# TACsy

Training Alliance for Computational systems chemistry

# Synthesis Rebalancing Framework

**Tieu-Long Phan & Klaus Weinbauer**

Date: 13.02.2024
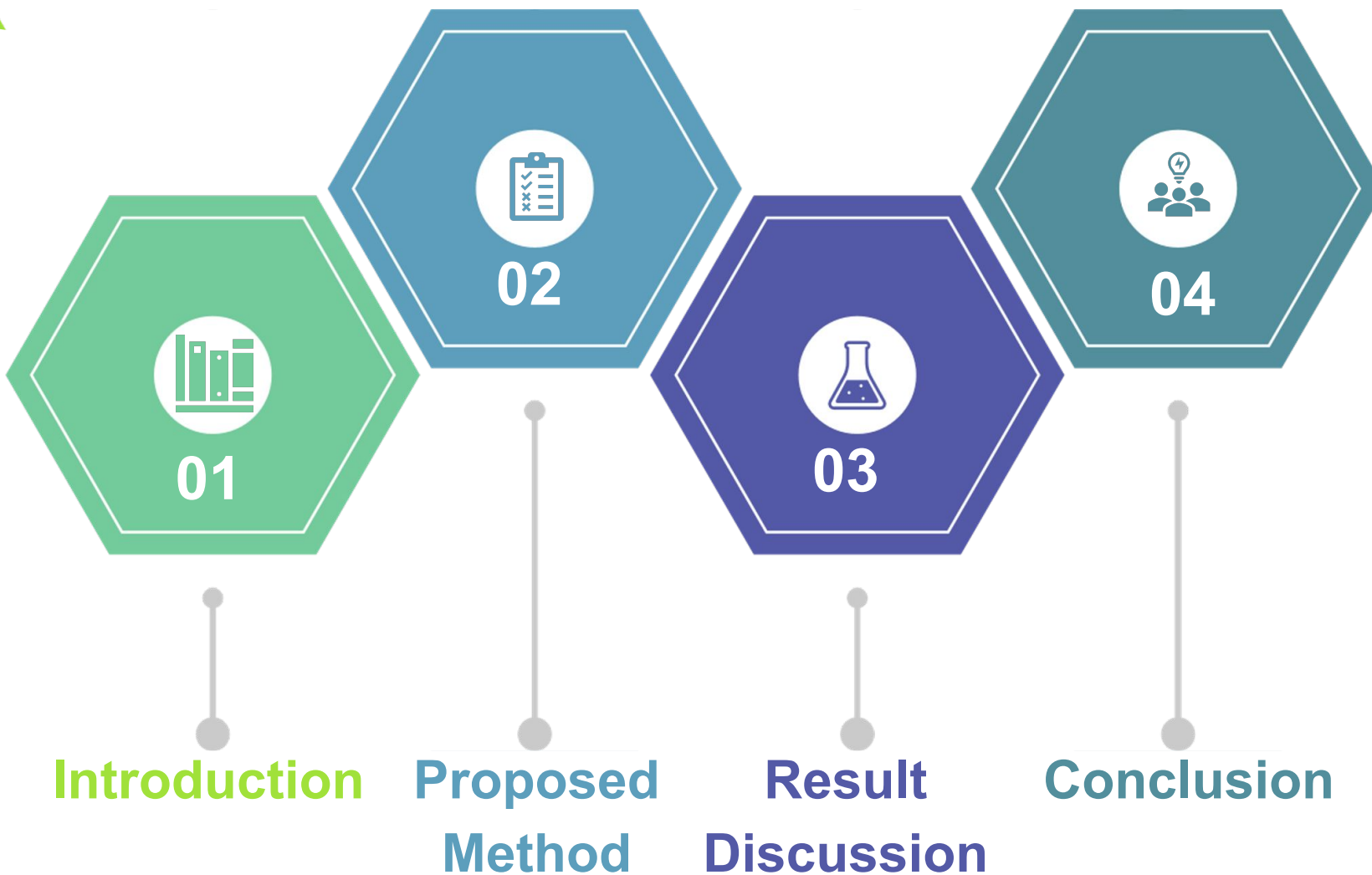
# AGENDA

**01**

**02**

**03**

**04**

**Introduction**

**Proposed Method**

**Result Discussion**

**Conclusion**

# 01

# Introduction

Trends in Chemistry

Wang, G., Ang, H. T., Dubbaka, S. R., O'Neill, P., & Wu, J. (2023). Multistep automated synthesis of pharmaceuticals. Trends in Chemistry.

