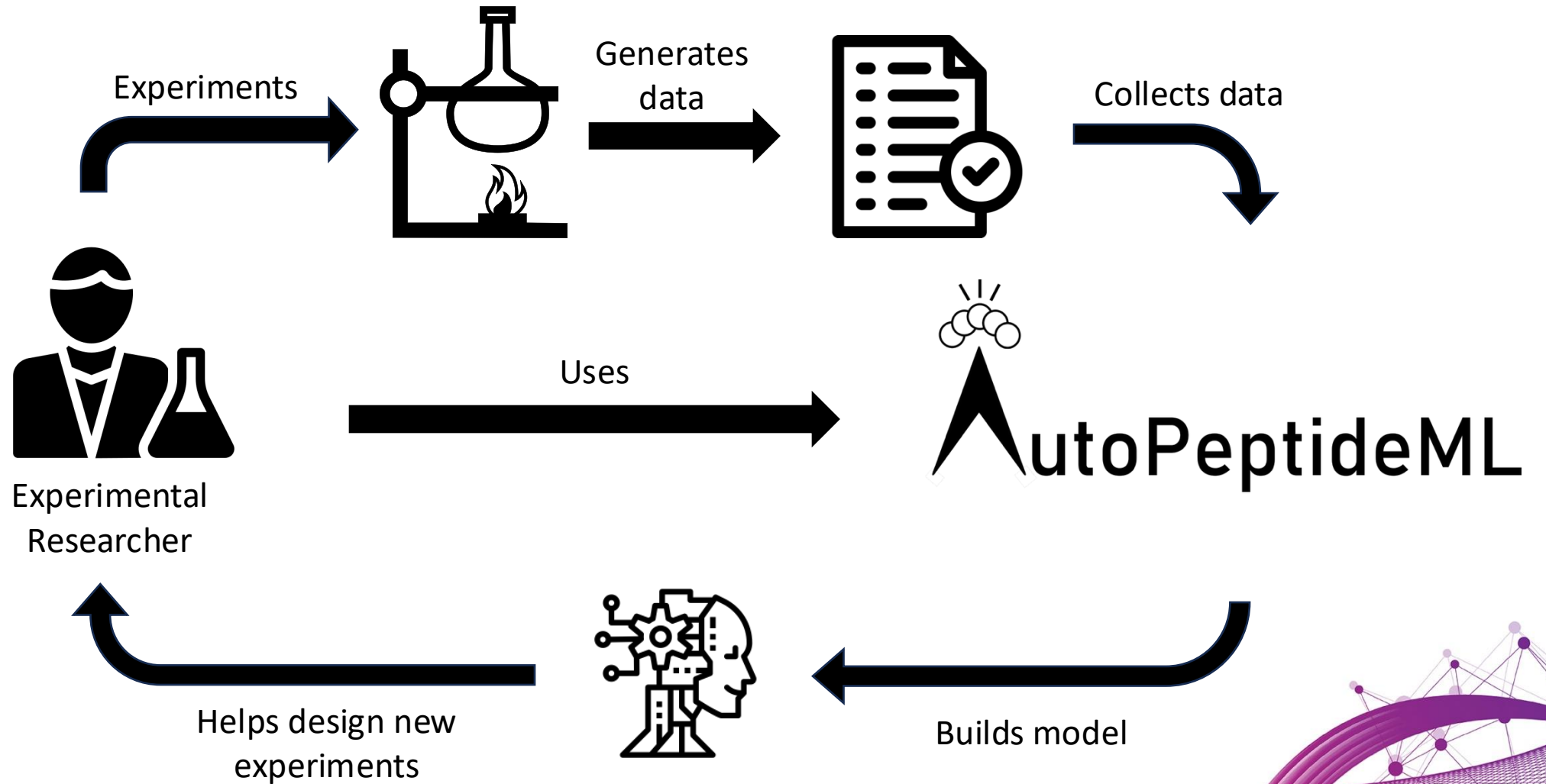
Part 0 - Introduction



More information
and contact info



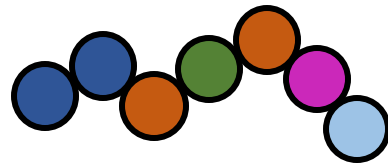
Main objective





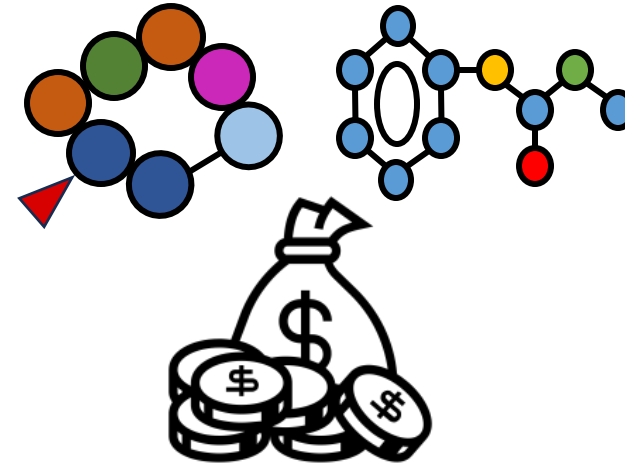
A tale of two peptides

Natural (or canonical) peptides
(cheap)



FOR SALE

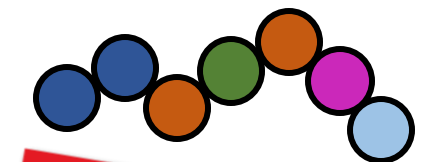
Synthetic peptides (or non-canonical) and
peptidomimetics
(expensive but better drugs)





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

Natural peptides
(cheap)



FOR SALE

Train

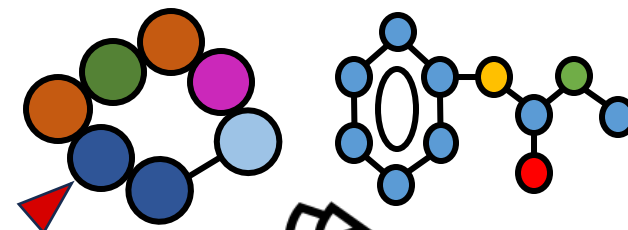


AutoPeptideML

Predict



Synthetic peptides and peptidomimetics
(expensive but better drugs)





Objectives

1. How to automatically build peptide property prediction models (and evaluate them)
2. How to extrapolate from standard to modified peptides or peptidomimetics

