

# Solubility prediction of lipid compounds using machine learning

A. Porley Santana<sup>a#</sup>, G. Gutierrez<sup>a#</sup>, S. Gutiérrez Parodi<sup>a</sup>, J. Ferreira<sup>a,b\*</sup>

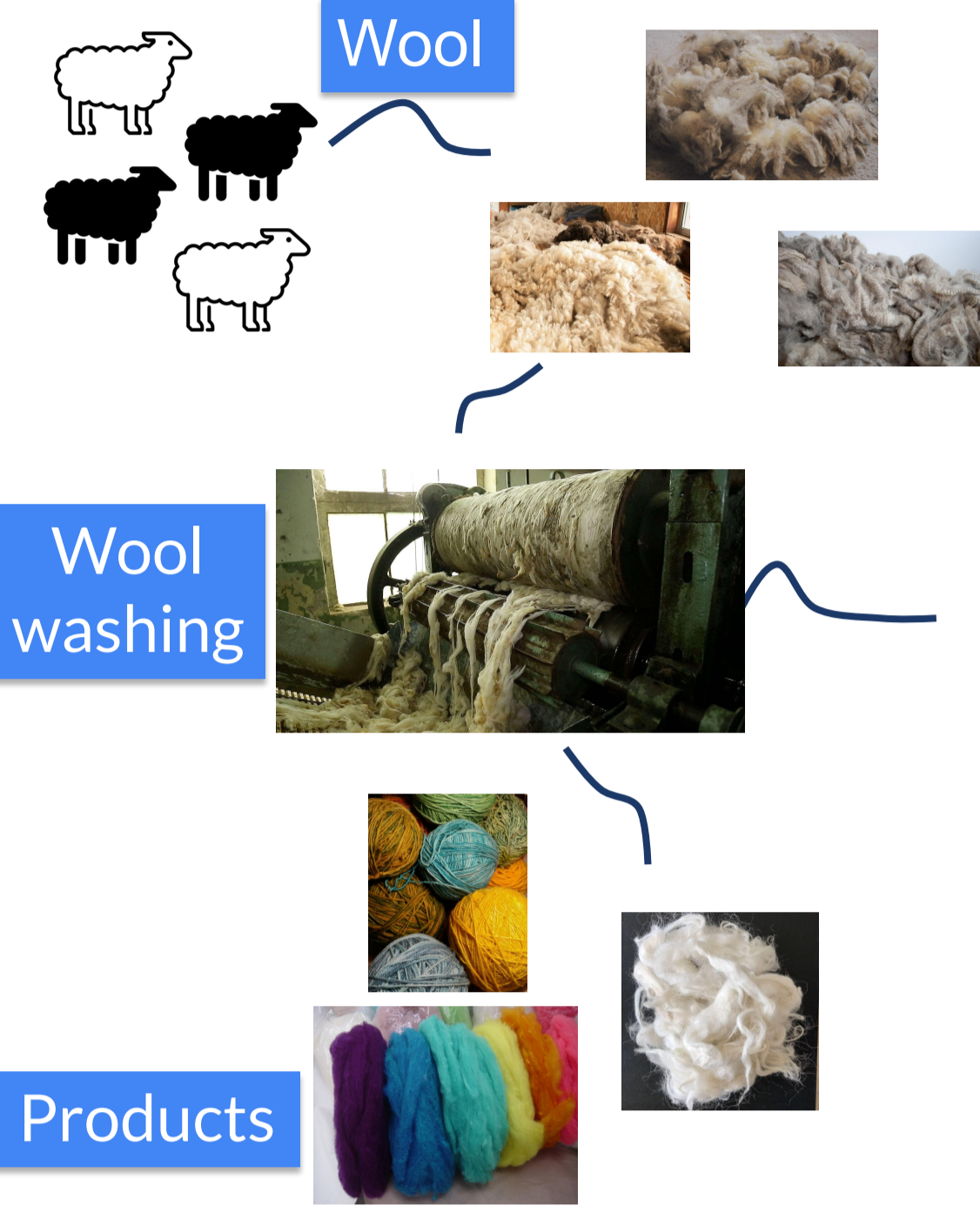
<sup>a</sup> Chemical & Process Systems Engineering Group, Faculty of Engineering, Universidad de la República, Montevideo, Uruguay

<sup>b</sup> Heterogeneous Computing Laboratory, Faculty of Engineering, Universidad de la República, Montevideo, Uruguay

