



TACsy

Training Alliance for
Computational systems
chemistry

SynCat: A light-weight model for reaction classification

Winterseminar-Bled

Presenter: Tieu-Long Phan

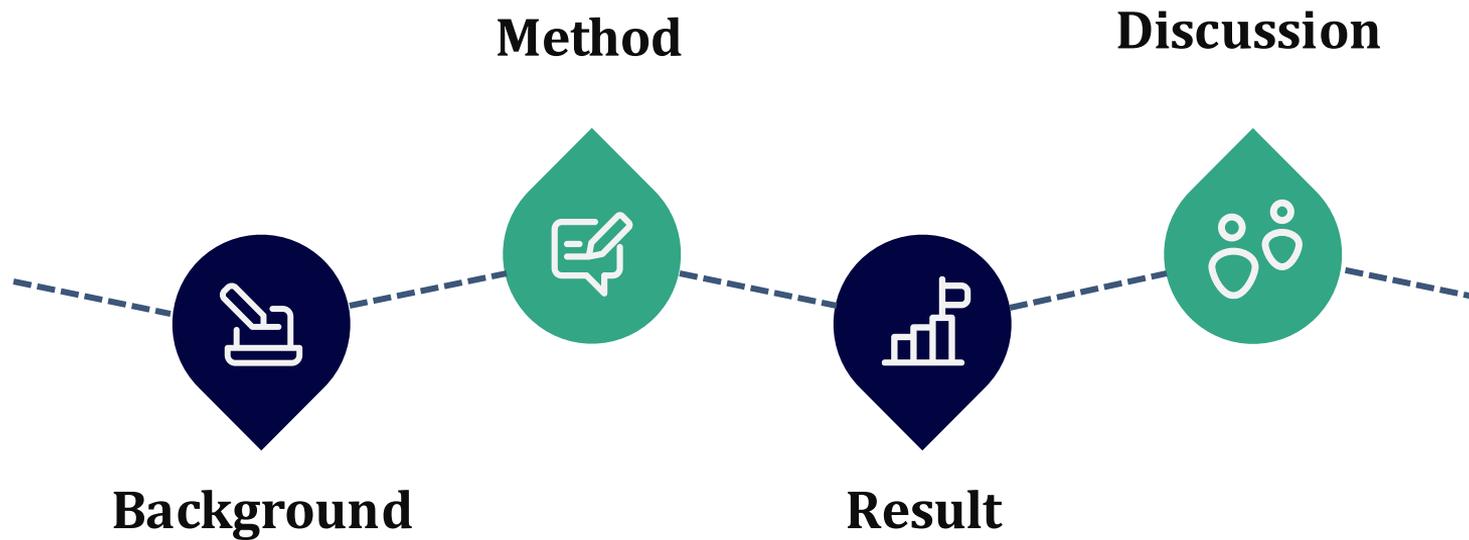
Date: 10.02.2025



Founded by the
European Union

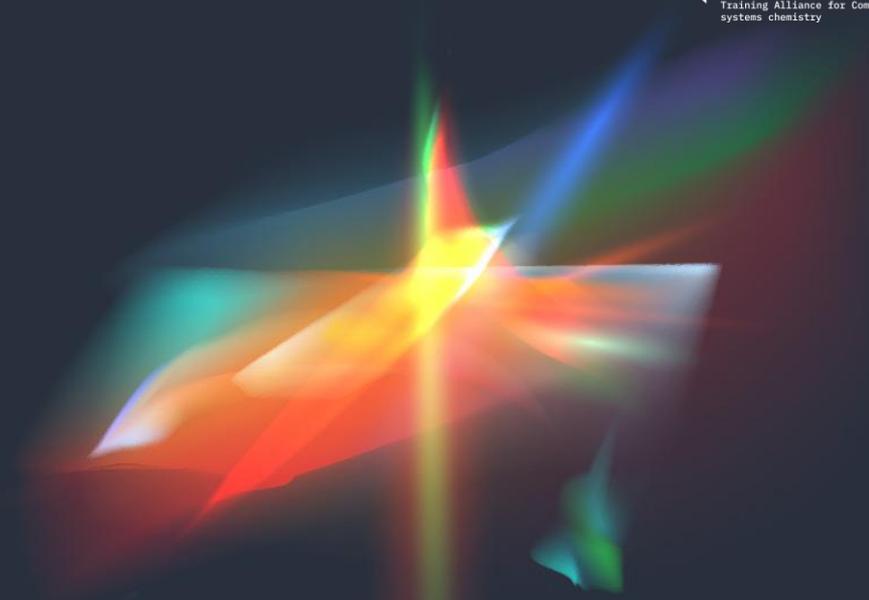
This project has received funding from the European Union's Horizon 2021 research and innovation programme under the Marie-Sklodowska-Curie grant agreement No 101072930