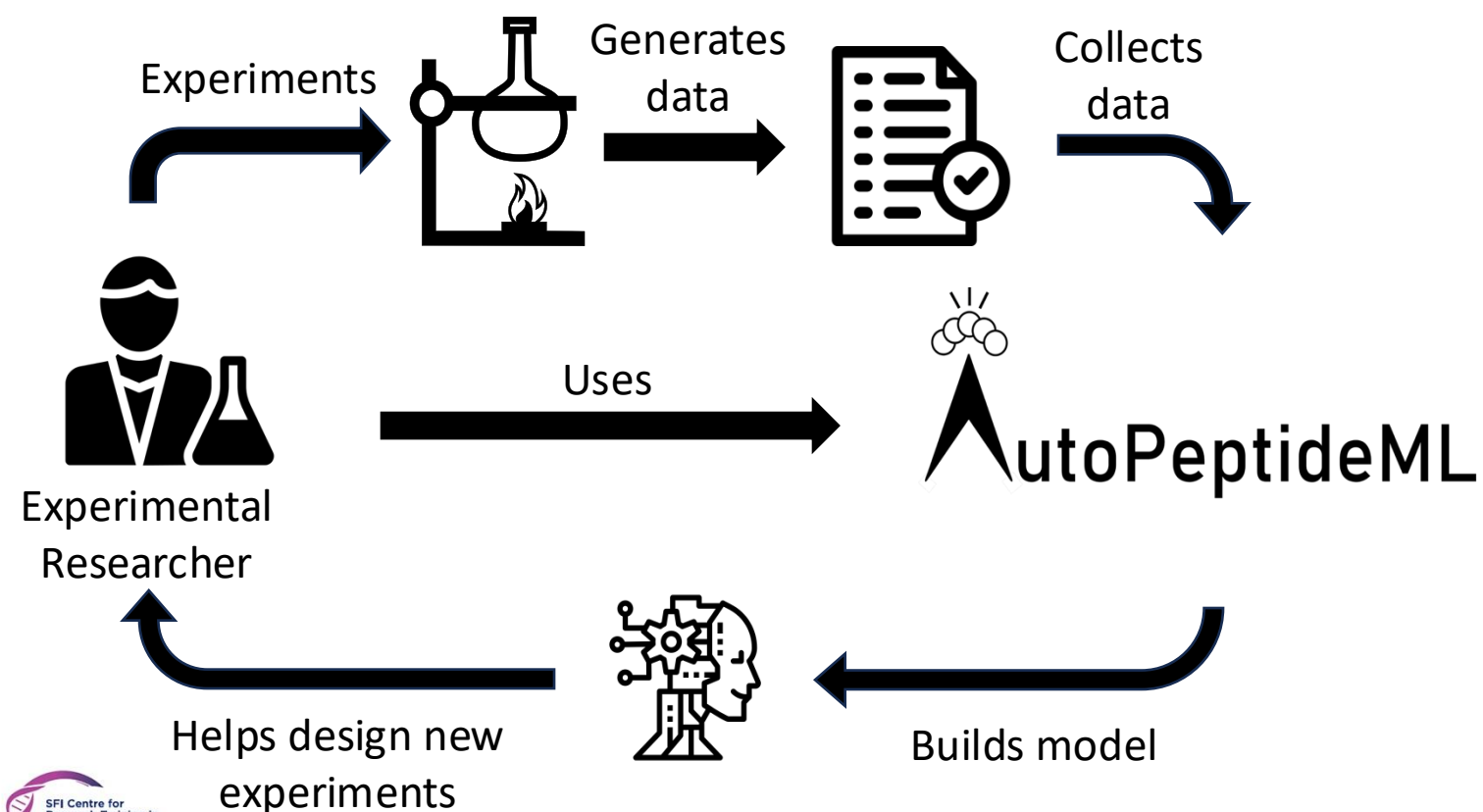


Part 1 – Automating ML for natural peptides





Objectives



Design Requirements

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

More information
and contact info



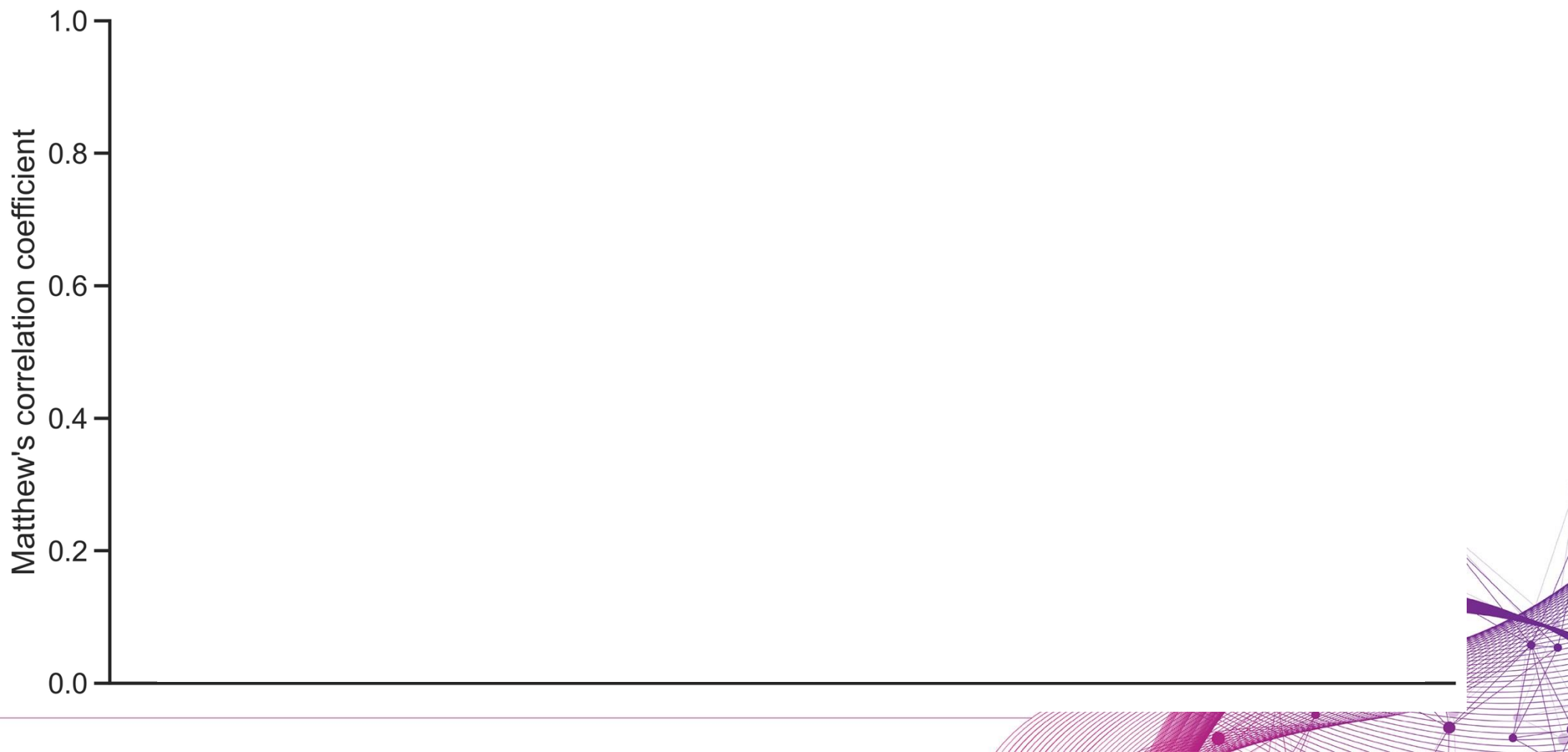
Collected **18**
datasets used for
building different
peptide bioactivity
predictors



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



More information
and contact info



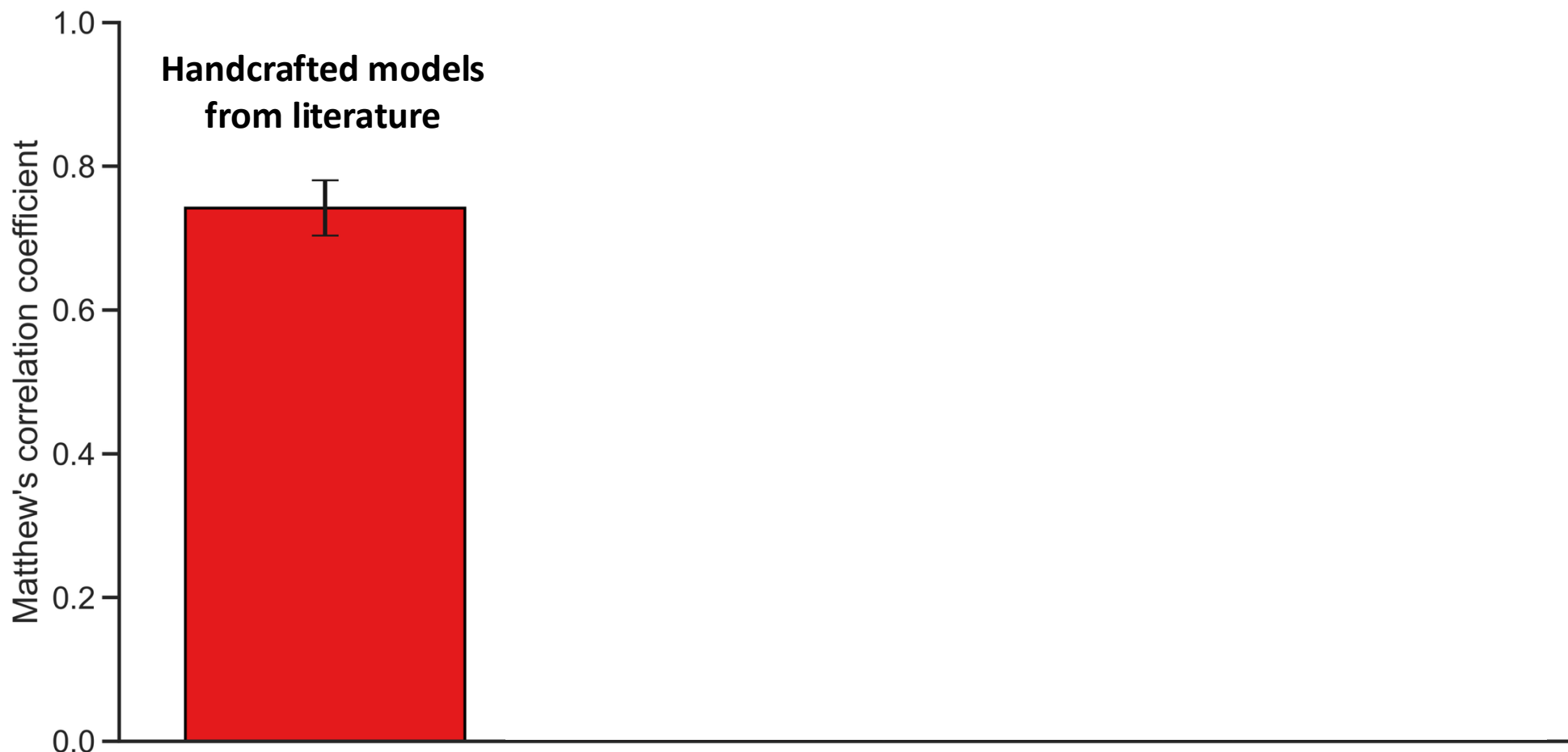
Collected **18**
datasets used for
building different
peptide bioactivity
predictors



