

Molecular Transformer-aided Biocatalysed Synthesis Planning



Alain Vaucher
15 August 2022

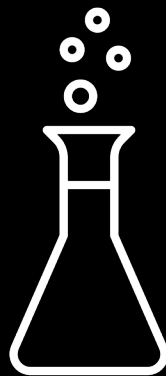
IBM Research

Introduction

Data

Chemistry

Artificial
intelligence



(Bio)catalysis

Automation

OUTLINE

1. AI models for chemistry
2. Synthesis automation – RoboRXN
3. AI for biocatalysis

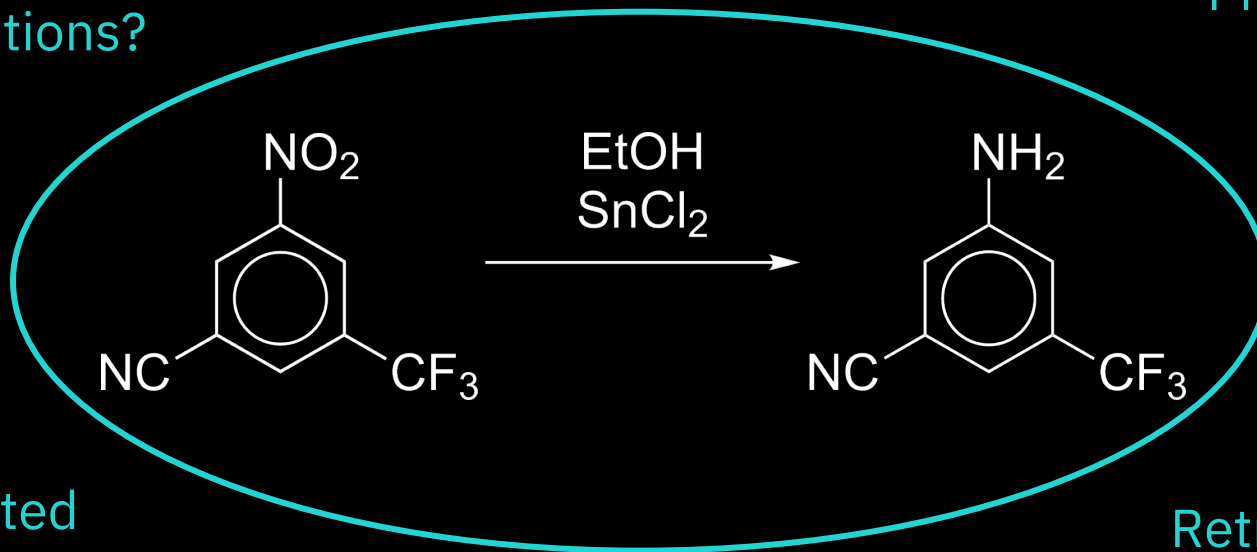
OUTLINE

1. AI models for chemistry
2. Synthesis automation – RoboRXN
3. AI for biocatalysis

AI and chemical reactivity

Experimental conditions?

Product?



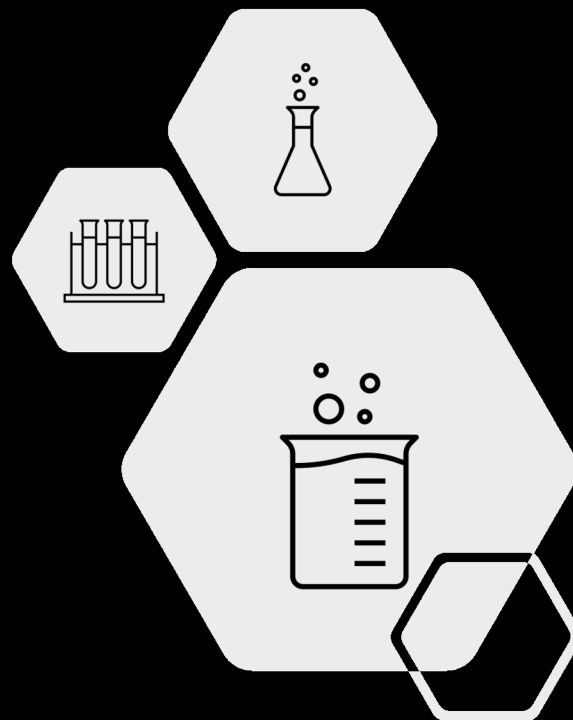
Related reactions?