OUTLINE



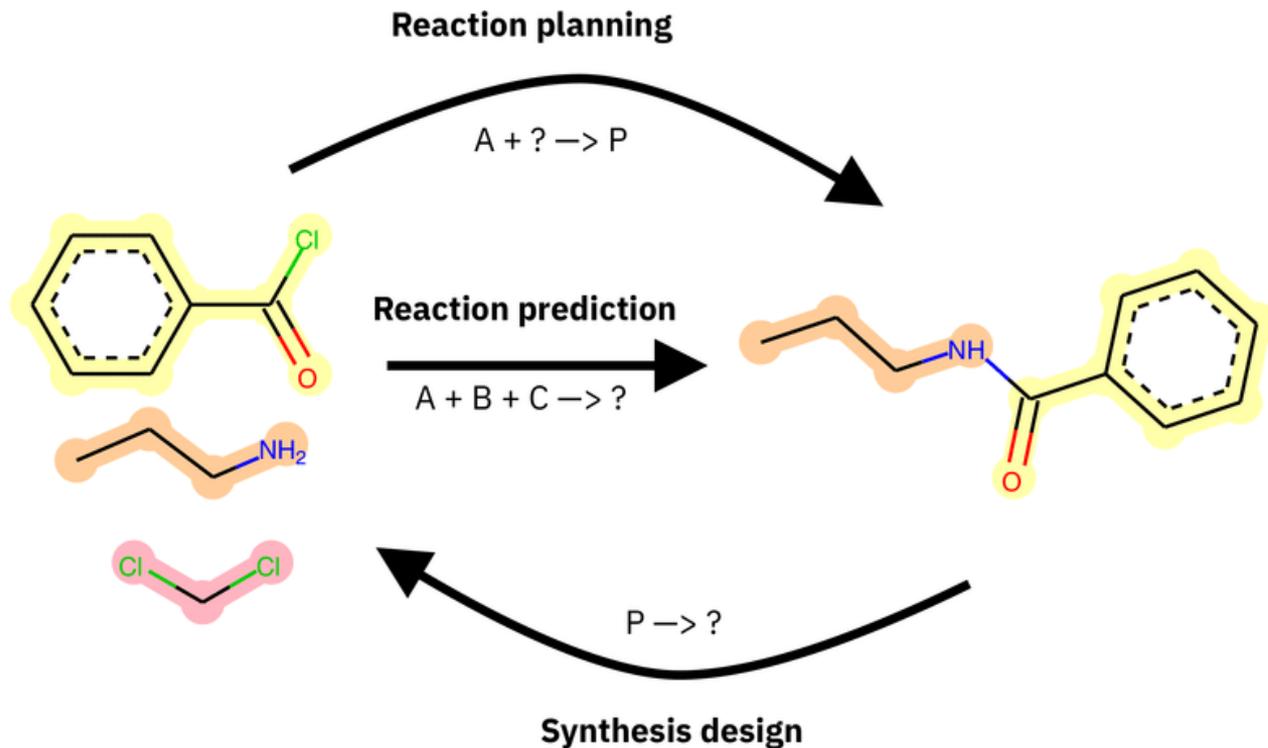
01

BACKGROUND



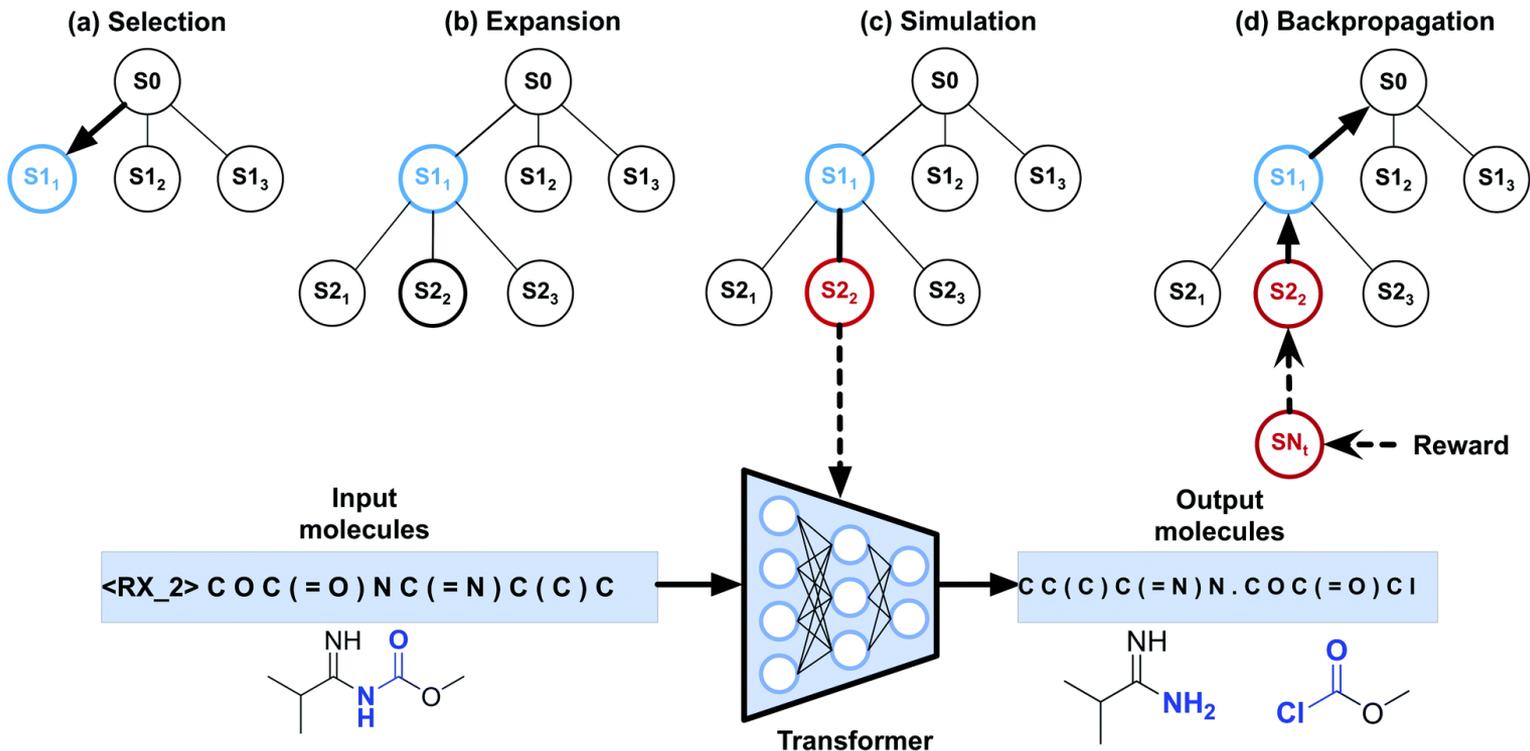
BACKGROUND

→ *Synthesis Planning*



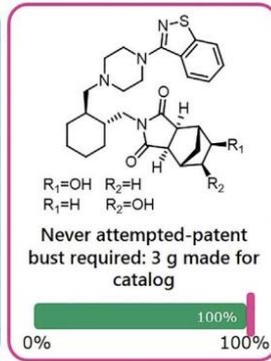
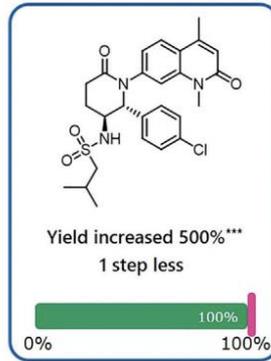
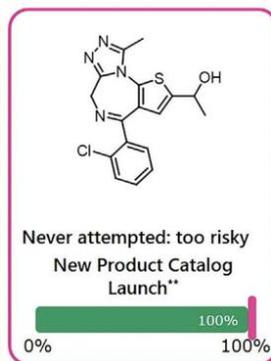
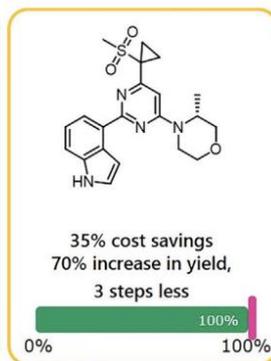
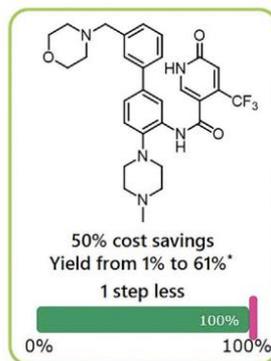
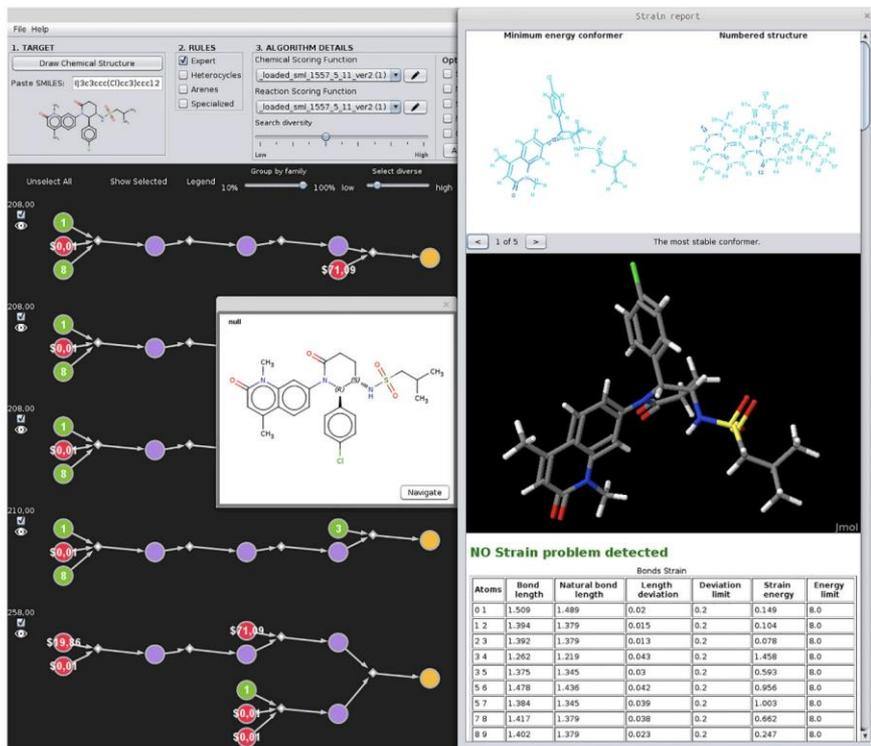
BACKGROUND