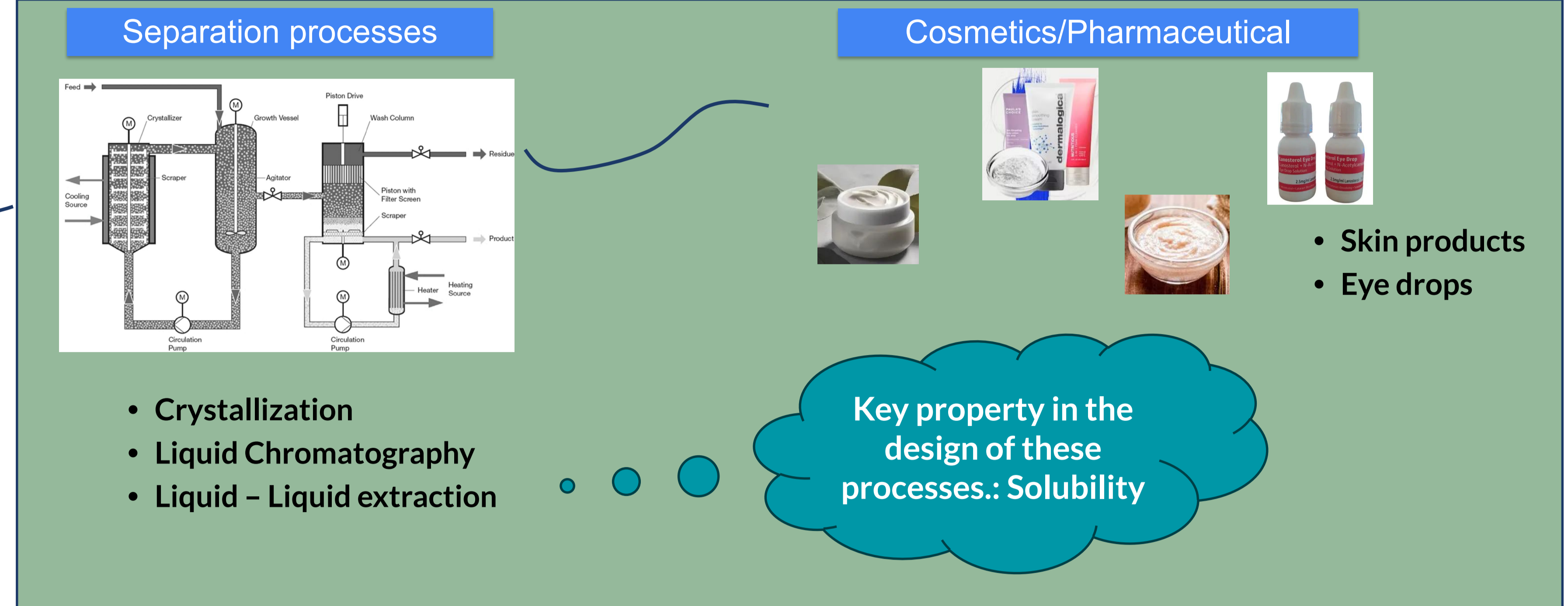
\* [jimenaf@fing.edu.uy](mailto:jimenaf@fing.edu.uy)