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



More information
and contact info

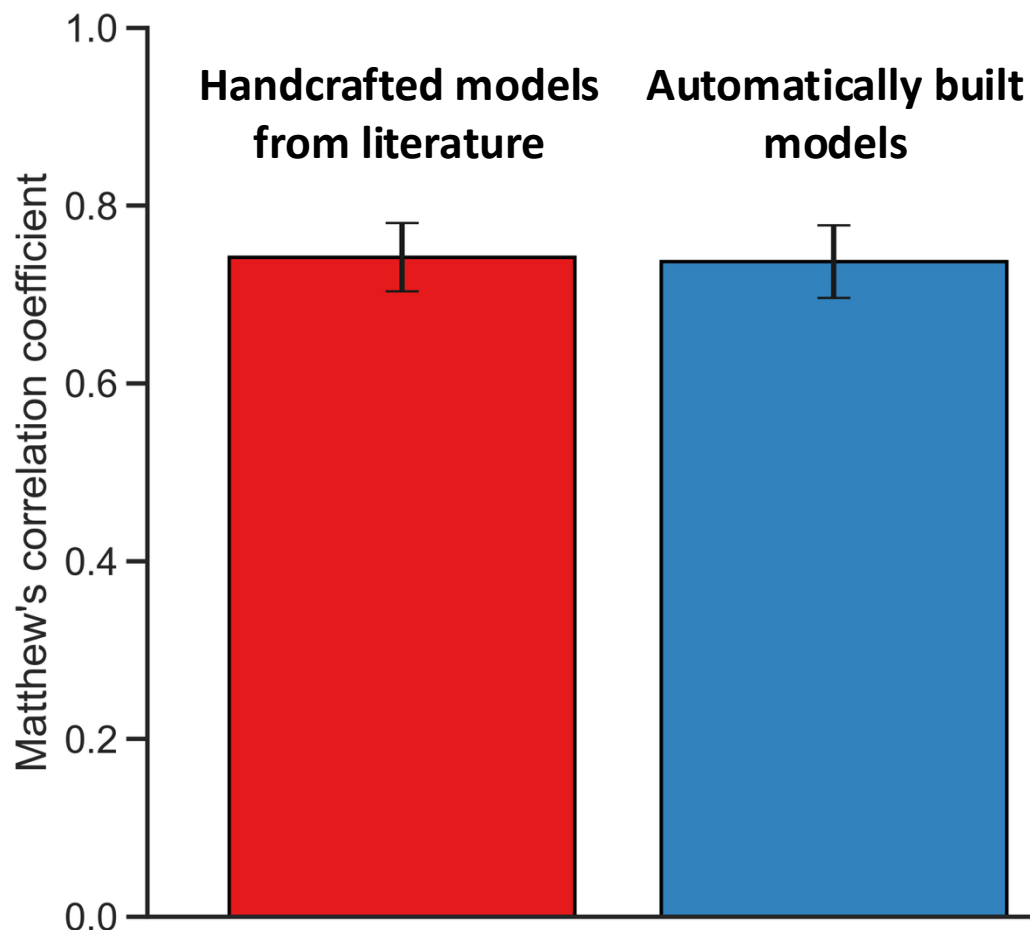


Automating ML for natural peptides



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

Collected **18**
datasets used for
building different
peptide bioactivity
predictors



Low intensity computing



Bayesian Optimization for
hyperparameter selection

**Protein Language
Models**

General representation/
featurization

More information
and contact info

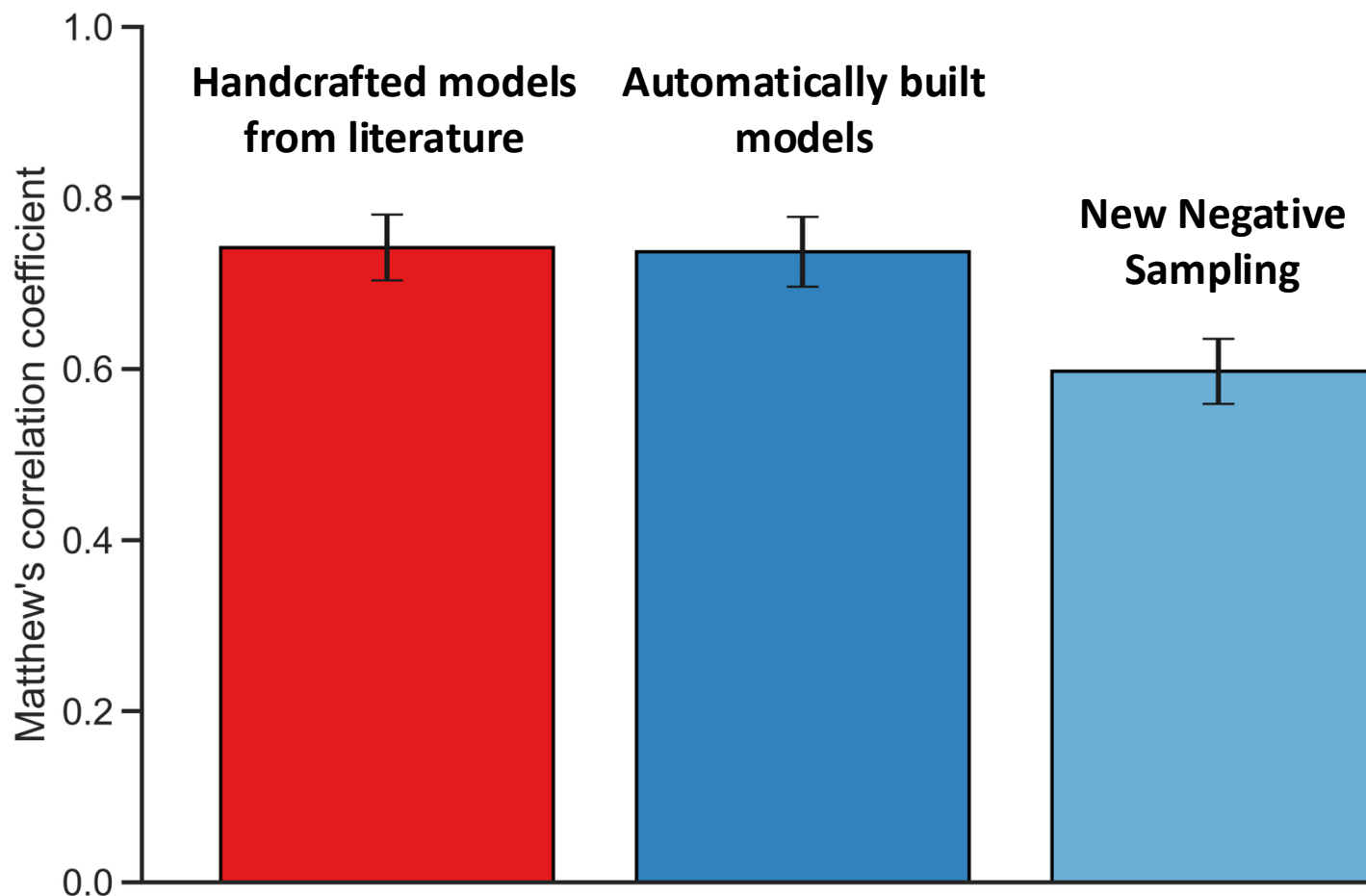


Automating ML for natural peptides



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

Collected **18**
datasets used for
building different
peptide bioactivity
predictors