Retrosynthesis?

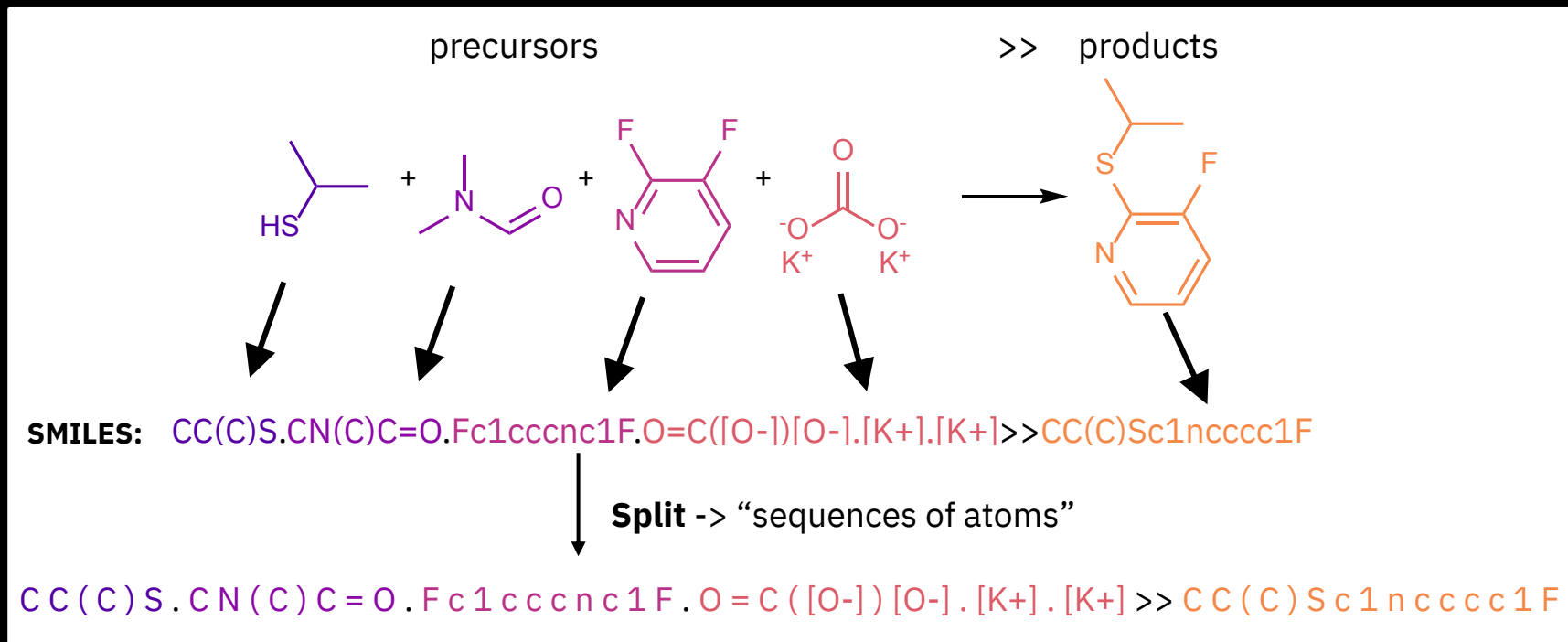
Yield?

Data sources

- Millions of reactions have been reported
- Sources:
 - Publicly available data: patents (USPTO, NextMove's Pistachio, etc.)
 - Scientific publications
 - Proprietary reactions (industry)
 - Publishers
 - Etc.