**Data Extraction Module**

1. Data sources
2. Conectors
3. ETL Process
4. Processed data
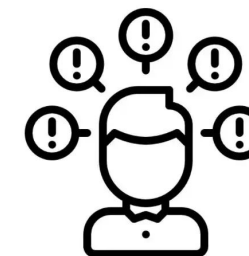
Insufficient data

## Bad data = Bad model

(bad data can mess up how companies decide things)

Incorrect algorithm selection
Incorrect hyperparameter tuning
Incorrect model deployment
Wrong evaluation metrics
Poorly collected requirements
...

**The effect of (bad) Data Quality on Model Accuracy in Supervised Machine Learning**
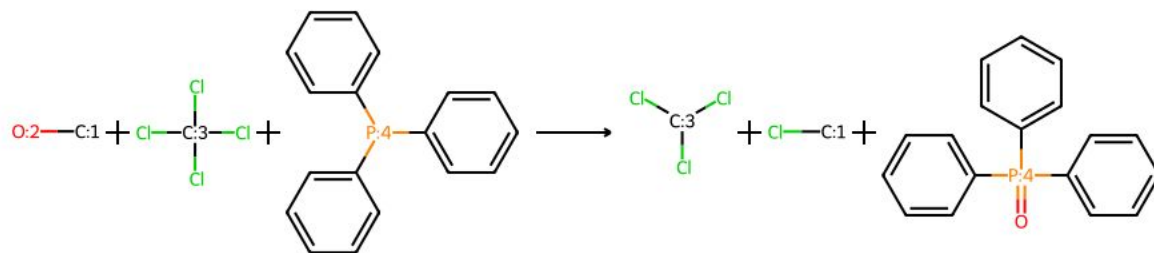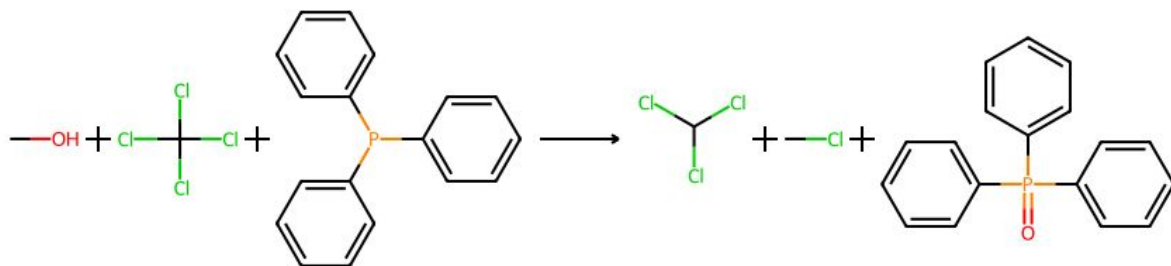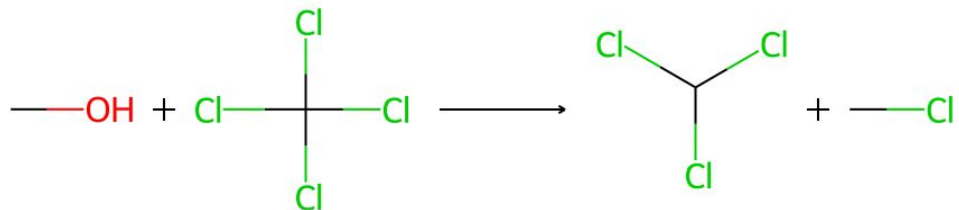
Saquicela, V., Baculima, F., Orellana, G., Piedra, N., Orellana, M., & Espinoza, M. (2018, March). Similarity Detection among Academic Contents through Semantic Technologies and Text Mining. In IWSW (pp. 1-12).