→ *Template-free synthesis planning*



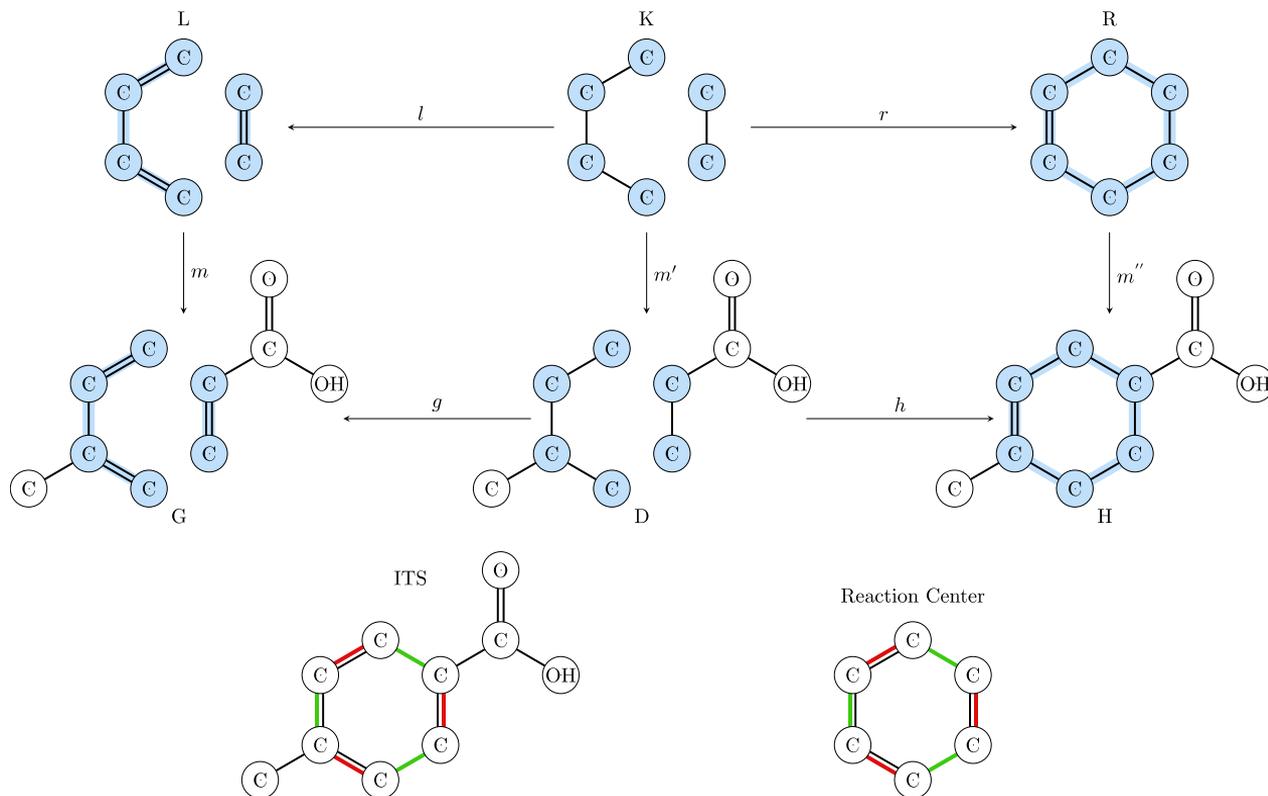
BACKGROUND

Template-based synthesis planning



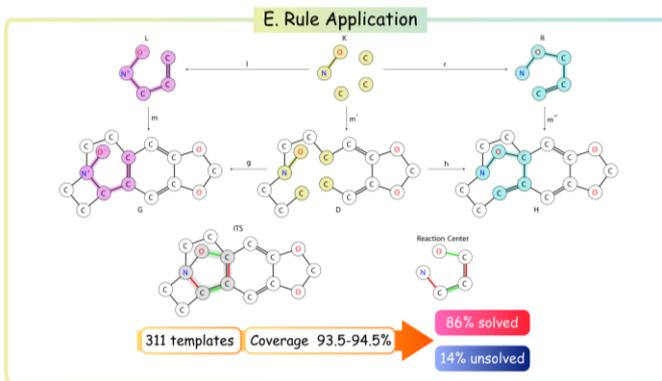
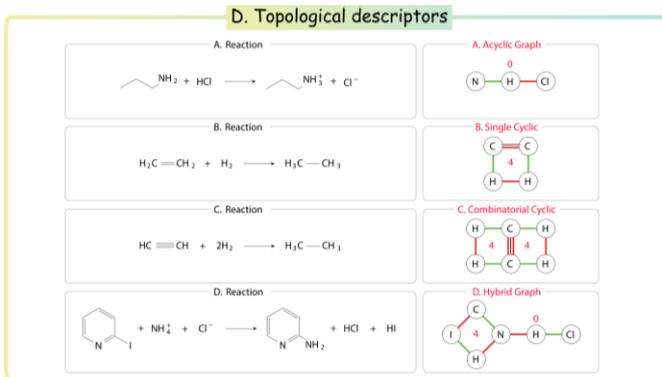
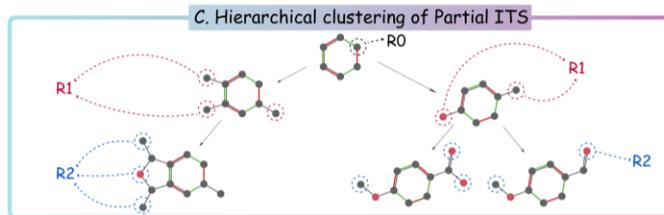
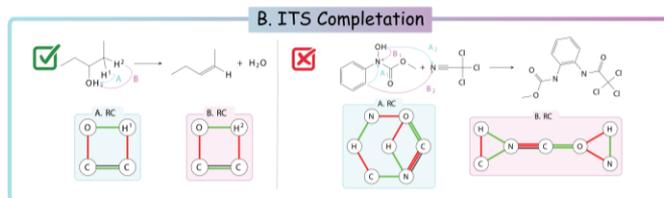
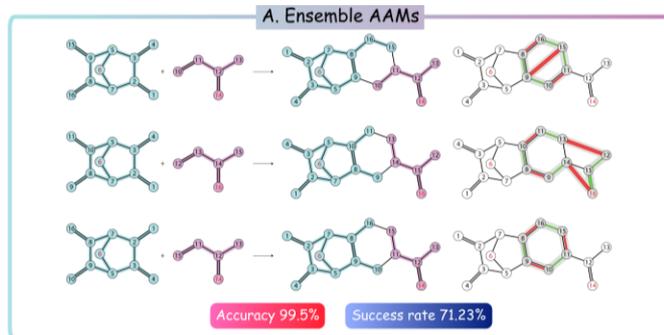
BACKGROUND