Introduction

## Wool Processing



## High-value products



## Objectives

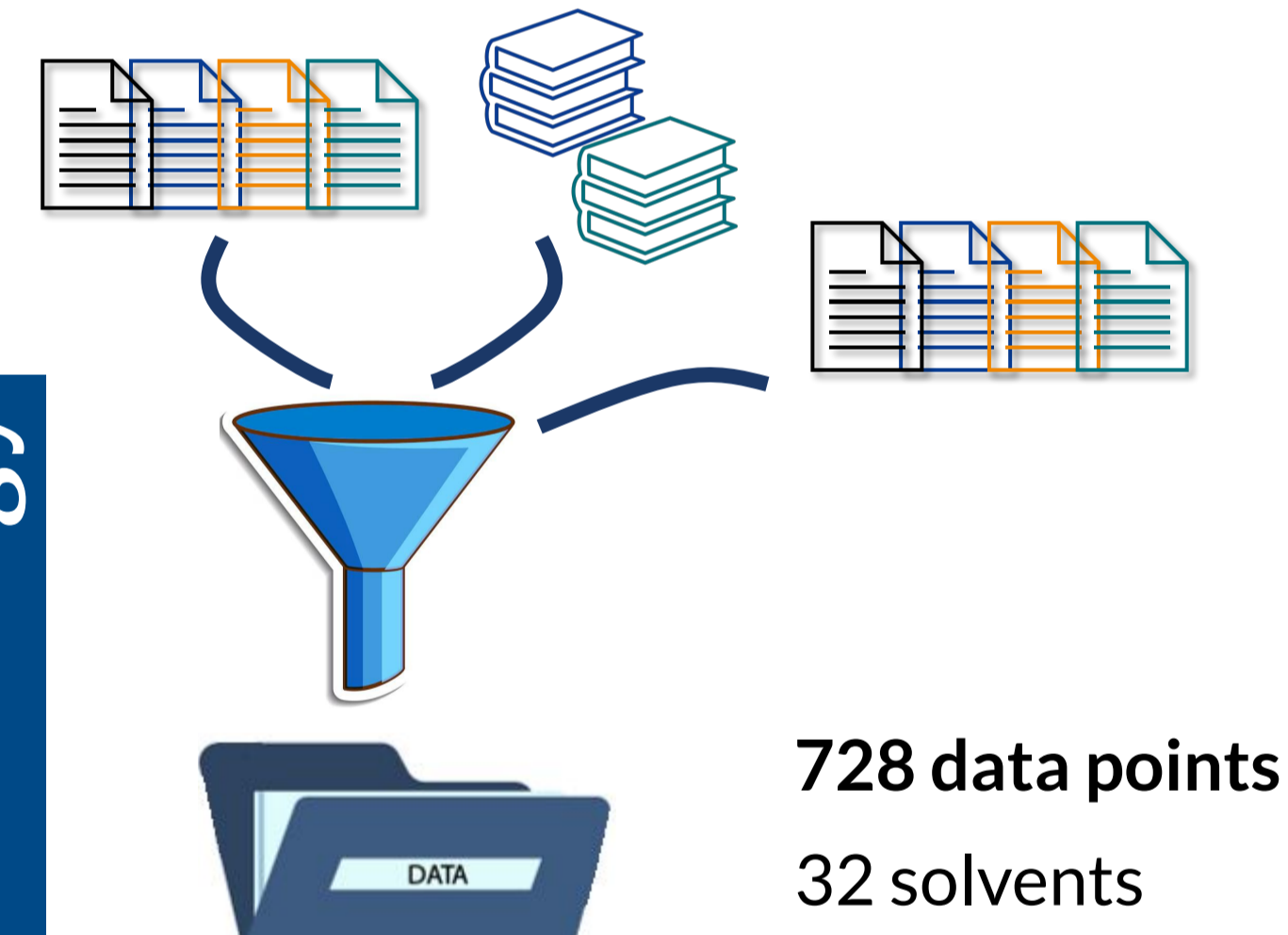
- Developing a **machine learning-based** predictive model to estimate the solubility of five key wool grease lipids (cholesterol, palmitic acid, oleic acid, stearic acid, and linolenic acid) along with naphthalene in a variety of organic solvents.
- Model is trained using **experimentally derived data** sourced from scientific articles and handbooks.