Rebalancing

Atom mapping
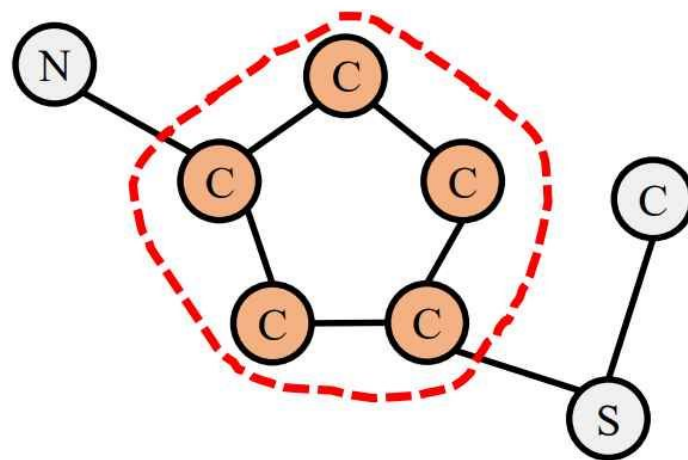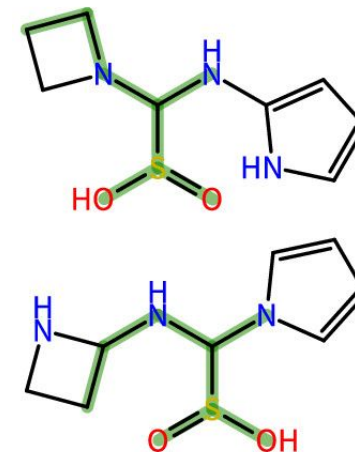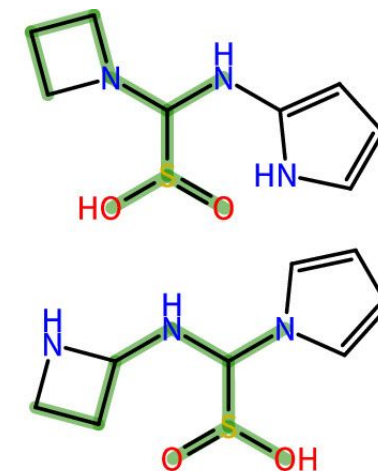
# Maximum-common-subgraph



$G_1$

$G_2$

(a) An example of the connected induced MCS (cMCIS).

(b) An example of the connected noninduced MCS (cMCES).

1. Bai, Y., Xu, D., Sun, Y., & Wang, W. (2021, July). Glsearch: Maximum common subgraph detection via learning to search. In International Conference on Machine Learning (pp. 588-598). PMLR.
2. Robert Schmidt, Florian Krull, Anna Lina Heinzke, and Matthias Rarey. Journal of Chemical Information and Modeling 2021 61 (1), 167-178 DOI: 10.1021/acs.jcim.0c00741
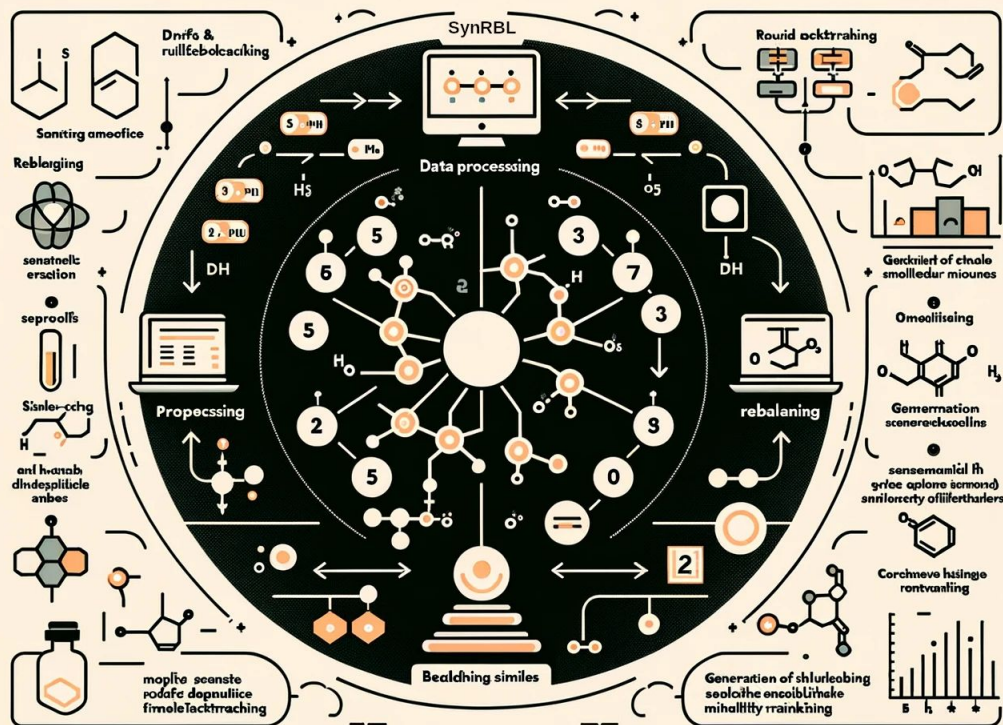
## SynRBL: Synthesis Rebalancing Framework