→ Double Pushout Rules



BACKGROUND

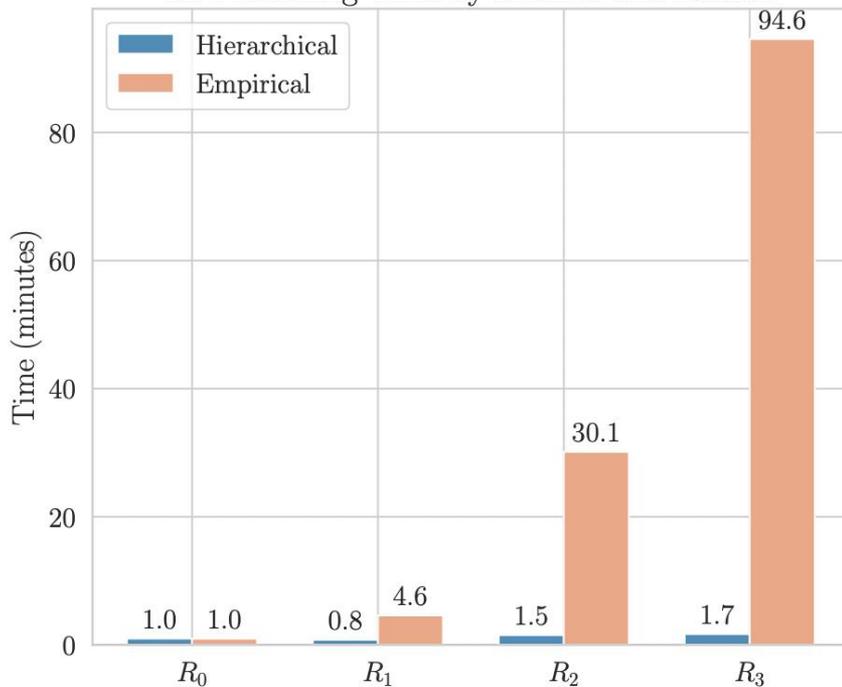
→ *SynTemp*



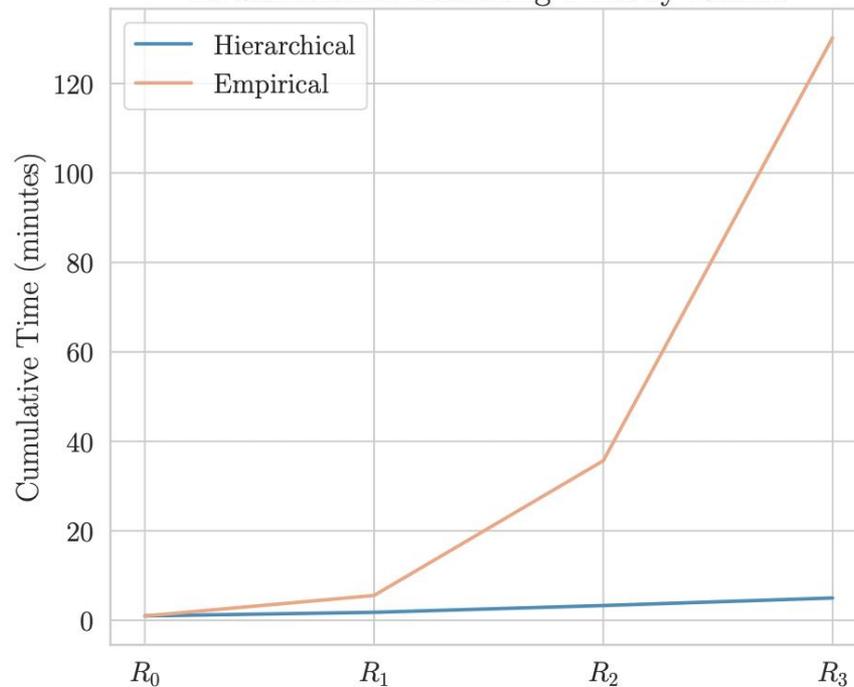
BACKGROUND

→ SynTemp - Challenge

A. Processing Time by Method and Radius

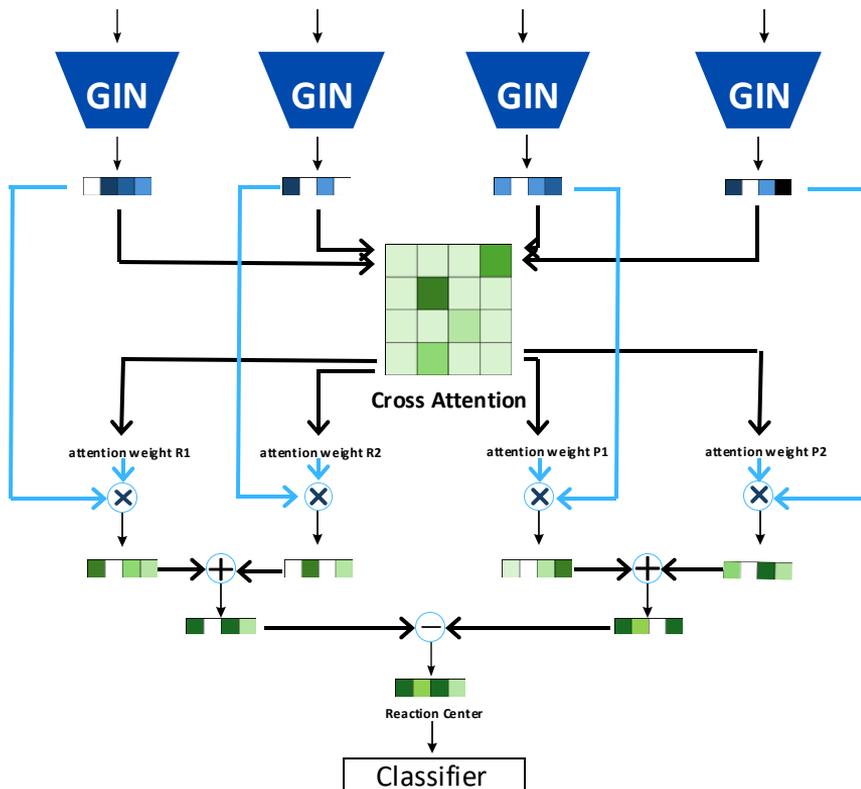


B. Cumulative Processing Time by Radius



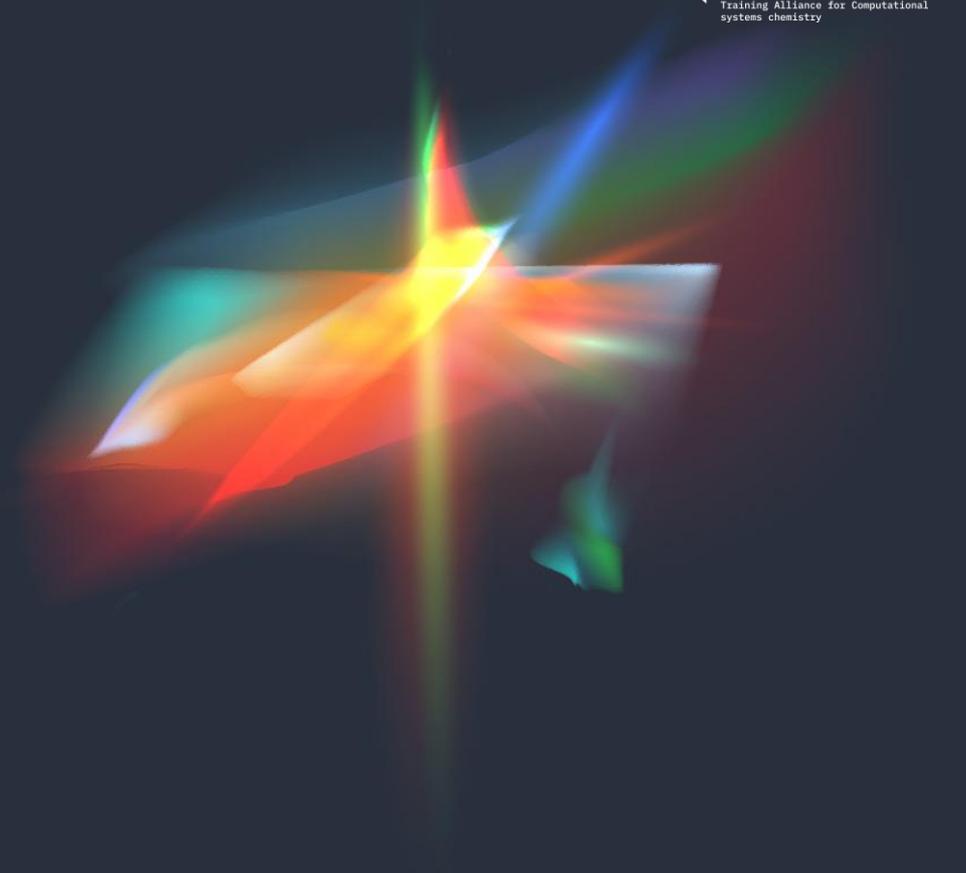
BACKGROUND

→ *SynCat*



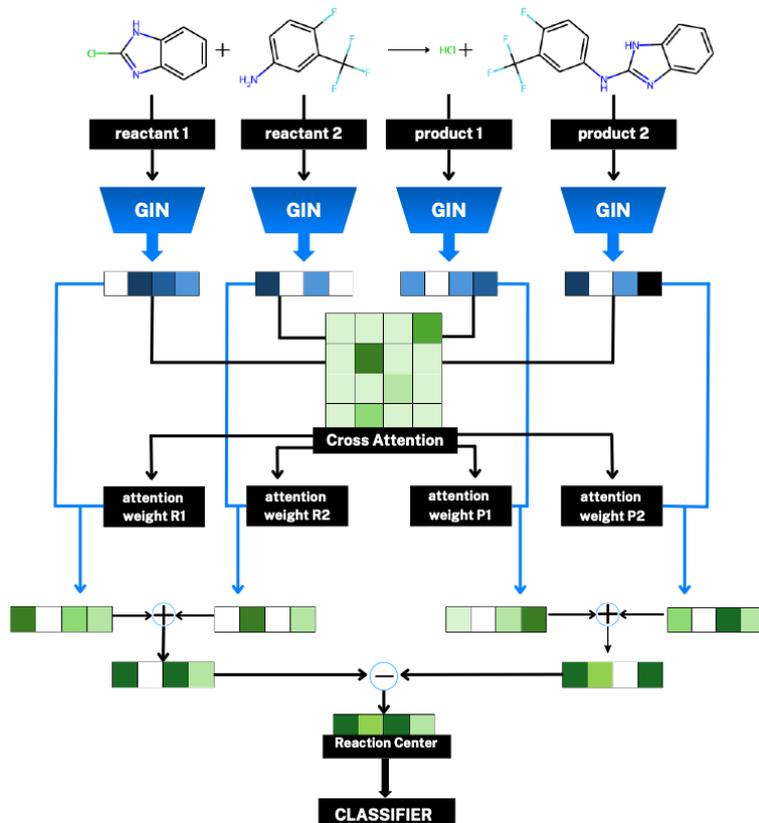
02

METHOD

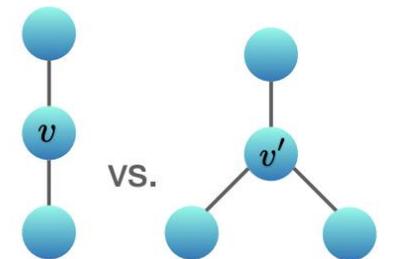


METHOD

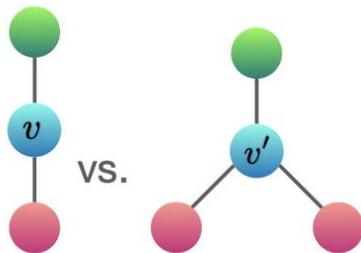
MODEL ARCHITECTURE



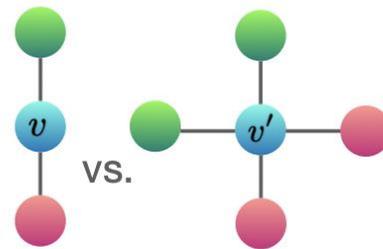
→ Graph isomorphism network



(a) Mean and Max both fail



(b) Max fails



(c) Mean and Max both fail

We utilize the Graph Isomorphism Network (GIN) to enhance graph representation accuracy, employing learnable transformations:

$$\mathbf{h}_v^{(0)} = \phi_n(\mathbf{v}), \quad \mathbf{h}_e^{(k)} = \phi_e(\mathbf{e}^{(k)}), \quad (1)$$

→ Graph isomorphism network

Vertex embeddings update through layers, integrating local and neighboring data:

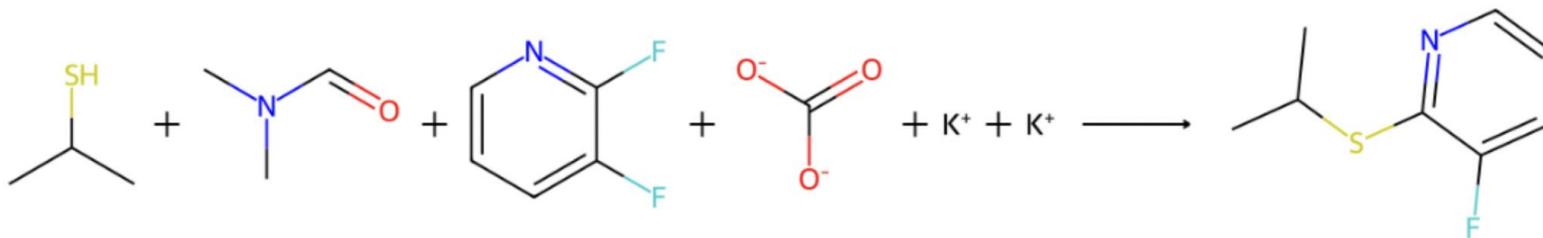
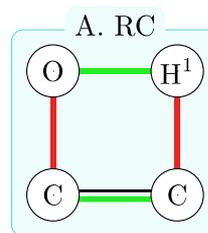
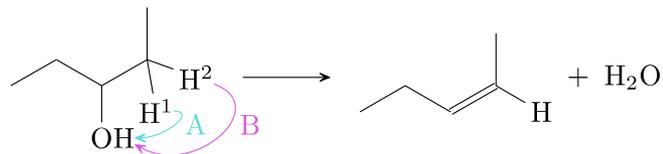
$$\mathbf{h}_v^{j,(l)} = \psi^{(l)} \left(\mathbf{h}_v^{j,(l-1)} + \sum_{k|e^{j,k} \in E} \text{ReLU}(\mathbf{h}_v^{j,(l-1)} + \mathbf{h}_e^{j,k}) \right), \quad (2)$$

culminating in a comprehensive graph embedding:

$$\mathbf{h}_g = \sum_{j|\mathbf{v}^j \in V} \mathbf{h}_v^{j,(L)}. \quad (3)$$

METHOD

→ Cross-attention



→ *Cross-attention*

The reaction center, Γ , is calculated by subtracting the aggregated product embeddings from the reactant embeddings:

$$\Gamma = \mathcal{R} - \mathcal{P}$$

where \mathcal{R} and \mathcal{P} denote the embeddings of reactants and products, respectively. This difference highlights the differential characteristics crucial for reaction classification.

→ *Cross-attention*

To address noise from redundant substances, we implement cross-attention mechanisms that assign significance to compounds based on their contribution to the reaction center. The matrices \mathcal{R}_{org} and \mathcal{P}_{org} represent original reactant and product embeddings, from which we derive:

$$\mathcal{R}_{\text{updated}} = \begin{bmatrix} \mathcal{R}_{\text{org}} \\ S_{\mathcal{R}_{i,j}} \end{bmatrix}, \quad \mathcal{P}_{\text{updated}} = \begin{bmatrix} \mathcal{P}_{\text{org}} \\ S_{\mathcal{P}_{i,j}} \end{bmatrix}$$

Pairwise sums are integrated to enhance embeddings.