## Solubility

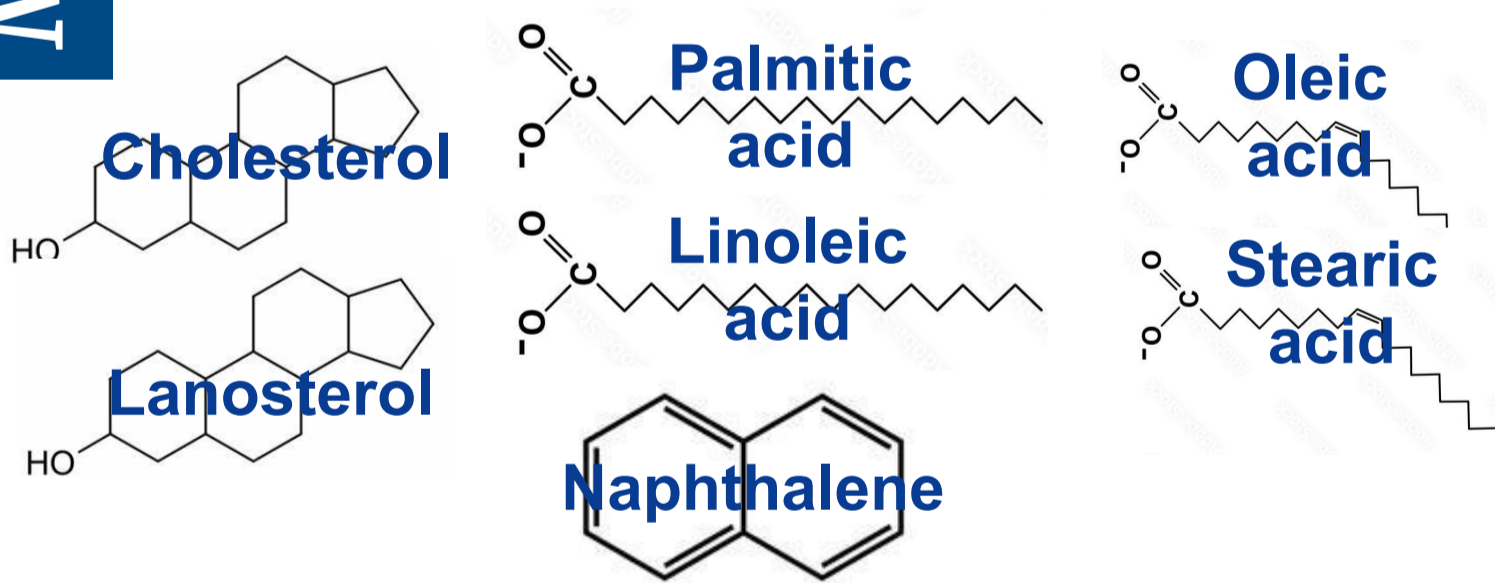
- Solubility can be performed by several experimental methods: **labor-intensive and time-consuming**.
- Solubility models: powerful **tool to accelerate** solvent selection processes in **process design**.
- Not specific for lipids:
  - Model based on thermodynamic properties: Hildebrand [7].
  - Model based on thermodynamic and molecular descriptor: Abraham [8].
  - Hybrid model based on free energy and solvation enthalpy using ANN [9].