SynRBL (Synthesis Rebalancing Framework) is a specialized toolkit designed for computational chemistry. Its primary focus is on rebalancing incomplete chemical reactions and providing rule-based methodologies for data standardization and analysis.
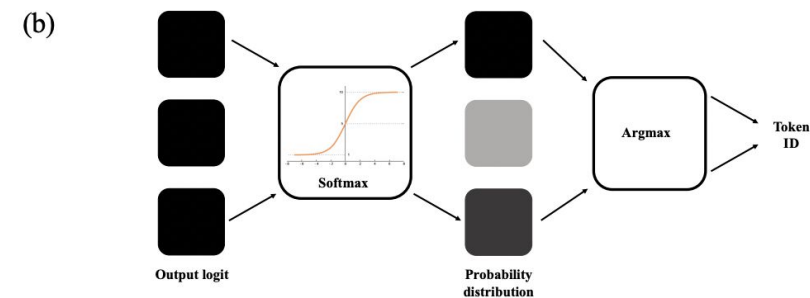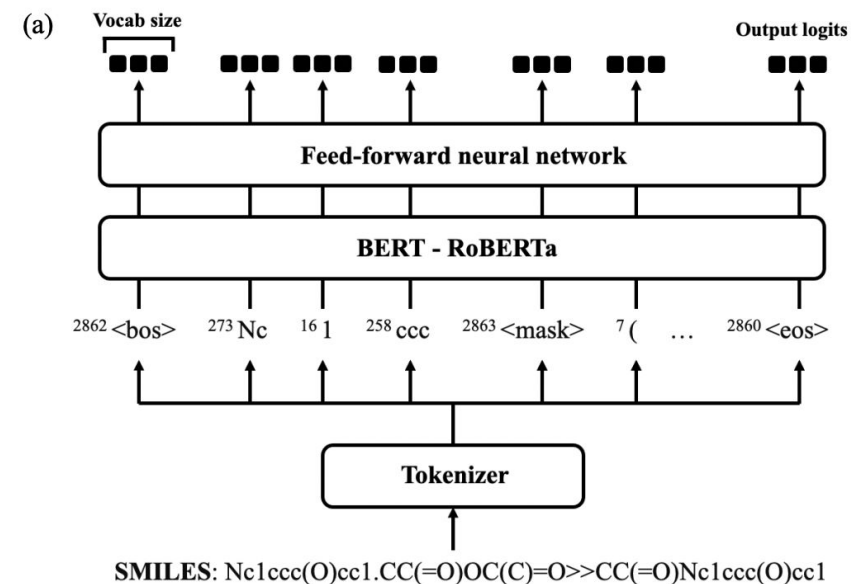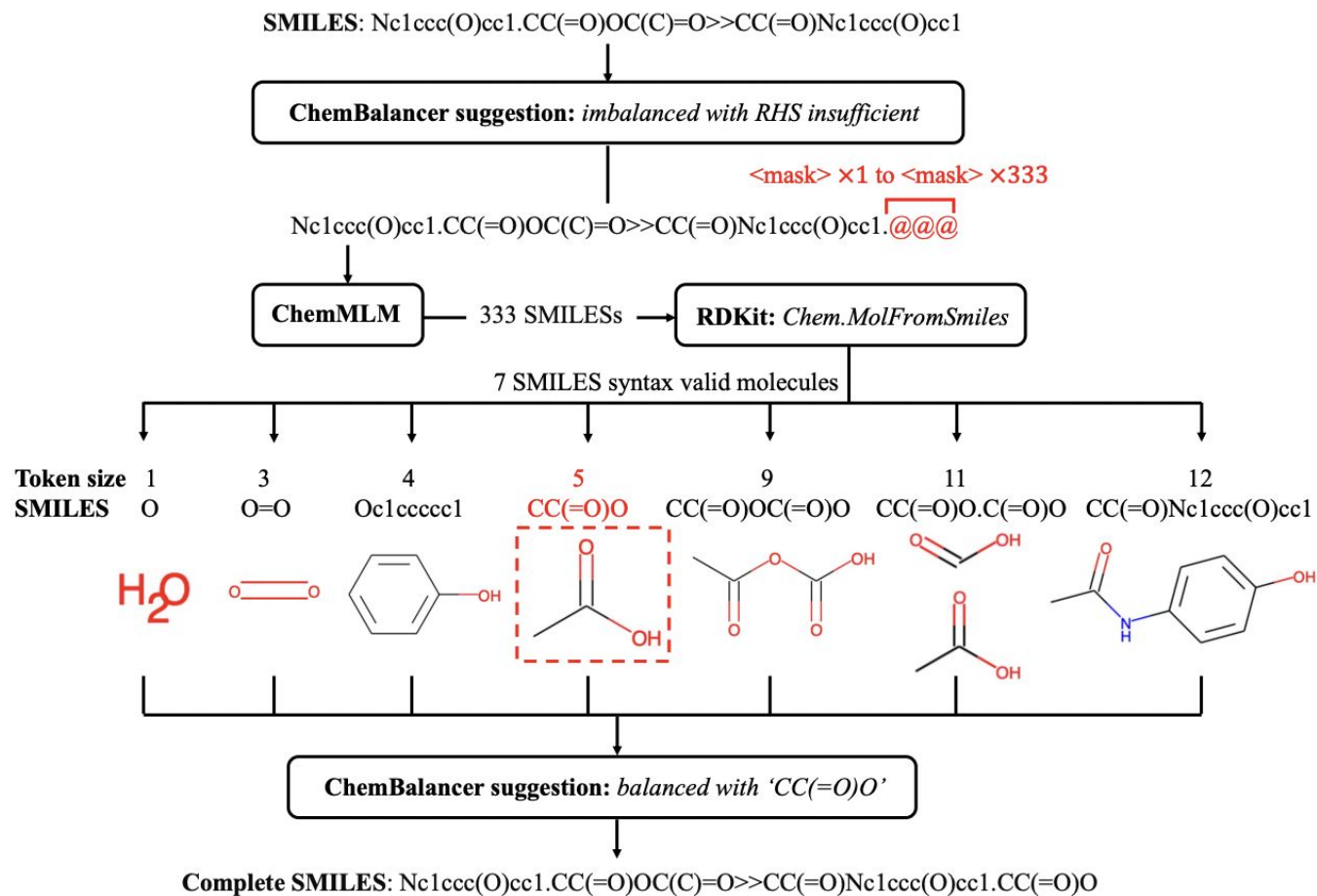
SMILES: Nc1ccc(O)cc1.CC(=O)OC(C)=O>>CC(=O)Nc1ccc(O)cc1