→ *Cross-attention*

To address noise from redundant substances, we implement cross-attention mechanisms that assign significance to compounds based on their contribution to the reaction center. The matrices \mathcal{R}_{org} and \mathcal{P}_{org} represent original reactant and product embeddings, from which we derive:

$$\mathcal{R}_{\text{updated}} = \begin{bmatrix} \mathcal{R}_{\text{org}} \\ S_{\mathcal{R}_{i,j}} \end{bmatrix}, \quad \mathcal{P}_{\text{updated}} = \begin{bmatrix} \mathcal{P}_{\text{org}} \\ S_{\mathcal{P}_{i,j}} \end{bmatrix}$$

Pairwise sums are integrated to enhance embeddings.

→ *Cross-attention*

Updated embeddings are processed to generate query and key vectors for attention calculations:

$$Q_{\mathcal{R}} = \mathcal{R}_{\text{updated}} \cdot W_q, \quad K_{\mathcal{R}} = \mathcal{R}_{\text{updated}} \cdot W_k$$

$$Q_{\mathcal{P}} = \mathcal{P}_{\text{updated}} \cdot W'_q, \quad K_{\mathcal{P}} = \mathcal{P}_{\text{updated}} \cdot W'_k$$

Attention weights are computed to emphasize the contributions of specific compounds.

→ *Cross-attention*

The final embeddings for reactants and products are derived by applying attention weights:

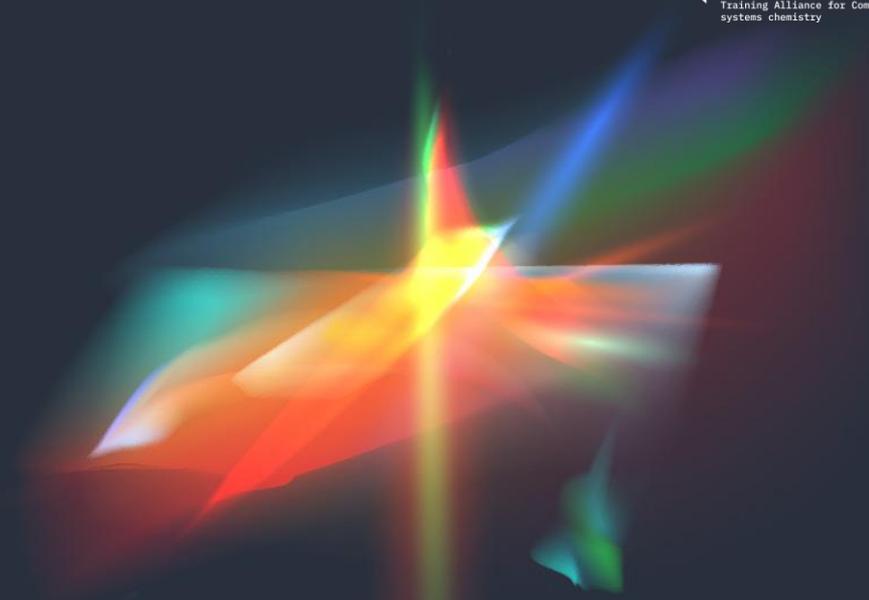
$$\mathcal{R} = \sum_{i=0}^{n_r-1} \mathcal{R}_{\text{updated}}[i, :] \cdot A_{\mathcal{R}_{\text{avg}}}[i], \quad \mathcal{P} = \sum_{i=0}^{n_p-1} \mathcal{P}_{\text{updated}}[i, :] \cdot A_{\mathcal{P}_{\text{avg}}}[i]$$

The reaction center Γ is then used in a neural network with softmax activation for classification, optimizing using cross-entropy loss:

$$\mathcal{L} = - \sum_{i=1}^N \sum_{c=1}^C y_{i,c} \log(p_{i,c})$$

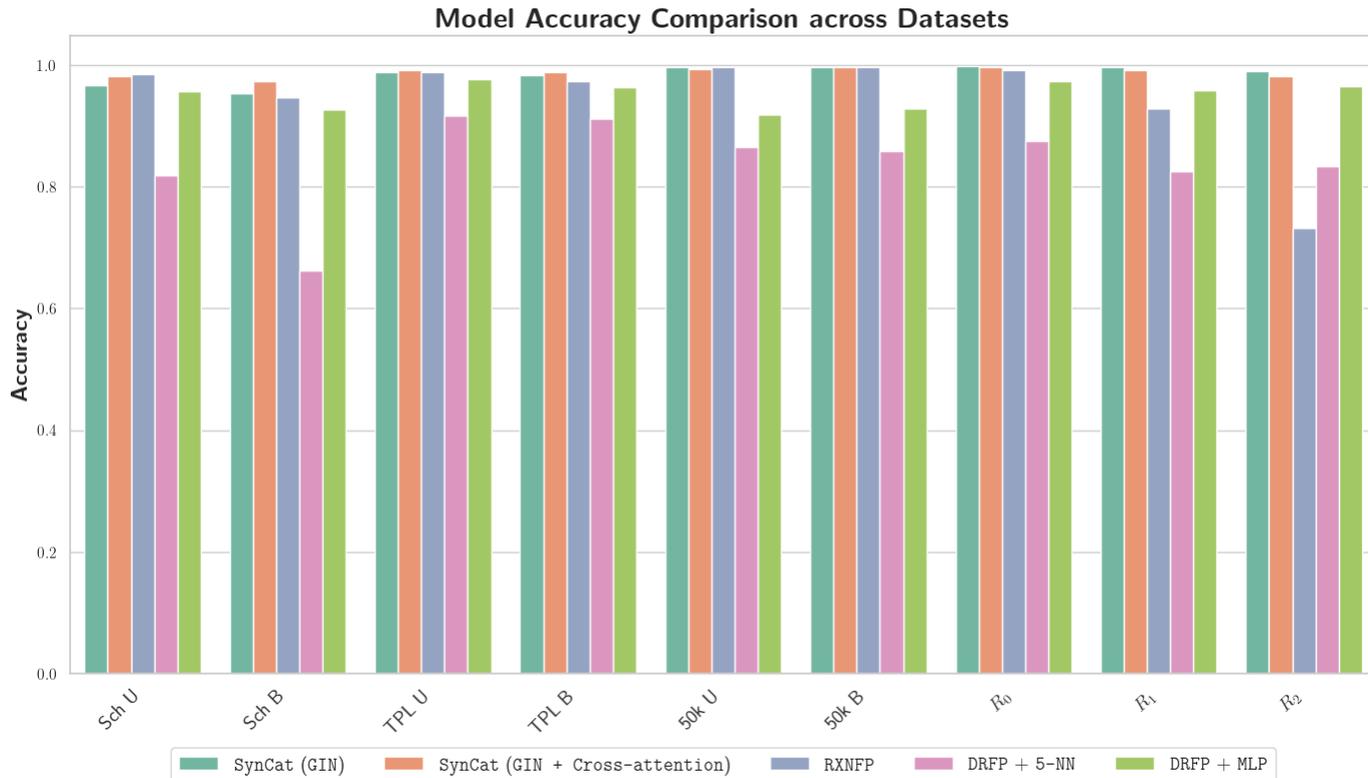
03

RESULT



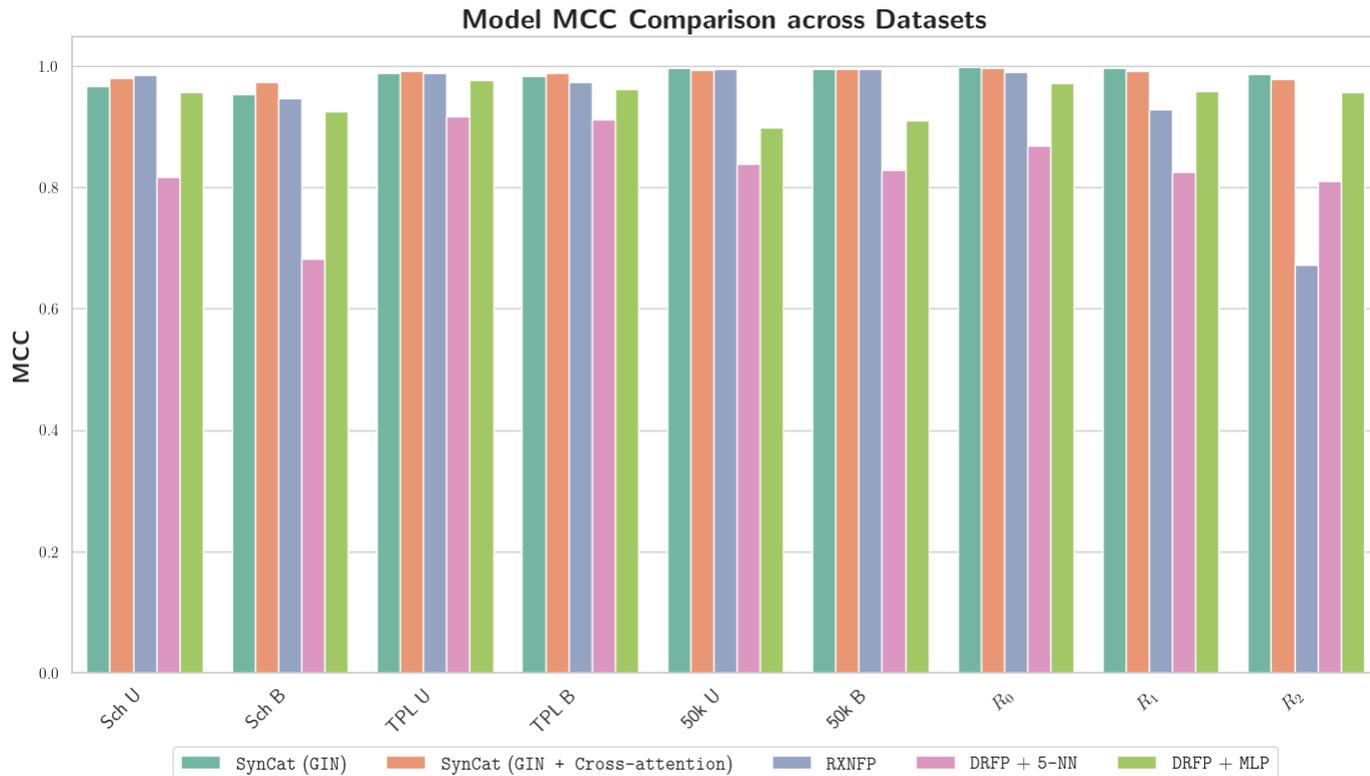
RESULT

Benchmarking



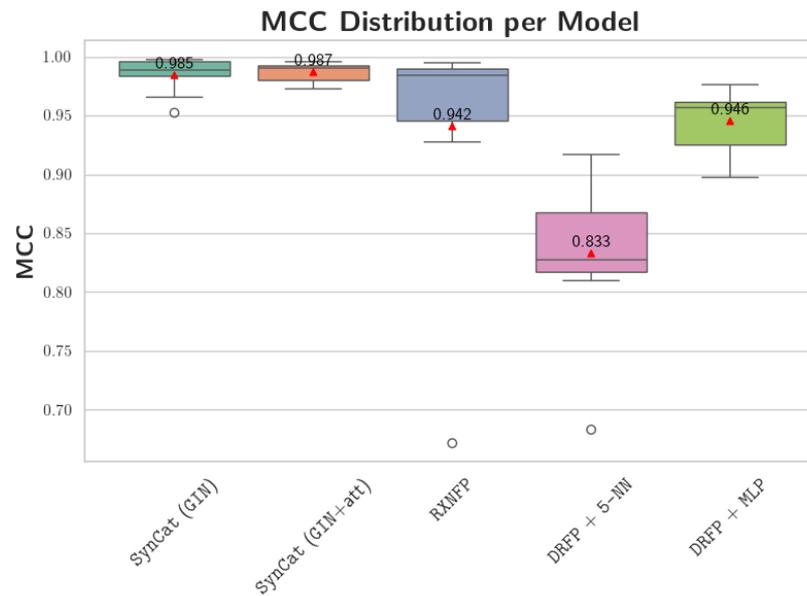
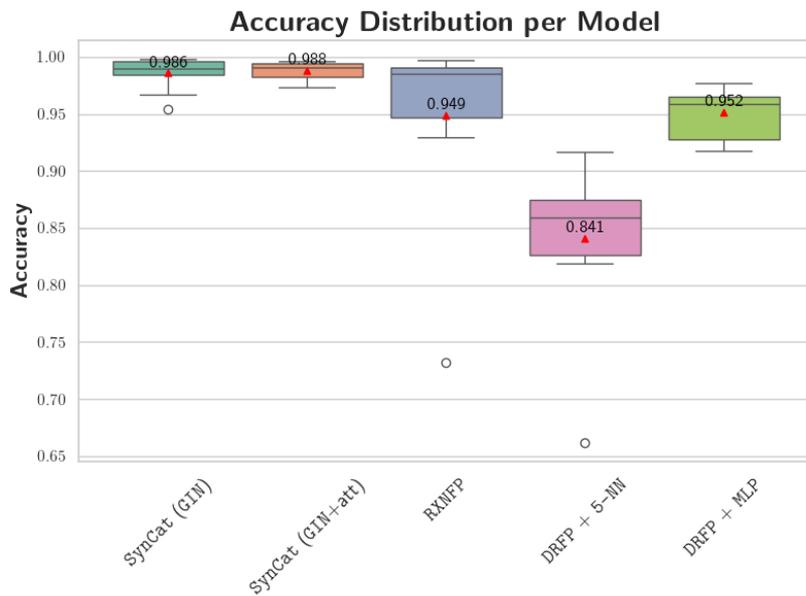
RESULT

Benchmarking



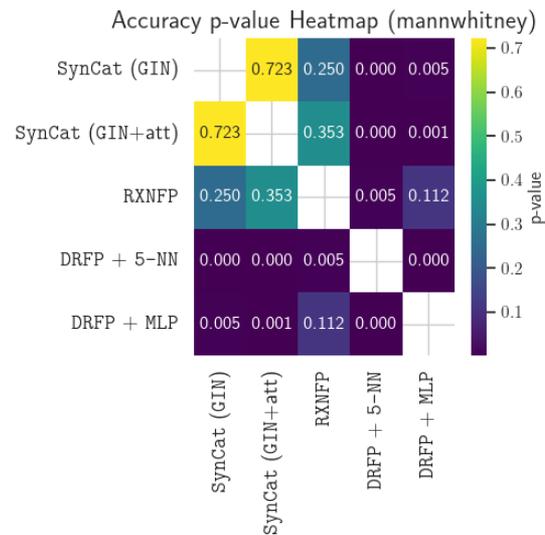
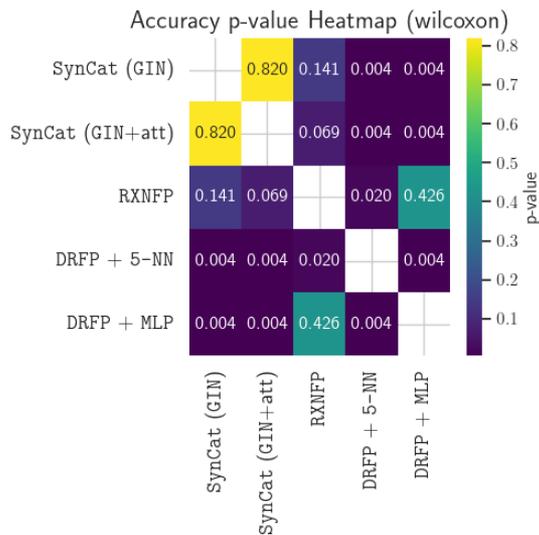
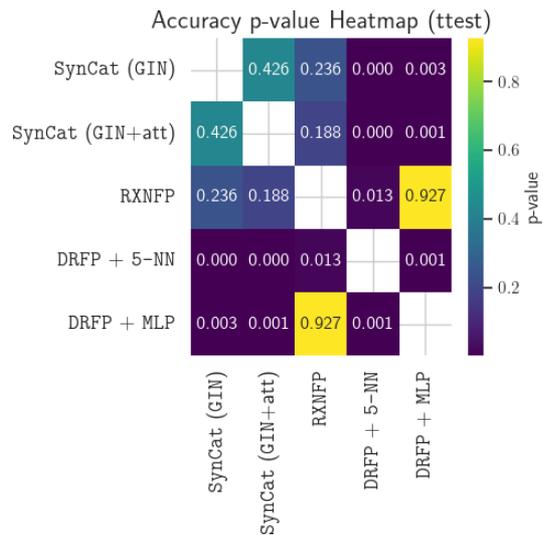
RESULT

Benchmarking



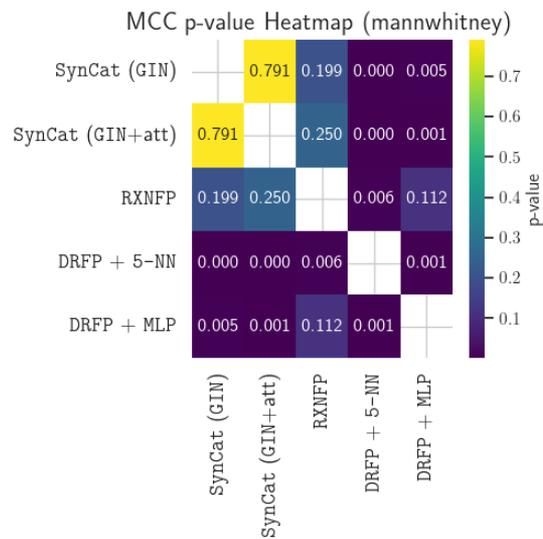
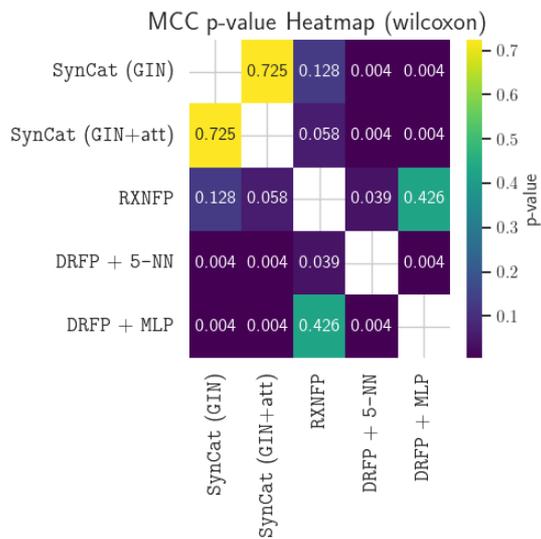
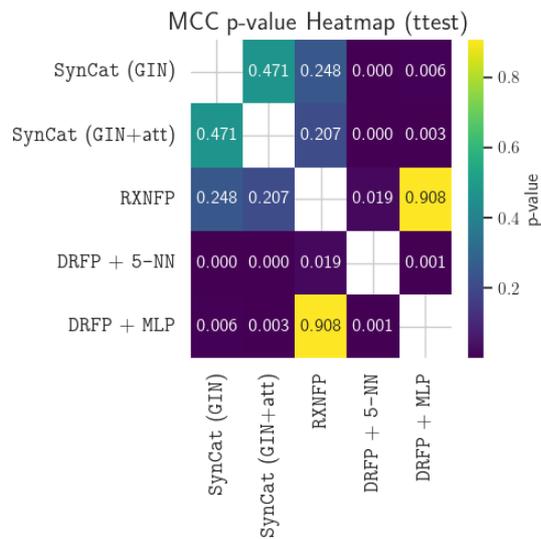
RESULT