Before:

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

Now:

- Other bioactive peptides

Peptipedia

More information
and contact info

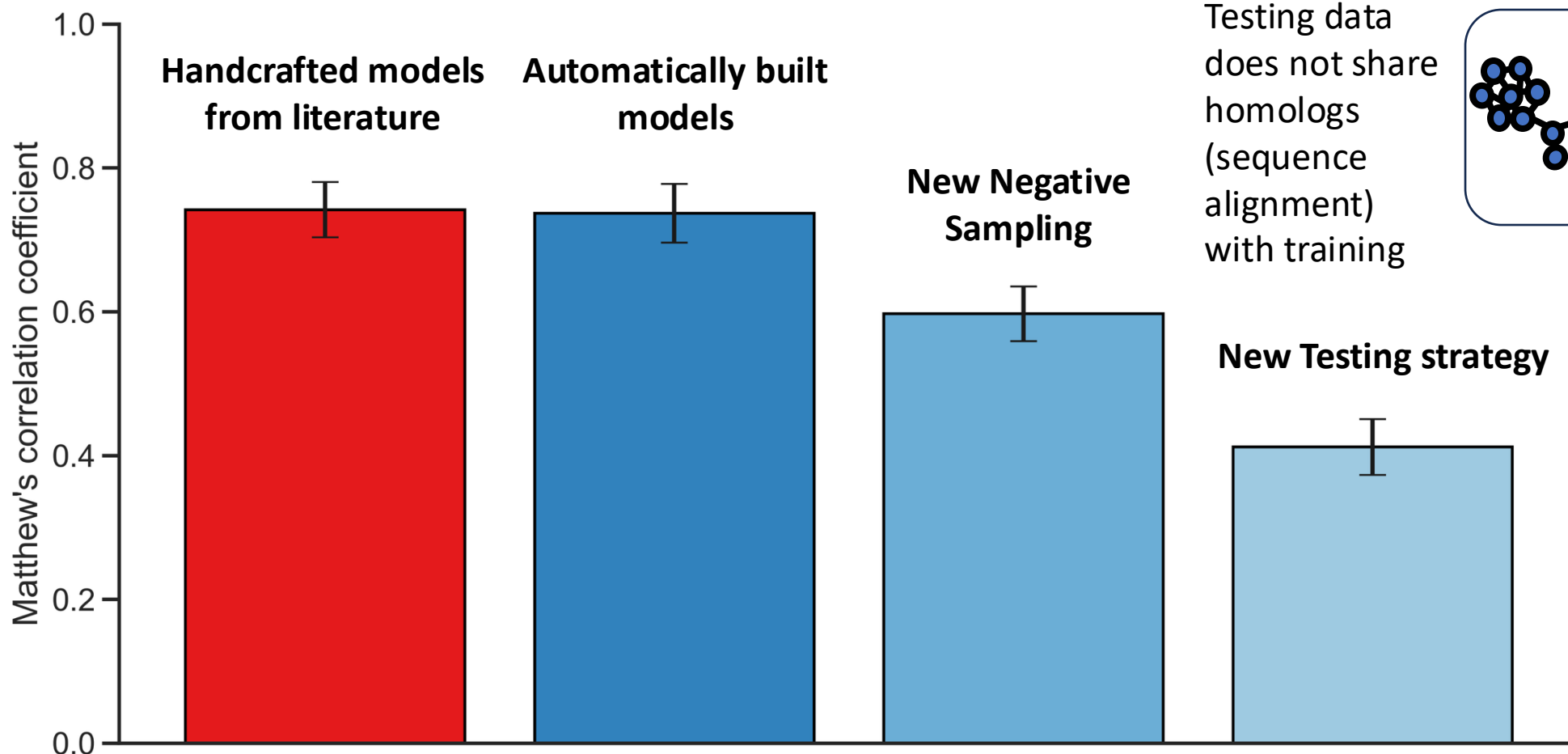


Automating ML for natural peptides



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

Collected **18**
datasets used for
building different
peptide bioactivity
predictors



More information
and contact info



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



BIOINFORMATICS PUBLICATION CHEMRXIV PREPRINT GITHUB REPOSITORY
DOCUMENTATION

Welcome to the AutoPeptideML webserver.

The next steps will help you build your own model.

0. Modelling task

First, start by defining the prediction task.

What is the prediction problem you are facing?

Classification (categorical values)

1. Inputs

In this section you will define the data from which you want the model to learn.

☐ Download sample dataset

Please upload dataset with your peptides and their labels if available

Drag and drop file here
Limit 200MB per file

Browse files

CLI tool

```
| AutoPeptideML v.2.0.3 |  
| By Raul Fernandez-Diaz |
```

Model builder

```
Part 1 - Define the data and preprocessing steps  
[?] What is the modelling problem you're facing?:  
> Classification (returning categorical value)  
Regression(returnin continuous value)
```