ChemBalancer suggestion: *imbalanced with RHS insufficient*

<mask> ×1 to <mask> ×333

Nc1ccc(O)cc1.CC(=O)OC(C)=O>>CC(=O)Nc1ccc(O)cc1.@@@

ChemMLM — 333 SMILESs → RDKit: *Chem.MolFromSmiles*

7 SMILES syntax valid molecules

| Token size | 1 | 3 | 4 | 5 | 9 | 11 | 12 |
| SMILES | O | O=O | Oc1ccccc1 | CC(=O)O | CC(=O)OC(=O)O | CC(=O)O.C(=O)O | CC(=O)Nc1ccc(O)cc1 |

ChemBalancer suggestion: *balanced with 'CC(=O)O'*

Complete SMILES: Nc1ccc(O)cc1.CC(=O)OC(C)=O>>CC(=O)Nc1ccc(O)cc1.CC(=O)O

(a) Vocab size ... Output logits

Feed-forward neural network

BERT - RoBERTa

2862 <bos>  273 Nc  16 1  258 ccc  2863 <mask>  7 (  ...  2860 <eos>

Tokenizer

SMILES: Nc1ccc(O)cc1.CC(=O)OC(C)=O>>CC(=O)Nc1ccc(O)cc1

(b)

Output logit — Softmax — Probability distribution — Argmax — Token ID

Zhang, C., Arun, A., & Lapkin, A. (2023). Completing and balancing database excerpted chemical reactions with a hybrid mechanistic-machine learning approach.

# 02

# METHOD

Flowchart:

Reaction → $n_{C,R} = n_{C,P}$

- No → (to MCS-based method path)
- Yes → Category?
  - Balance → Result
  - Bothside → (path down)
  - Reactant or Product → Rule-based method → Balanced?
    - No → MCS-based method → Rule-based method → Result
    - Yes → Result



The Golden Dataset

UNITED STATES PATENT AND TRADEMARK OFFICE
uspto

11

# Rule-based approach

| Molecular | Empirical |
|-----------|-----------|
| $N_2O_4$ | $NO_2$ |
| $C_6H_6$ | $CH$ |
| $C_2H_6O_2$ | $CH_3O$ |

Molecular Representation
Eg: $CH_3CHOOH$
{C:2, H:4, O : 2, Q : 0}.

# Rule-based approach



What is the current scale and comprehensiveness of the template library within this context?

# Rule-based approach



**DFS search**

# Case Study



'Unbalance': 'Products'