Benchmarking



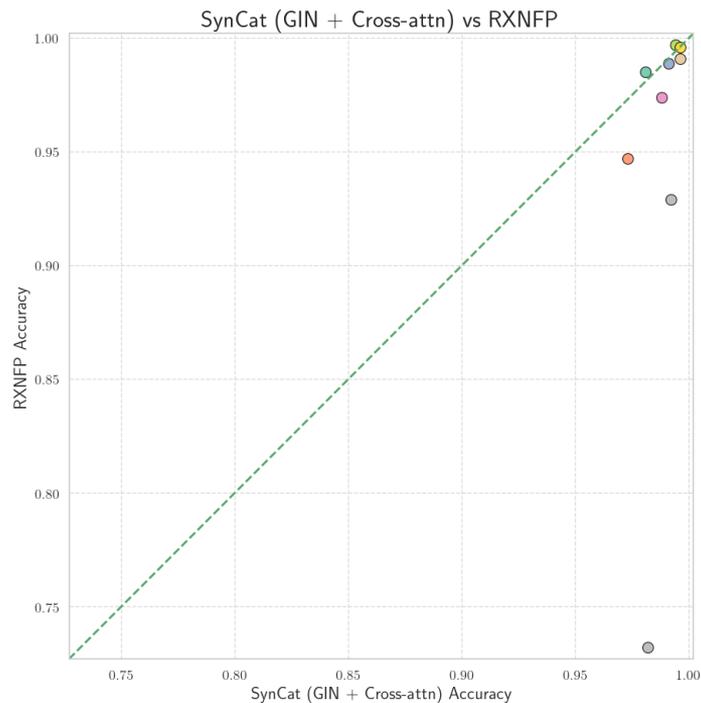
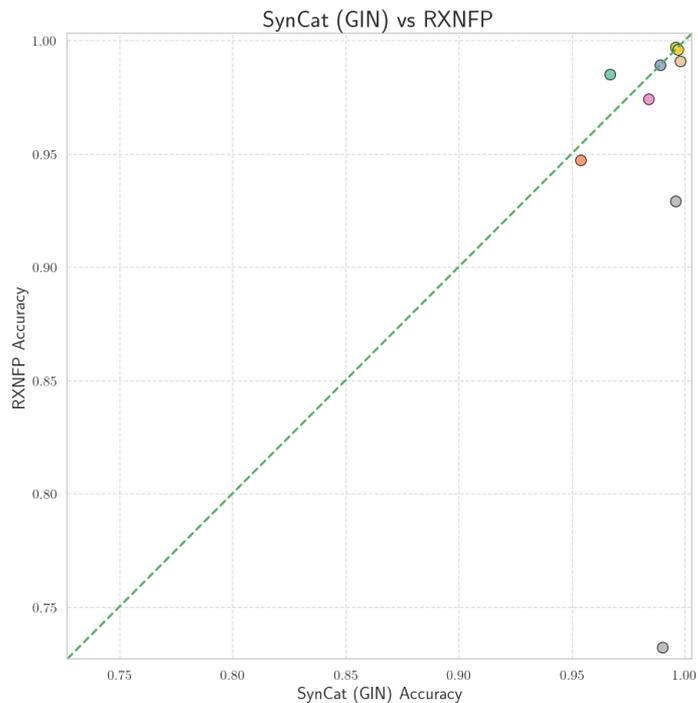
RESULT

Benchmarking



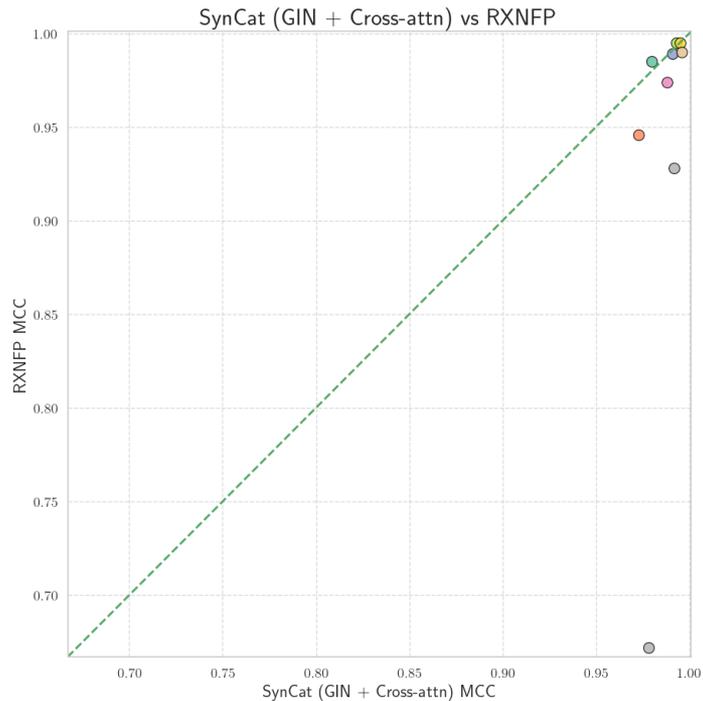
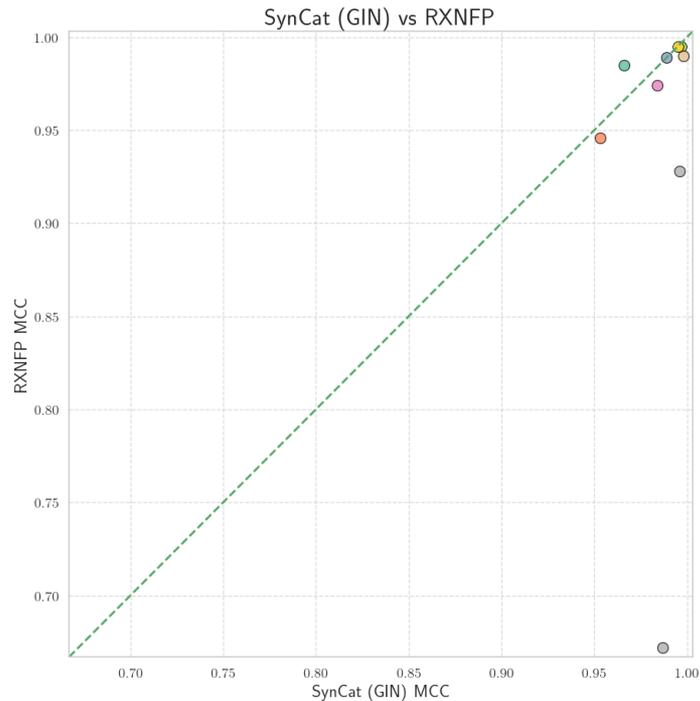
RESULT