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
```



Automating ML for natural peptides - Conclusions

1. Automation achieves model performance on par with manually engineered previous studies
2. Proper automation leads to more robust model evaluation
3. Previous studies tended to overestimate model performance, due to:
 - a) Negative sampling strategy
 - b) Data leakage from similar peptides in training and testing



Part 2 – Natural to synthetic peptides extrapolation

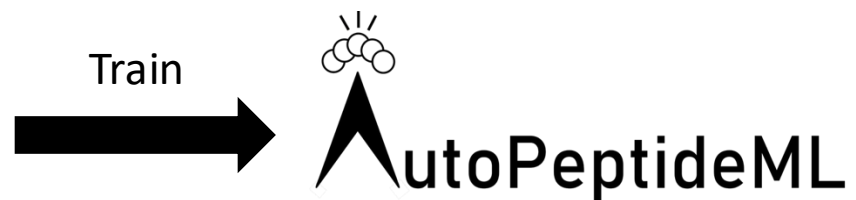




Can we leverage data on cheaper experiments to prioritise more expensive experiments?



Train



Predict

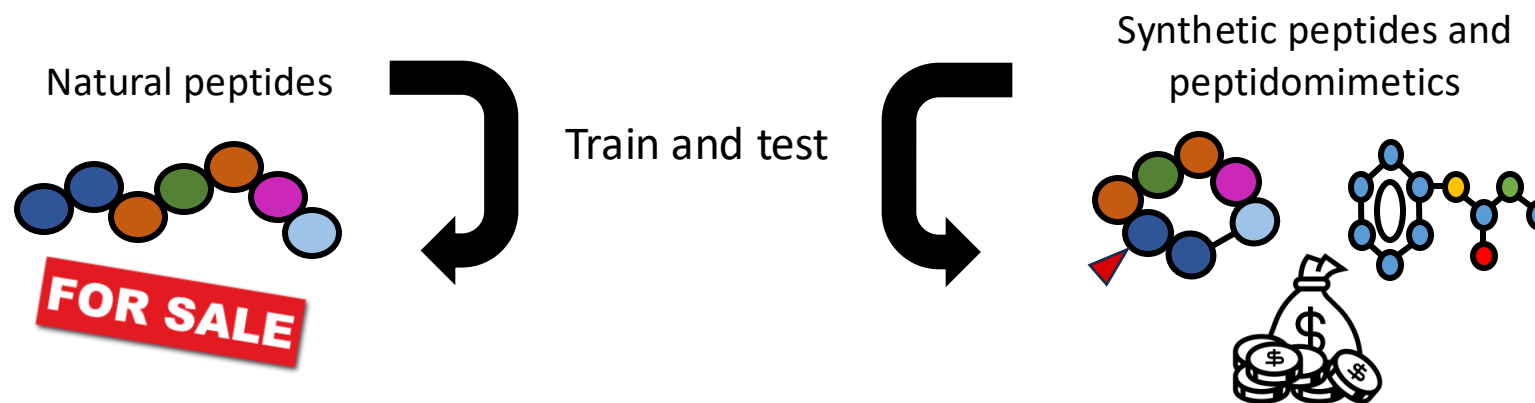
Synthetic peptides and peptidomimetics
(expensive but better drugs)



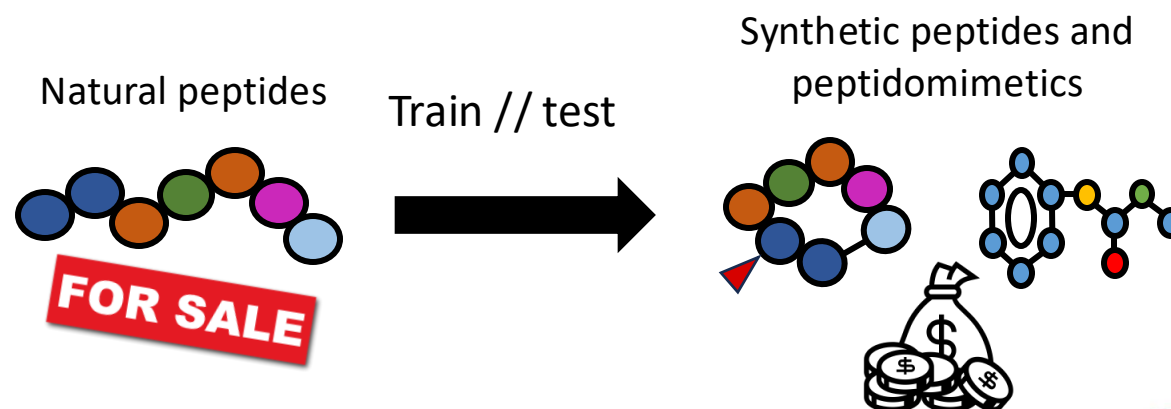


Computational experiments

Interpolation



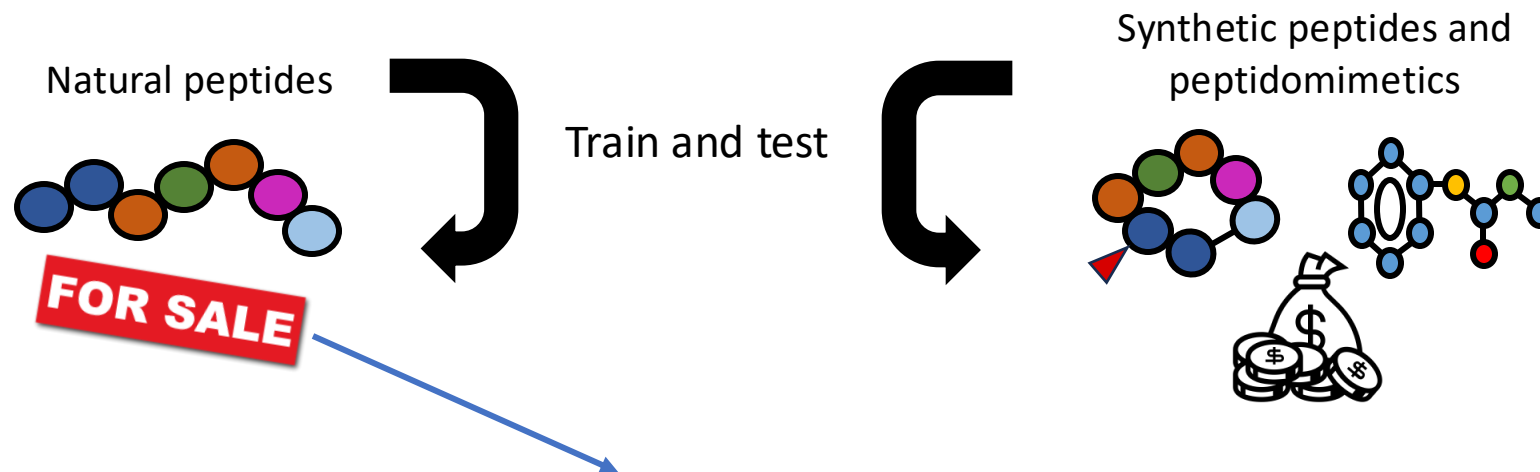
Extrapolation





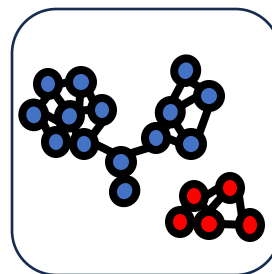
Computational experiments

Interpolation



Small digression: Should we use sequence alignment for measuring peptide similarity?

Similarity-informed
train/test split





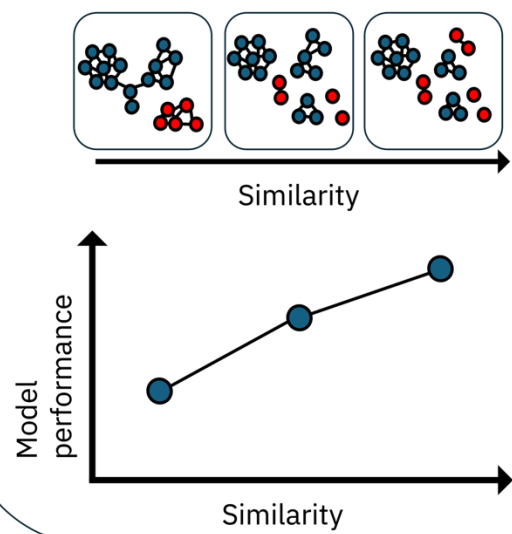
Finding the best similarity metric