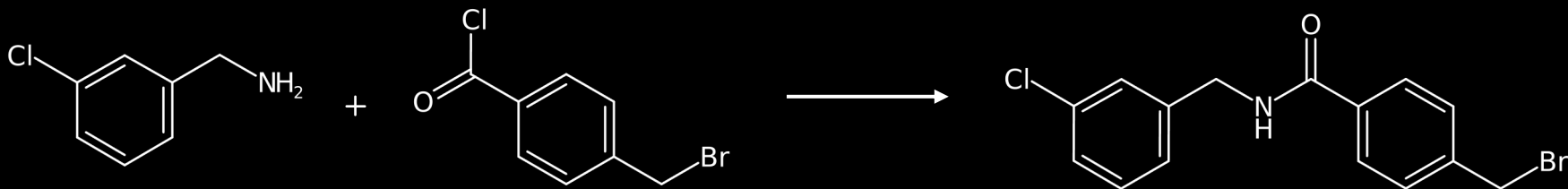


Atoms as letters, molecules as words



→ Borrow methods developed for human languages

Forward reaction prediction



Textual representation (SMILES)

NCc1ccc(Cl)c1

O=C(Cl)c1ccc(CBr)cc1

O=C(Cl)c1ccc(CBr)cc1

“Sentence of atoms”

N C c 1 c c c c (C l) c 1 . O = C (C l) c 1 c c c (C B r) c c 1

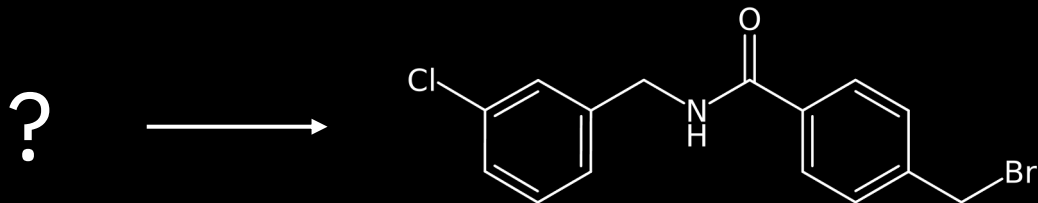
“Translation”

Transformer

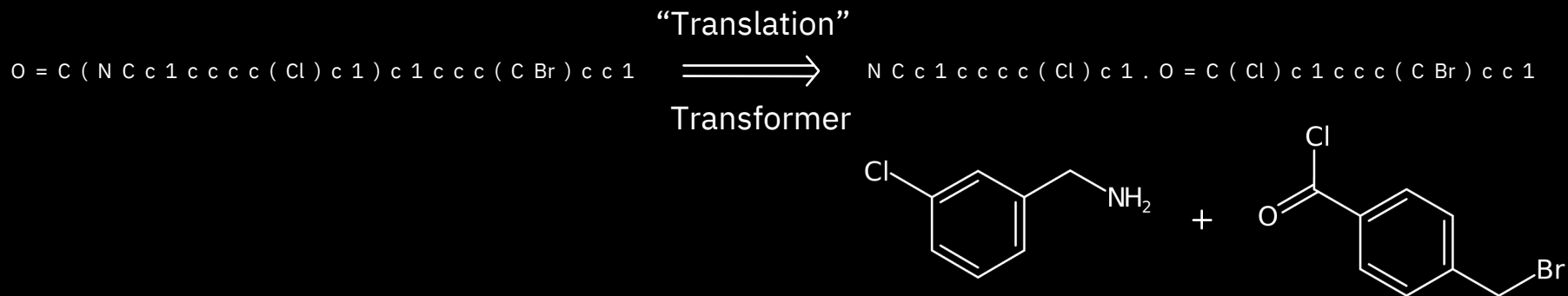
O = C (N C c 1 c c c c (C l) c 1) c 1 c c c (C B r) c c 1

Schwaller, P.; Laino, T.; Gaudin, T.; Bolgar, P.; Hunter, C. A.; Bekas, C. & Lee, A. A., *ACS Cent. Sci.*, **2019**, 5, 1572-1583.

Retrosynthesis



Similar approach, both sides switched

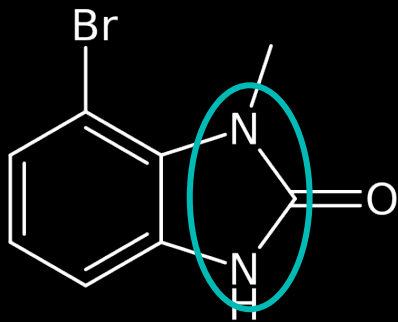


Schwaller, P.; Petraglia, R.; Zullo, V.; Nair, V. H.; Haeuselmann, R. A.; Pisoni, R.; Bekas, C.; Iuliano, A. & Laino, T., *Chem. Sci.*, **2020**, *11*, 3316-3325.