## Data Collection

15 scientific papers and handbooks



Methodology



## Data Preprocessing

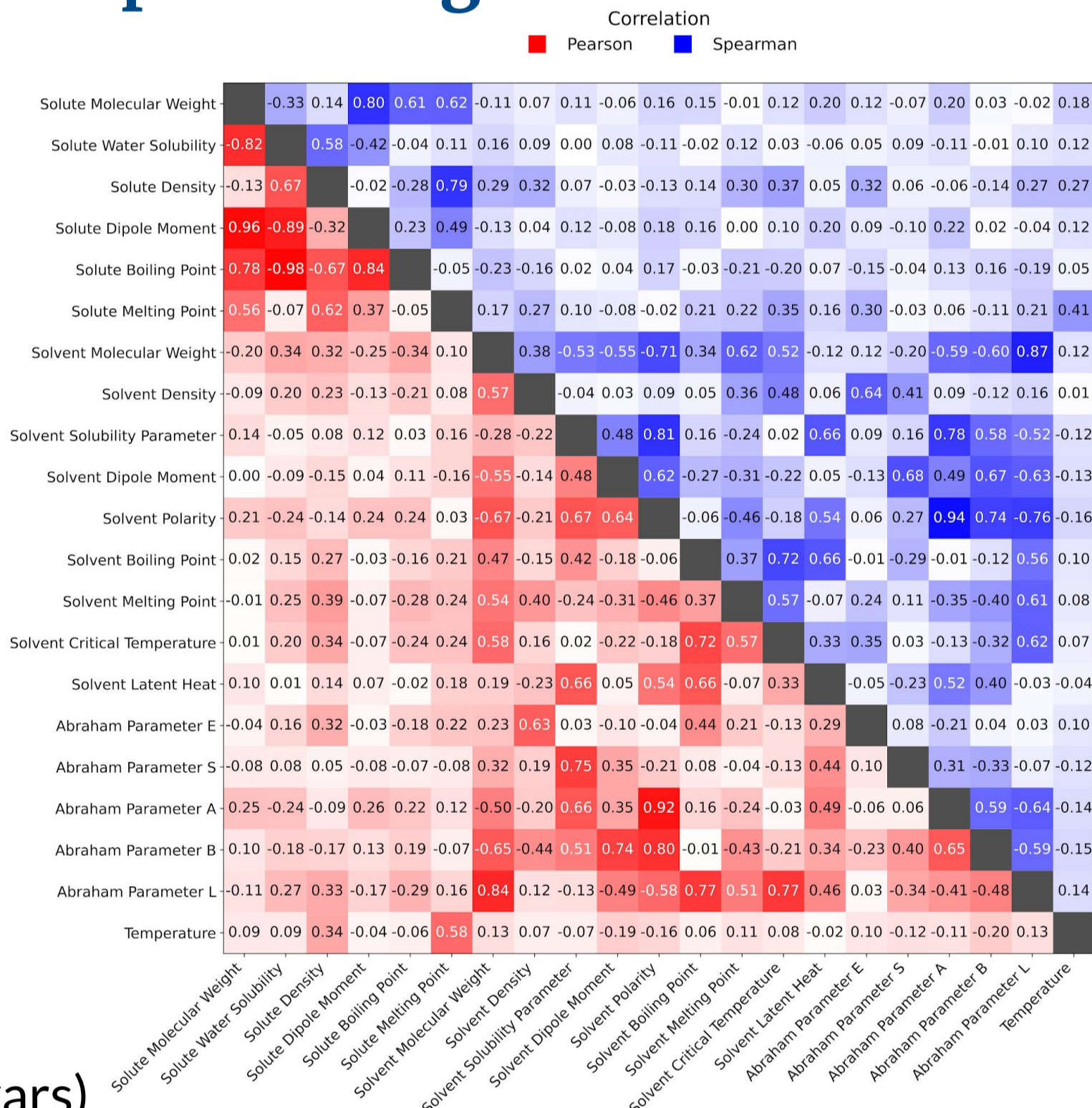
- 21 input variables
  - Variable normalization
  - Variable selection:
    - Pearson correlation
    - Spearman correlation
- Temperature  
Solute properties  
Solvent properties  
Abraham's solvation parameter