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

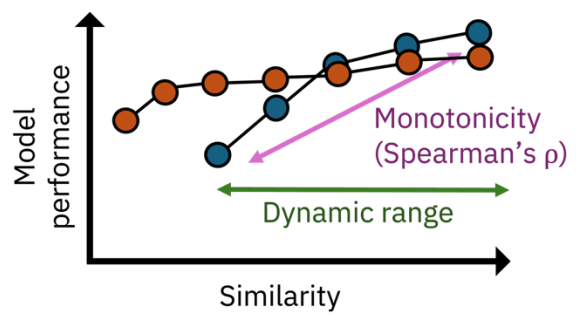
GOOD curve

Model performance as a function of train-test similarity



Similarity metric selection

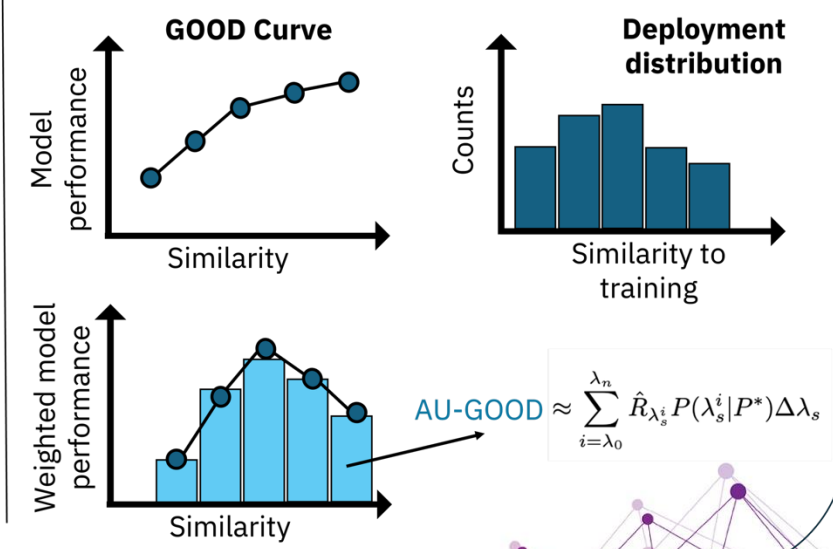
Quantitative analysis of best similarity function for a given task/dataset



- **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



Best metrics for each dataset

Dataset	Peptide type	Task	Similarity Type	Similarity	Dynamic range (↑) [a]	Monotonicity (↑) [b]
Protein-peptide binding affinity	Standard	Regression	Chemical FP	MAPc-8	70 %	0.8 ± 0.1
Protein-peptide binding affinity	Modified	Regression	Chemical FP	MAPc-20	80 %	0.95 ± 0.03
Cell penetration	Standard	Classification	Chemical FP	MAPc-8	60 %	0.98 ± 0.04
Cell penetration	Modified	Classification	Chemical FP	MAPc-12	60 %	0.5 ± 0.2
Antibacterial	Standard	Classification	Chemical FP	MAPc-8	60 %	0.97 ± 0.02
Antibacterial	Modified	Classification	Chemical FP	ECFP-12	50 %	0.9 ± 0.1
Antiviral	Standard	Classification	Sequence Alignment	MMSeqs2	80 %	0.96 ± 0.05
Antiviral	Modified	Classification	Chemical FP	MAPc-12	70 %	0.6 ± 0.2

Metrics explored:

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