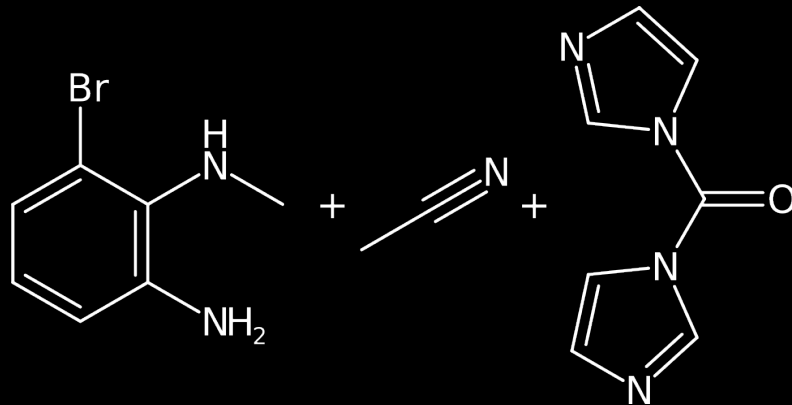
Disconnection-aware retrosynthesis

- Let the chemists decide where to break the compound?

Input (target compound)

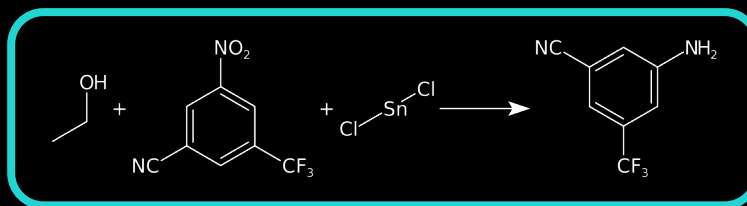
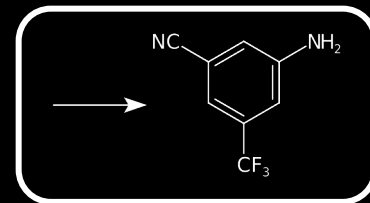
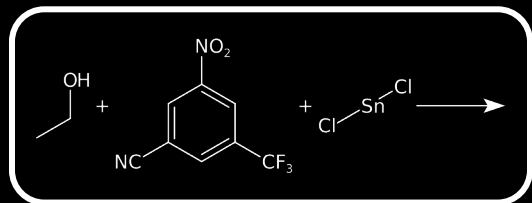
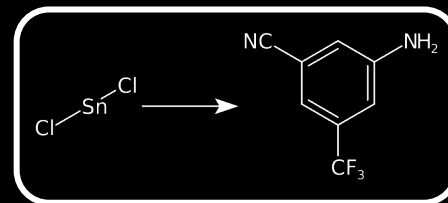
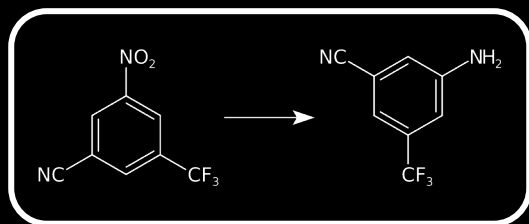
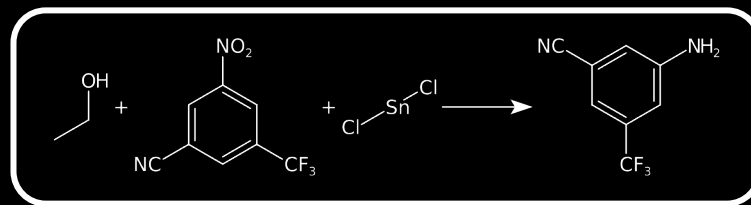
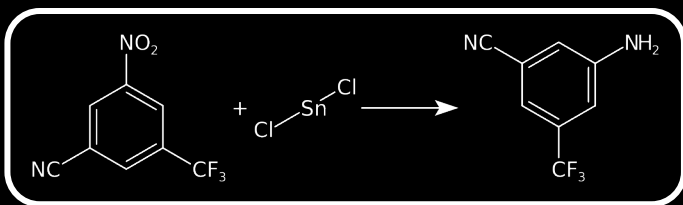