Randomly split

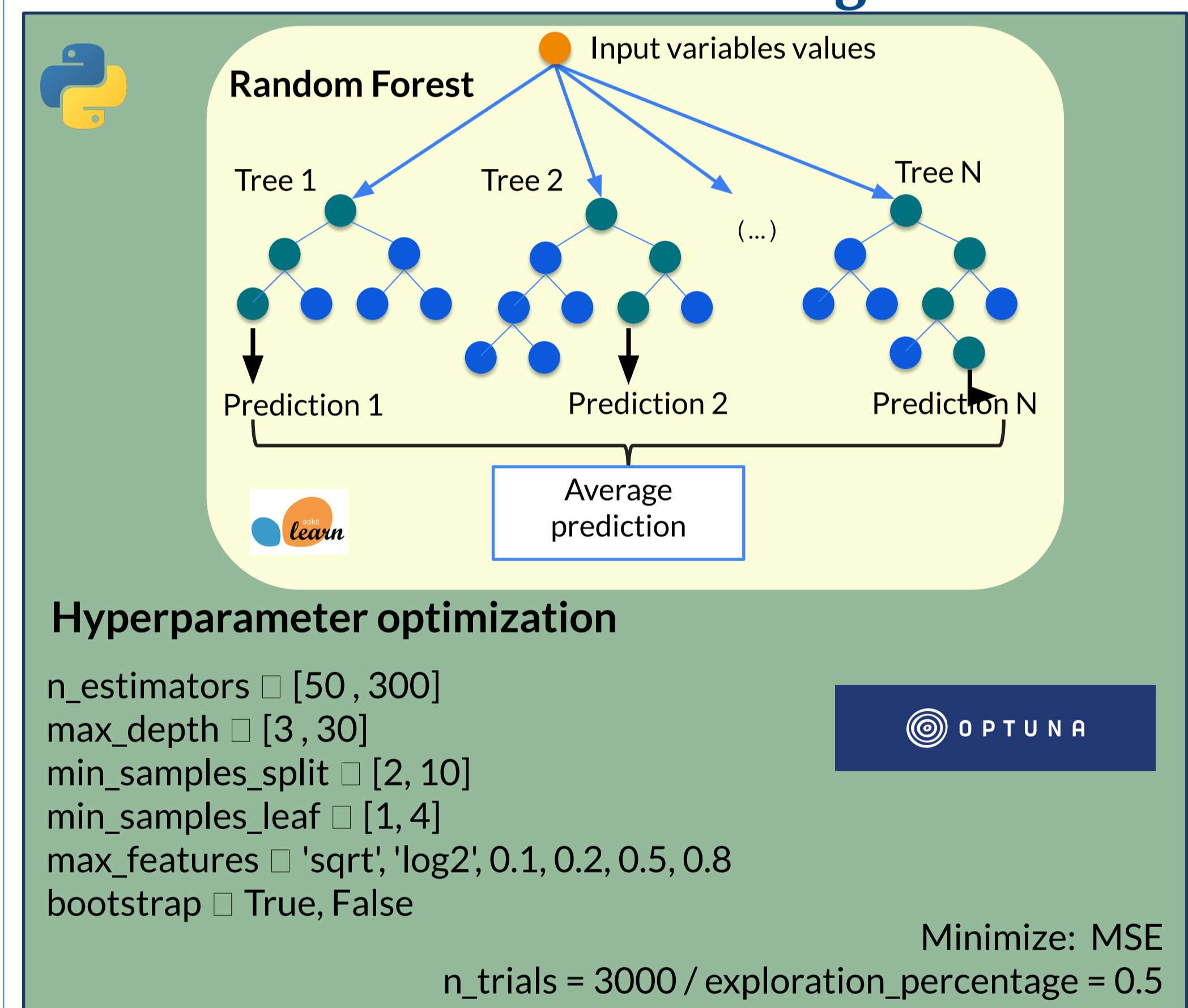
- Training: 70 %
- Testing: 30 %

Variations:

- Simple, S (9 uncorrelated vars)
- Without Abraham parameters, A, (16 vars)
- Complete (21 vars)