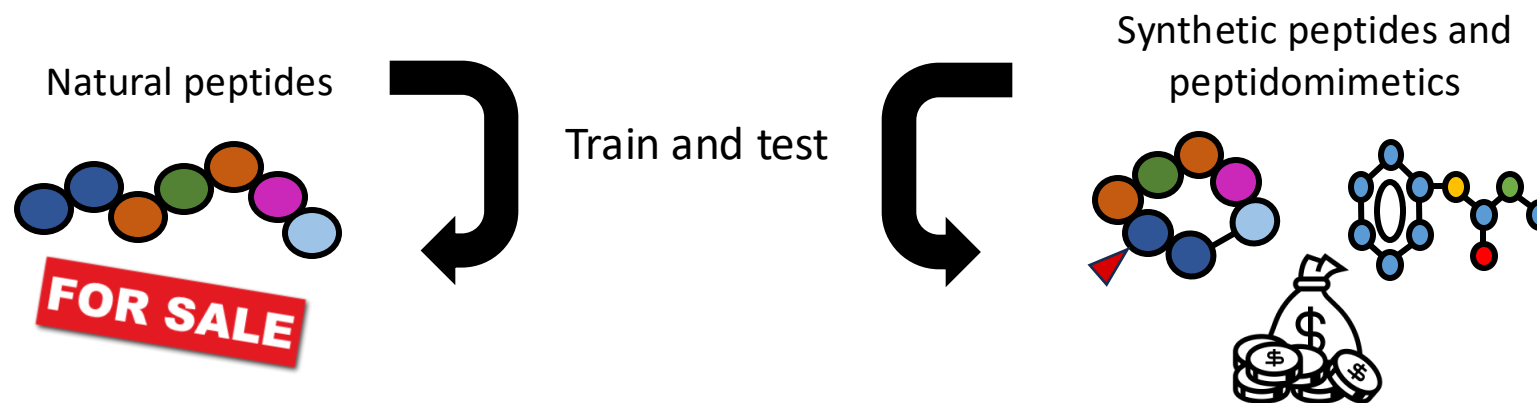
More information
and contact info

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

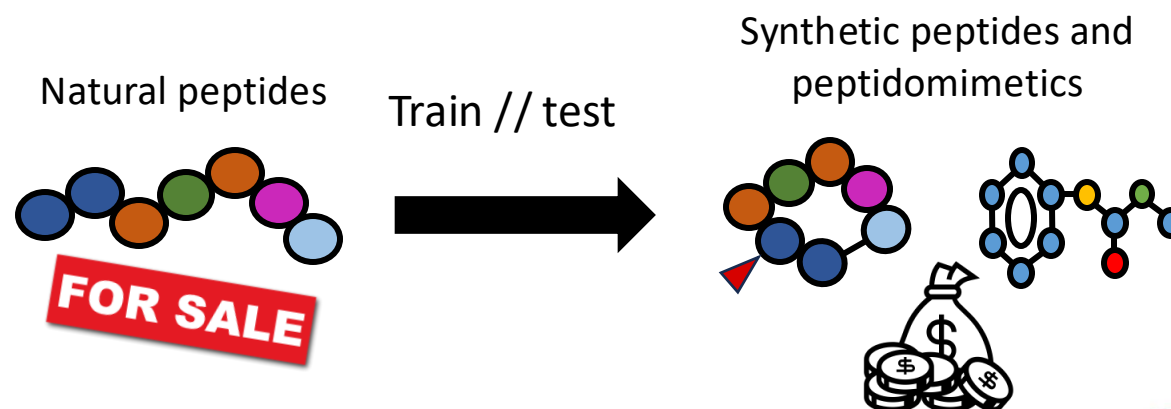


Computational experiments

Interpolation



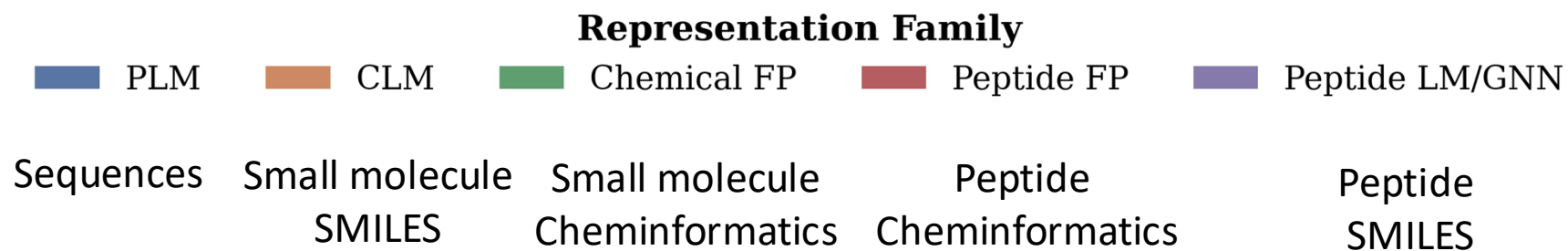
Extrapolation





Interpolation experiments

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

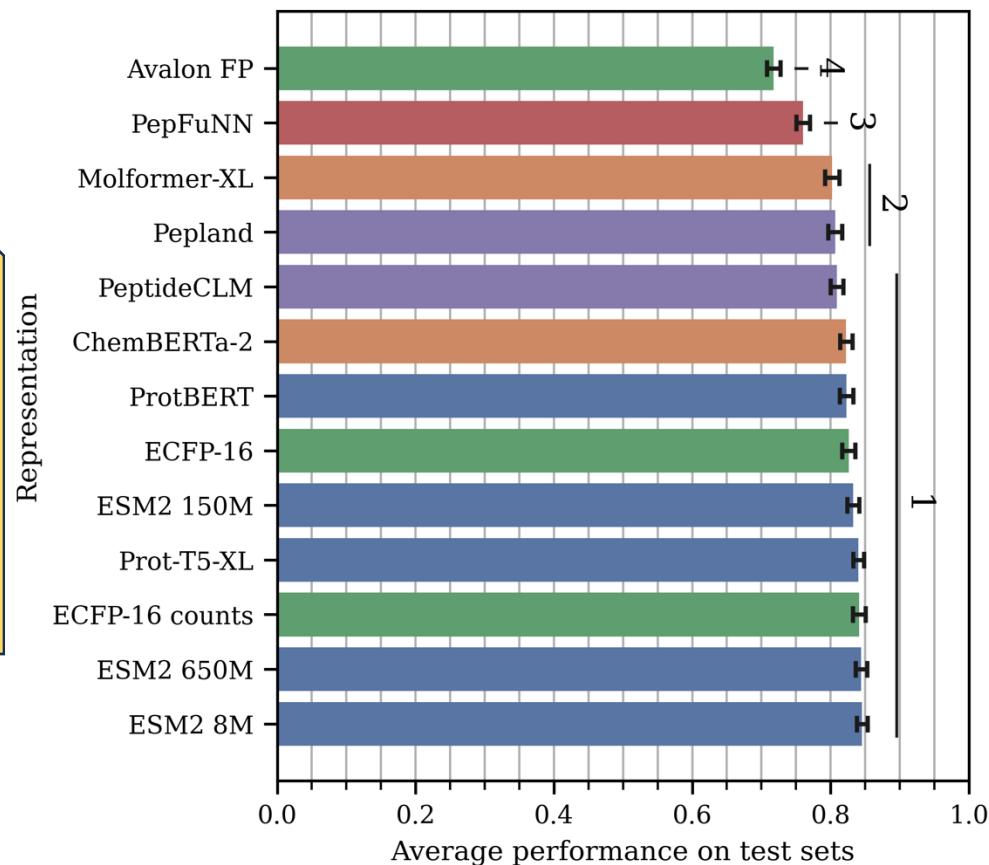




Interpolation experiments

Synthetic to synthetic

Natural to natural



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



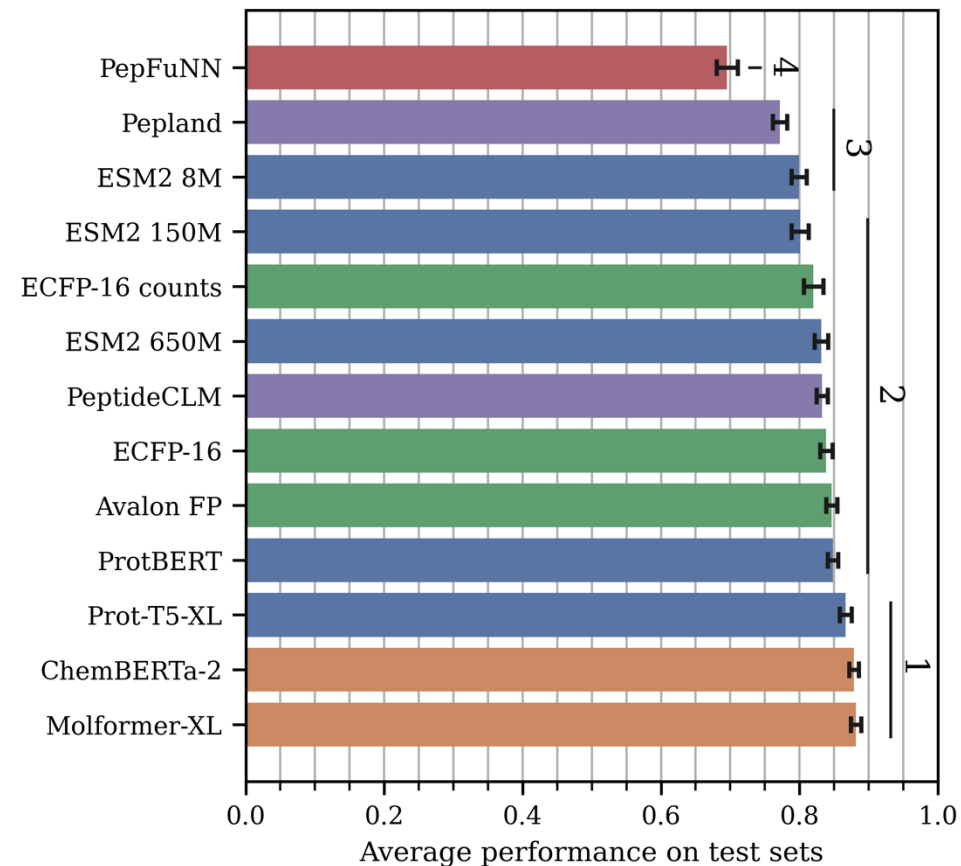


Interpolation experiments

Natural to natural

Synthetic to synthetic

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

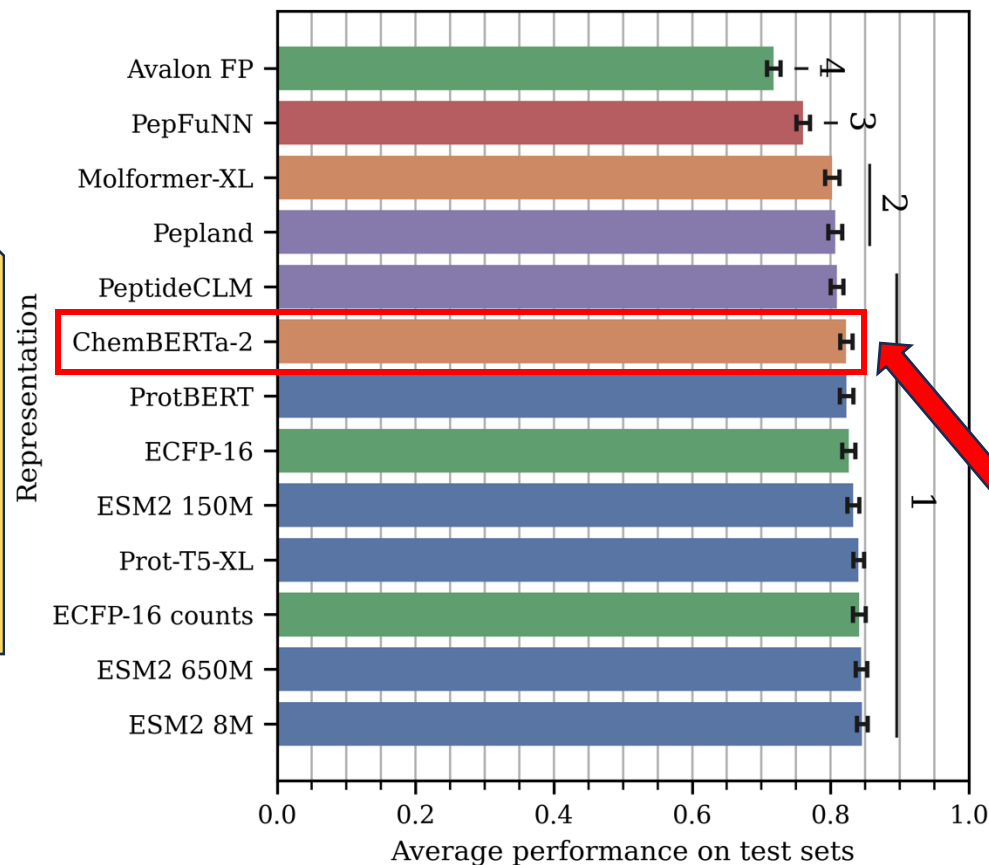




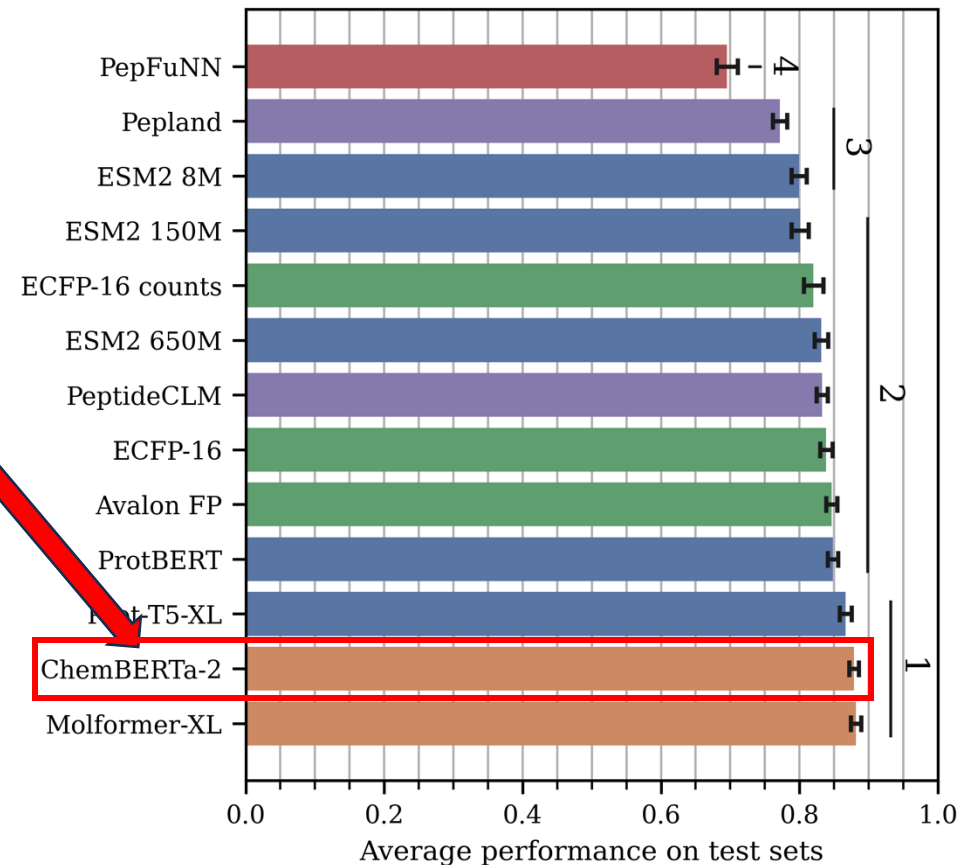
Interpolation experiments

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

Natural to natural



Synthetic to synthetic



Representation Family

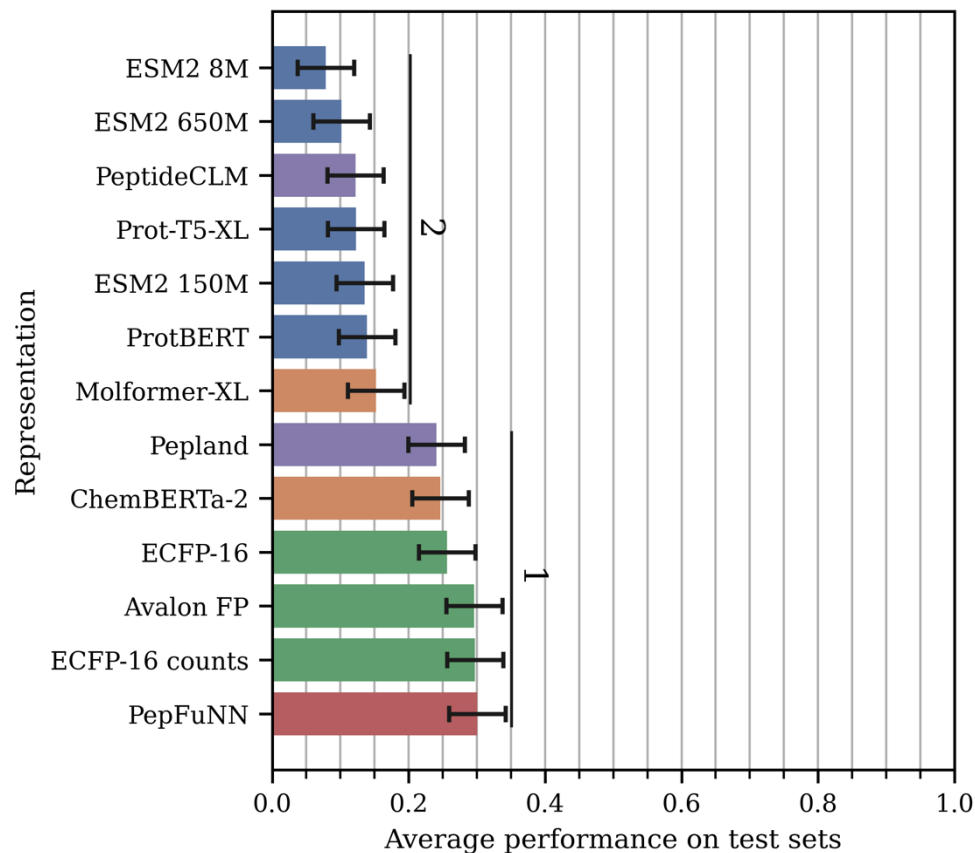
PLM CLM Chemical FP Peptide FP Peptide LM/GNN





Natural to synthetic extrapolation

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



Representation Family

PLM CLM Chemical FP Peptide FP Peptide LM/GNN