'Diff_formula': {'S': 1, 'O': 3, 'H': 1, 'Q': -1},

# Case Study



Check length: 4

Search rules from length 4: $SO_3^{2-}$ {'S': 1, 'O': 3, 'Q': -2}

Substrate: {'H':1, 'Q':+1}

Check length: 2

Search rules with length 2: $H^+$ {'H':1, 'Q':1}

## MCS-based approach

# MCS-based approach

Reactants



Products

Reactants

Products

Reactants

Products

Reactants

Products

$\Longrightarrow$

Expand Rule

Reactants

Products

$\Longrightarrow$

Expand Rule

Merge Rule

Reactants

Products

$\Longrightarrow$

Expand Rule    Merge Rule

Reactants

Products

$\Longrightarrow$

$+ H_2O$

Expand Rule

Merge Rule

# 03

# RESULT - DISCUSSION

# Rule-based approach

# MCS-based approach

## Limitations

**04**

# CONCLUSION

# Thank you for your attendance

UNIVERSITÄT LEIPZIG

# Appendix

# Equivariant Isomorphism

# Applicability Domain



**A** Confusion Matrix

**B** Classification Report

**C** ROC Curve

**D** Precision-Recall Curve

Original set — SMOTE + Tomek

Class #0   Class #1

1.   Chen, T., He, T., Benesty, M., Khotilovich, V., Tang, Y., Cho, H., ... & Zhou, T. (2015). Xgboost: extreme gradient boosting. R package version 0.4-2, 1(4), 1-4.
2.   LemaÃŽtre, G., Nogueira, F., & Aridas, C. K. (2017). Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. Journal of machine learning research, 18(17), 1-5.
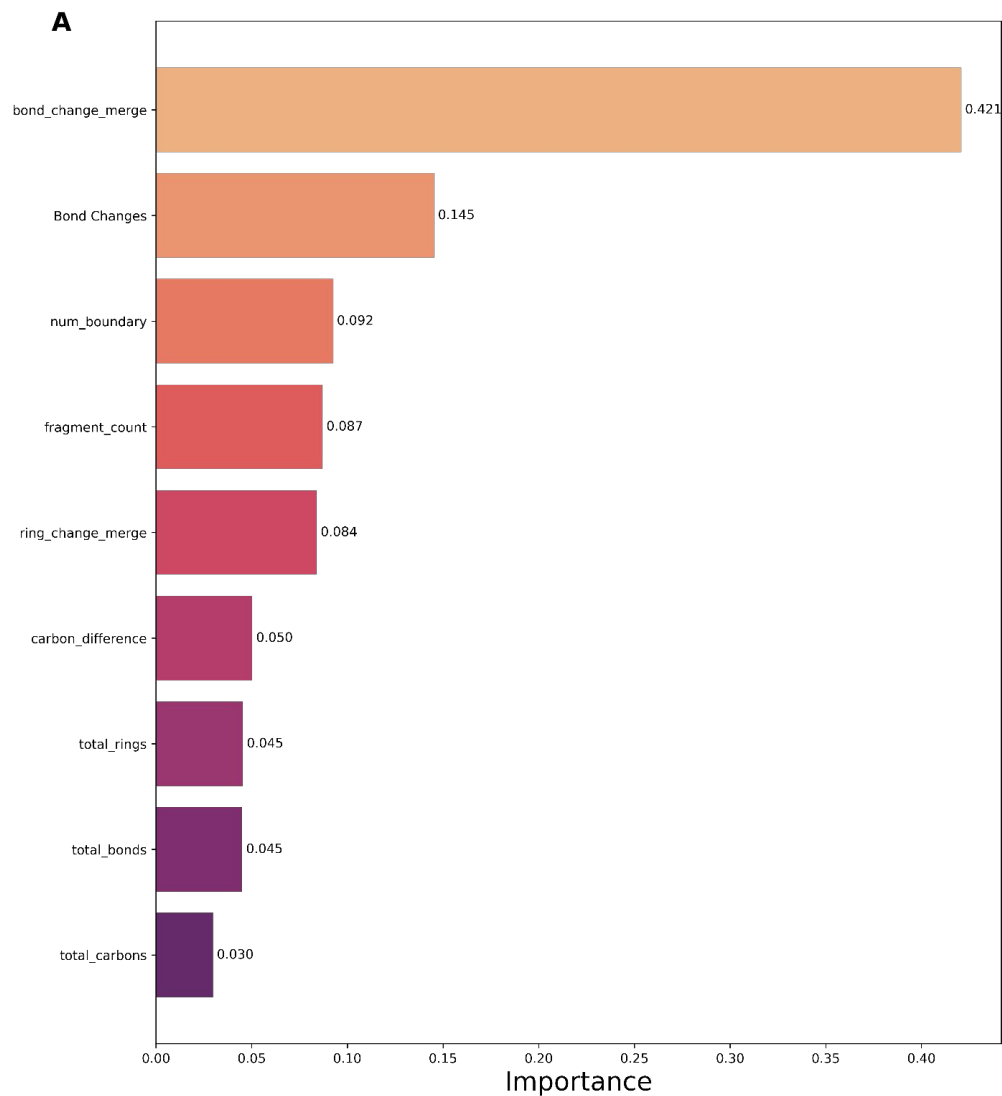
Thank you for your attention!

# INTRODUCTION

| INTRODUCTION | METHOD | RESULT-DISCUSSION | CONCLUSION |

| INTRODUCTION | METHOD | RESULT-DISCUSSION | CONCLUSION |

| INTRODUCTION | METHOD | RESULT-DISCUSSION | CONCLUSION |

Zhang, C., Arun, A., & Lapkin, A. (2023). Completing and balancing database excerpted chemical reactions with a hybrid mechanistic-machine learning approach.