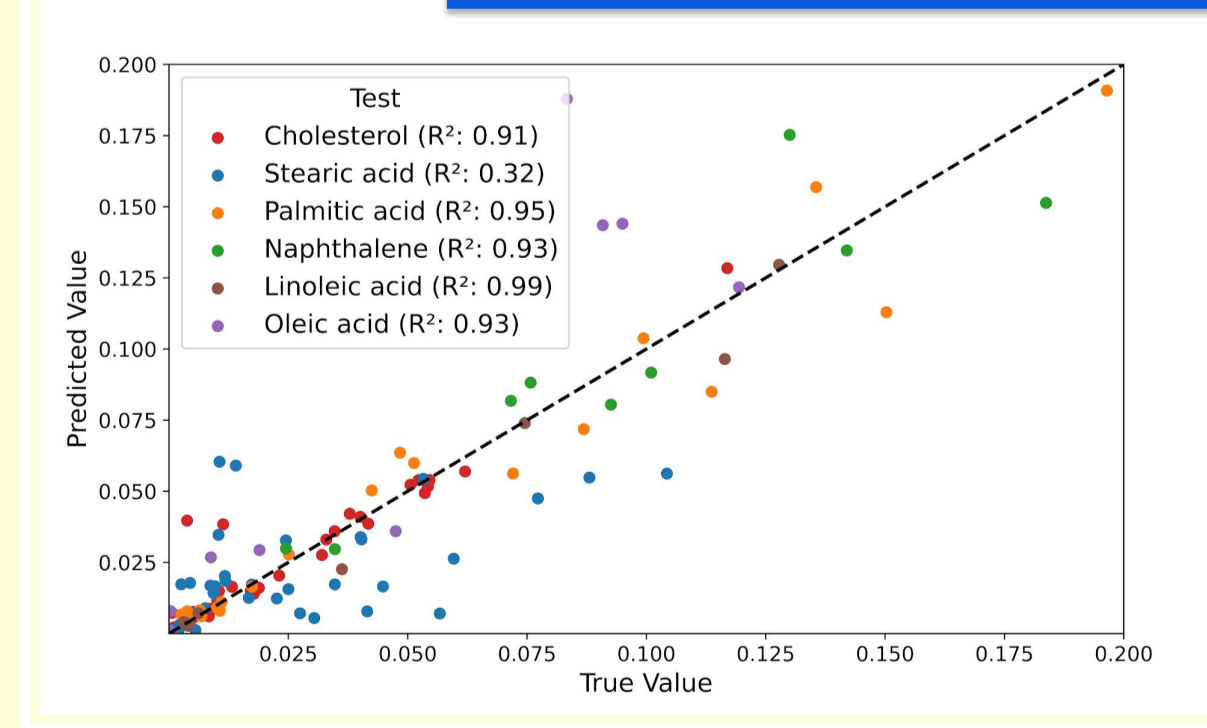
## Data-driven modeling



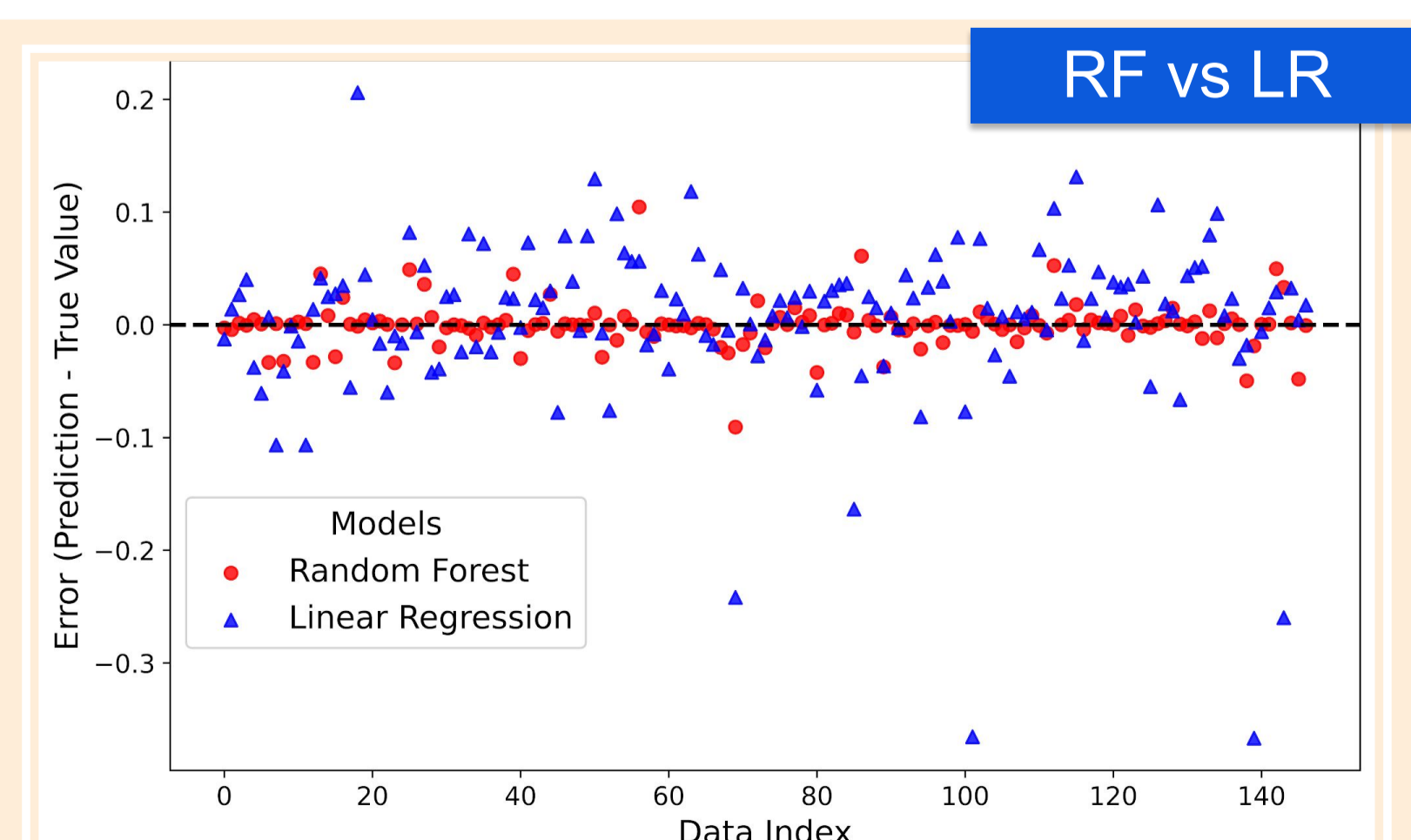
## Results and discussion

Variation	Linear Regression		Random Forest	
	R <sup>2</sup>	MSE	R <sup>2</sup>	MSE
S	0.30	6.7 x 10 <sup>-3</sup>	0.83	1.7 x 10 <sup>-3</sup>
A	0.45	5.3 x 10 <sup>-3</sup>	<b>0.97</b>	<b>0.3 x 10<sup>-3</sup></b>
C	0.47	5.1 x 10 <sup>-3</sup>	0.94	0.5 x 10 <sup>-3</sup>