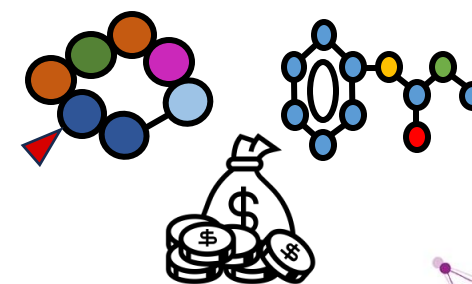
Natural peptides



FOR SALE

Train // test

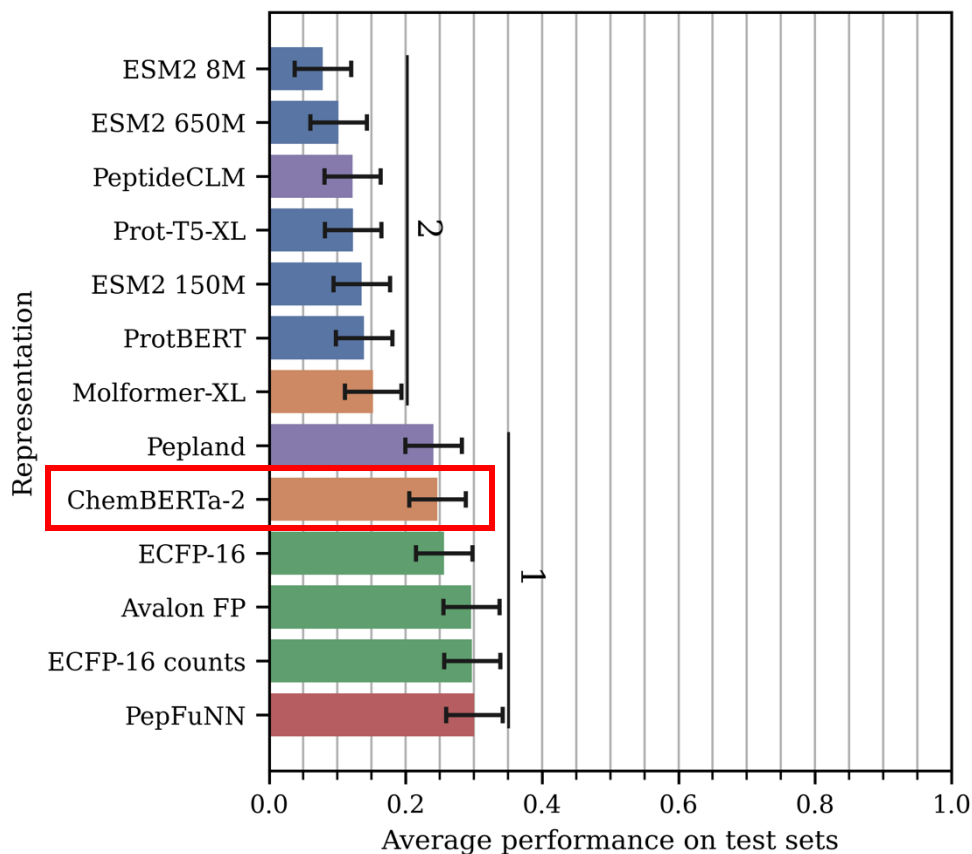
Synthetic peptides and peptidomimetics





Natural to synthetic extrapolation

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



Conclusions

1. 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
3. Chemical and Peptide Fingerprints are robust options as well

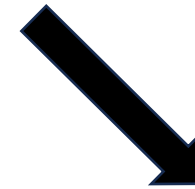




Conclusions

1. AutoPeptideML **empowers** experimental scientist to build their own models
2. Dataset building (negative definition) and partitioning (train/test split) are key for proper model evaluation
3. Chemical fingerprints are better for partitioning natural and synthetic datasets than sequence alignment.
4. Natural to synthetic extrapolation is possible, but there is room for improvement
5. ChemBERTa-2 appears to be the most versatile tool, closely followed by chemical fingerprints

Contact info, papers, and
slides of the presentation



More information
and contact info



How to generalize machine learning models to both canonical and non-canonical peptides

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

More information
and contact info