Benchmarking



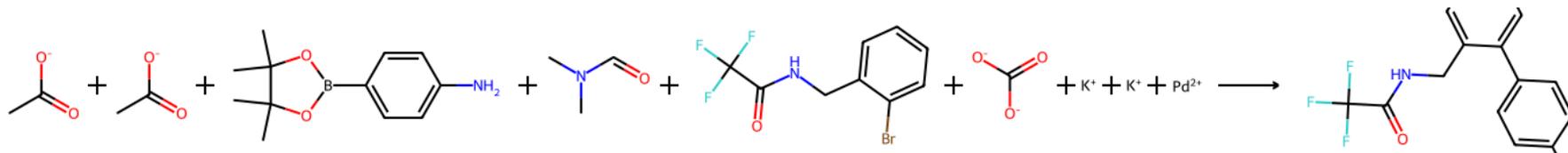
RESULT

Benchmarking



RESULT

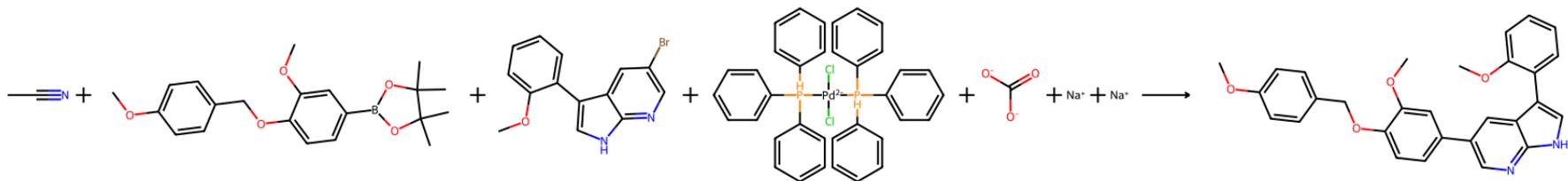
Benchmarking



Dataset: Schneider

Label: Bromo Suzuki Coupling

Predict: Bromo Suzuki-type Coupling



Dataset: Schneider

Label: Bromo Suzuki Coupling

Predict: Bromo Suzuki-type Coupling

RESULT

→ Reagent detection

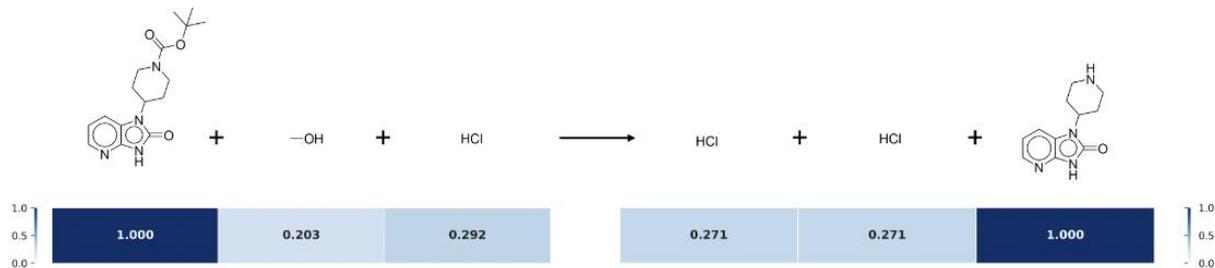
Dataset: USPTO TPL

TPL label 43: [C:4][S:1](=[O:2])(=[O:3])Cl.[C:5][N:6]>>[C:5][N:6][S:1](=[O:2])(=[O:3])



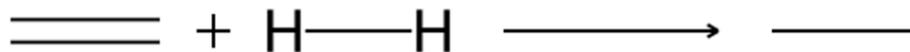
Dataset: Schneider

Label: N-Boc deprotection



RESULT

→ *Reagent detection*



04

DISCUSSION

DISCUSSION

→ Database Organization

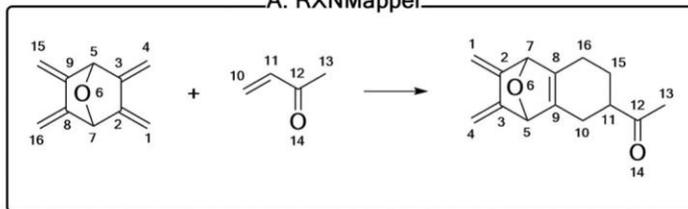


SynCat + SynTemp

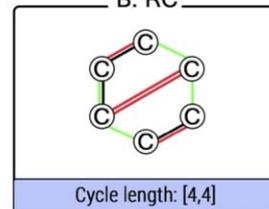
DISCUSSION

→ Inequivalent AAMs

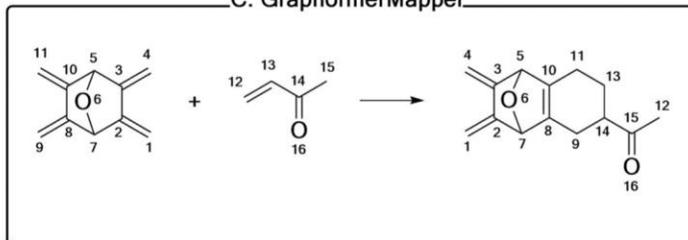
A. RXNMapper



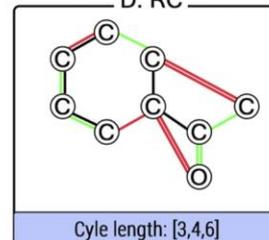
B. RC



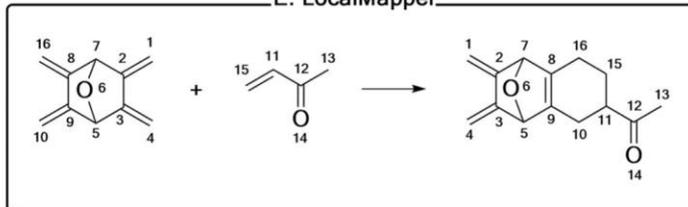
C. GraphormerMapper



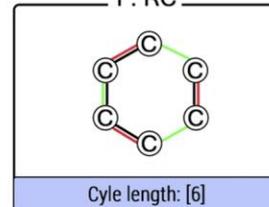
D. RC



E. LocalMapper

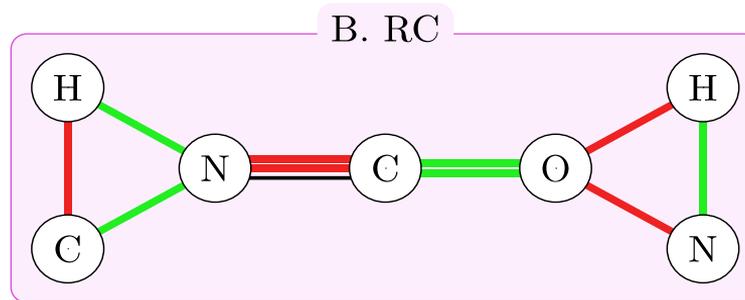
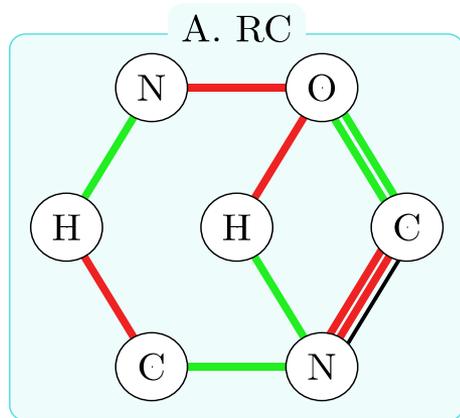
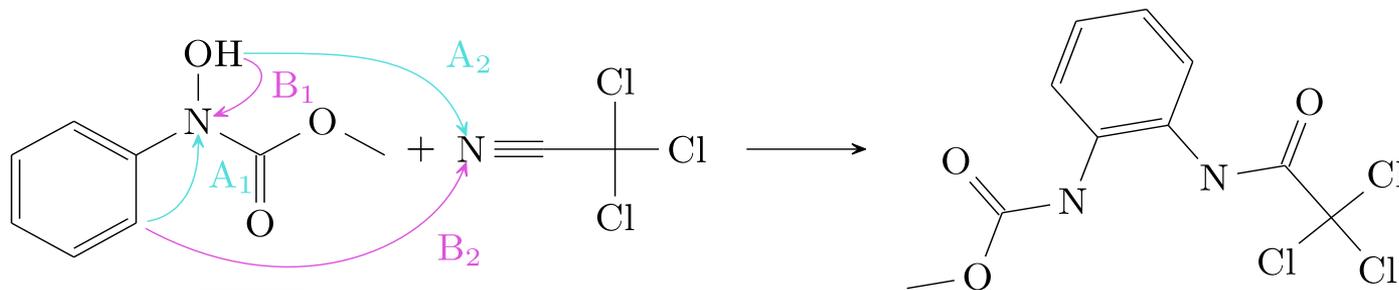


F. RC



DISCUSSION

Ambiguous Hydrogen



Our team



Phuoc-Chung Van Nguyen
UMP



Ngoc-Vi Nguyen Tran
Uppsala University

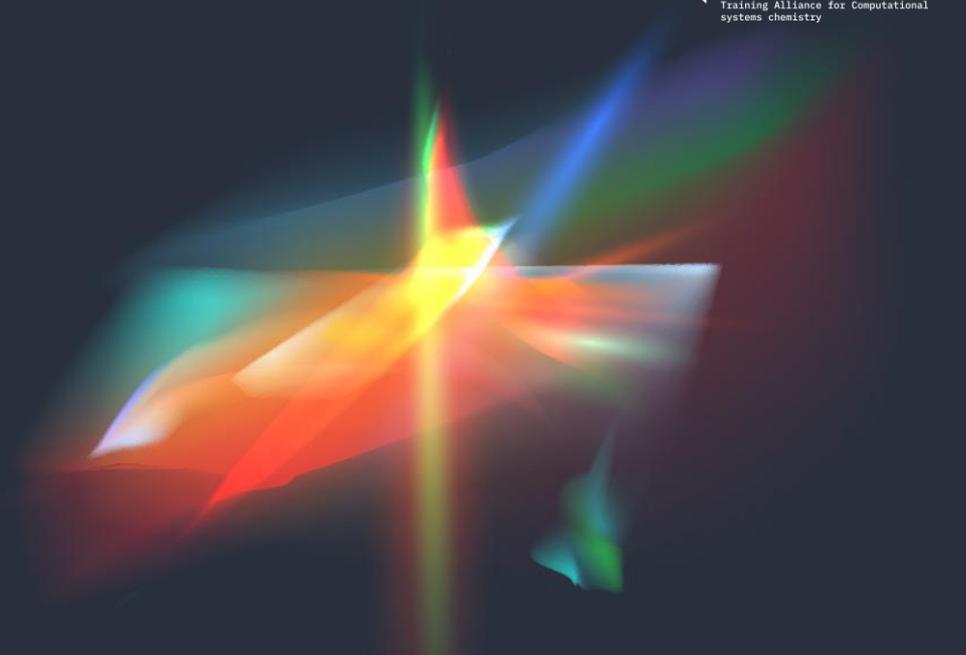


Tieu-Long Phan
Leipzig University



Peter F. Stadler
Leipzig University

Thank you for your attendance



Founded by the
European Union

This project has received funding from the European Union's Horizon 2021 research and innovation programme under the Marie-Sklodowska-Curie grant agreement No 101072930