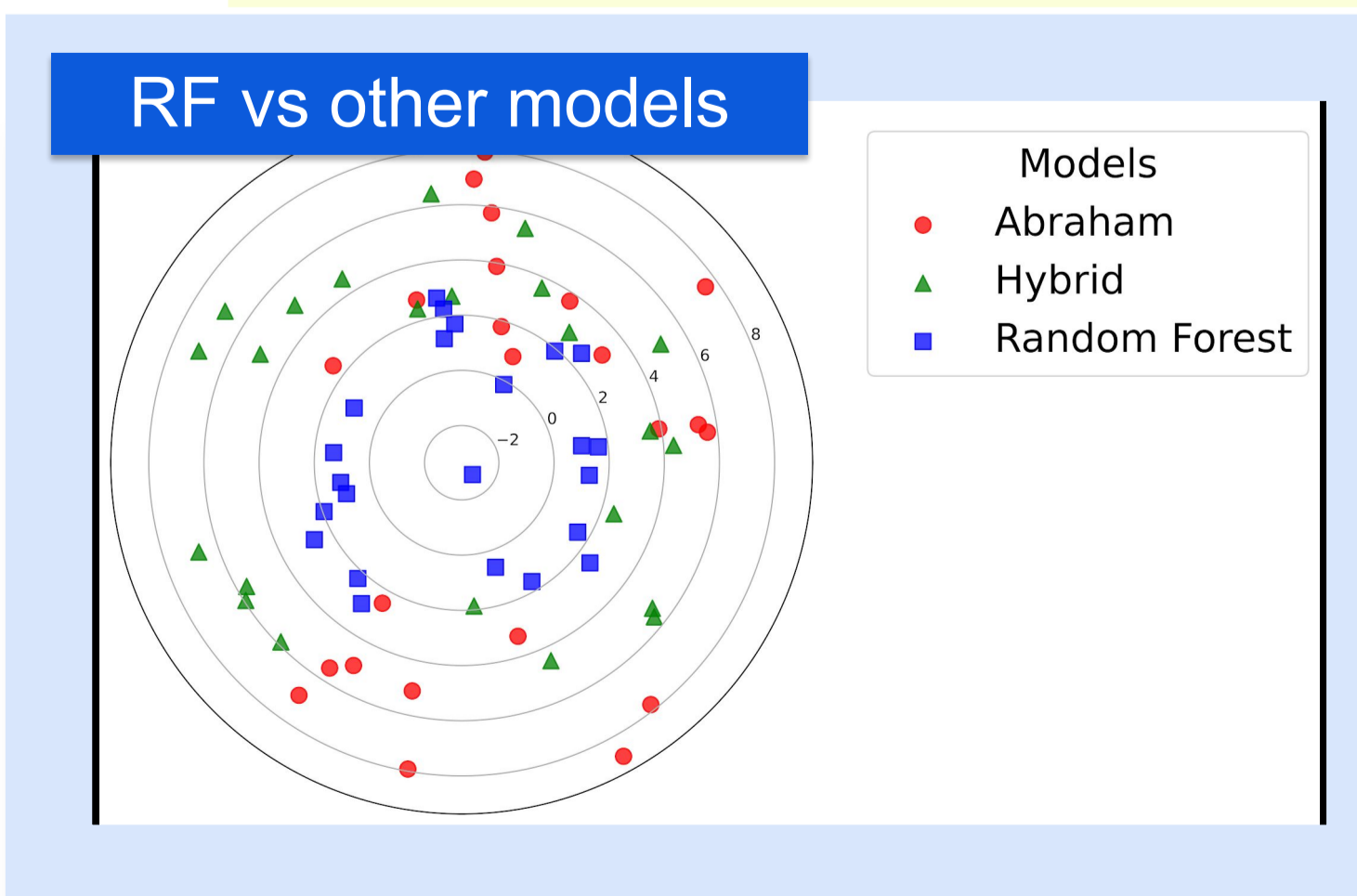
### Prediction by solute



### RF vs LR



### RF vs other models



## Future work

- Extend the dataset to enhance prediction accuracy and robustness.
- Other machine learning methods with strong extrapolation.
- The obtained model will be used as a basis for **designing biomass lipid separation processes**.

## References

- Gutiérrez, S., Viñas M. (2003). Water Sci. and Tech. <https://doi.org/10.2166/wst.2003.0379>
- Ding, H., et al. (2017). Chem. Pap. <https://doi.org/10.1007/s11696-016-0043-1>
- Pei, H., et al. (2019). J. of Separation Sci. <https://doi.org/10.1002/jssc.201900063>
- Diorazio, L. J., et al. (2016). Organic Proc. Res. & Dev. <https://doi.org/10.1021/acs.oprd.6b00015>
- Atkins, P., de Paula J. (2014). Physical Chemistry. 10th edition.
- Hermans, A., et al. (2022). The AAPS J. <https://doi.org/10.1208/s12248-022-00760-8>
- Hildebrand, J. H. (1936). Solubility of Non-Electrolytes. 2nd edition.
- Abraham, M. H., et al. (1987). J. of Chromatography A. [https://doi.org/10.1016/S0021-9673\(01\)86779-0](https://doi.org/10.1016/S0021-9673(01)86779-0)
- Chung, Y., et al. (2022). J. of Chem. Info. and Modeling. <https://doi.org/10.1021/acs.jcim.1c01103>

Acknowledgement CSIC