Output (precursors)

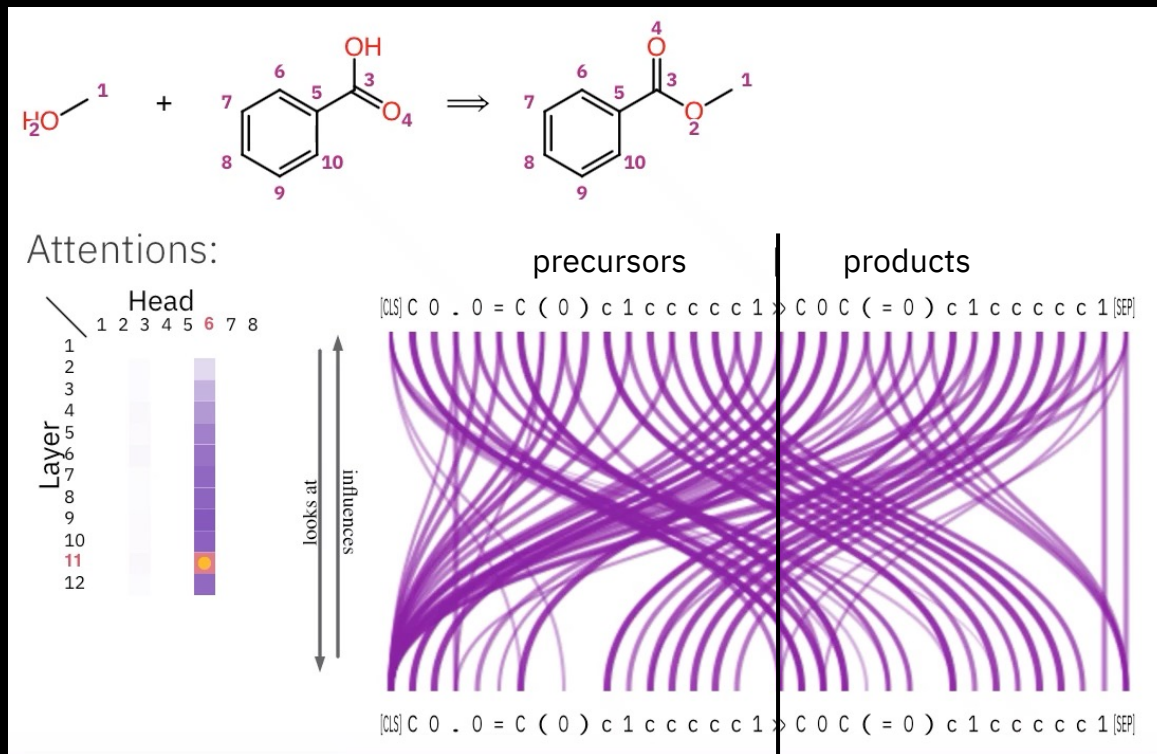
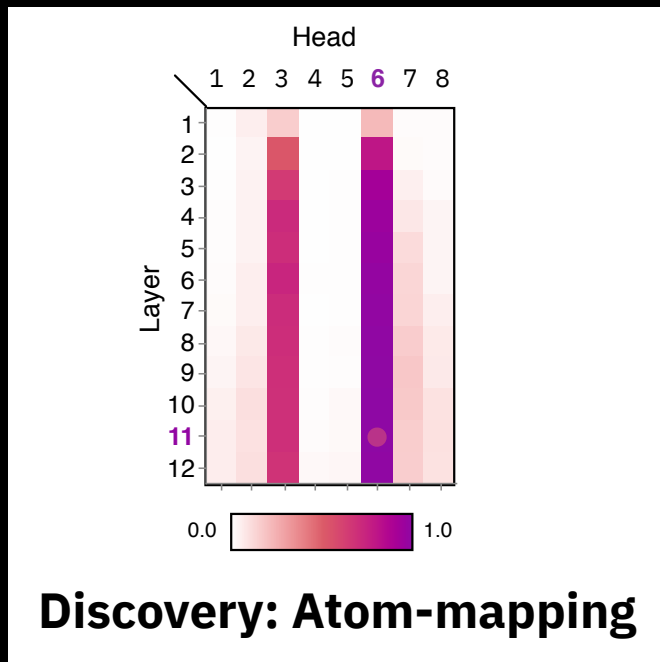


Byekwaso, A.; Schwaller, P.; Vaucher, A. C.; Toniato, A.; Laino, T.; "AI for Science" workshop @ NeurIPS 2021.

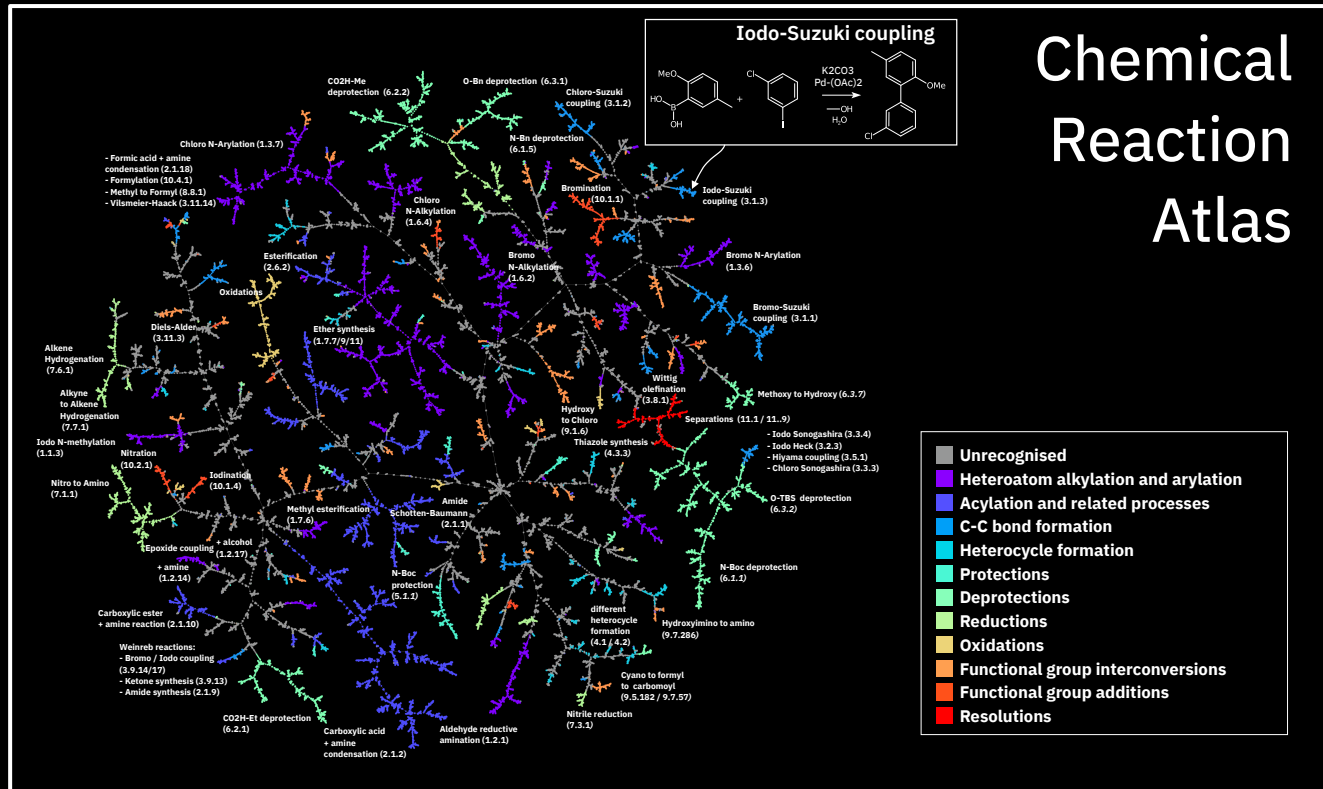
Completing partial chemical equations



Atom-Mapping and the learning of chemical reaction grammar

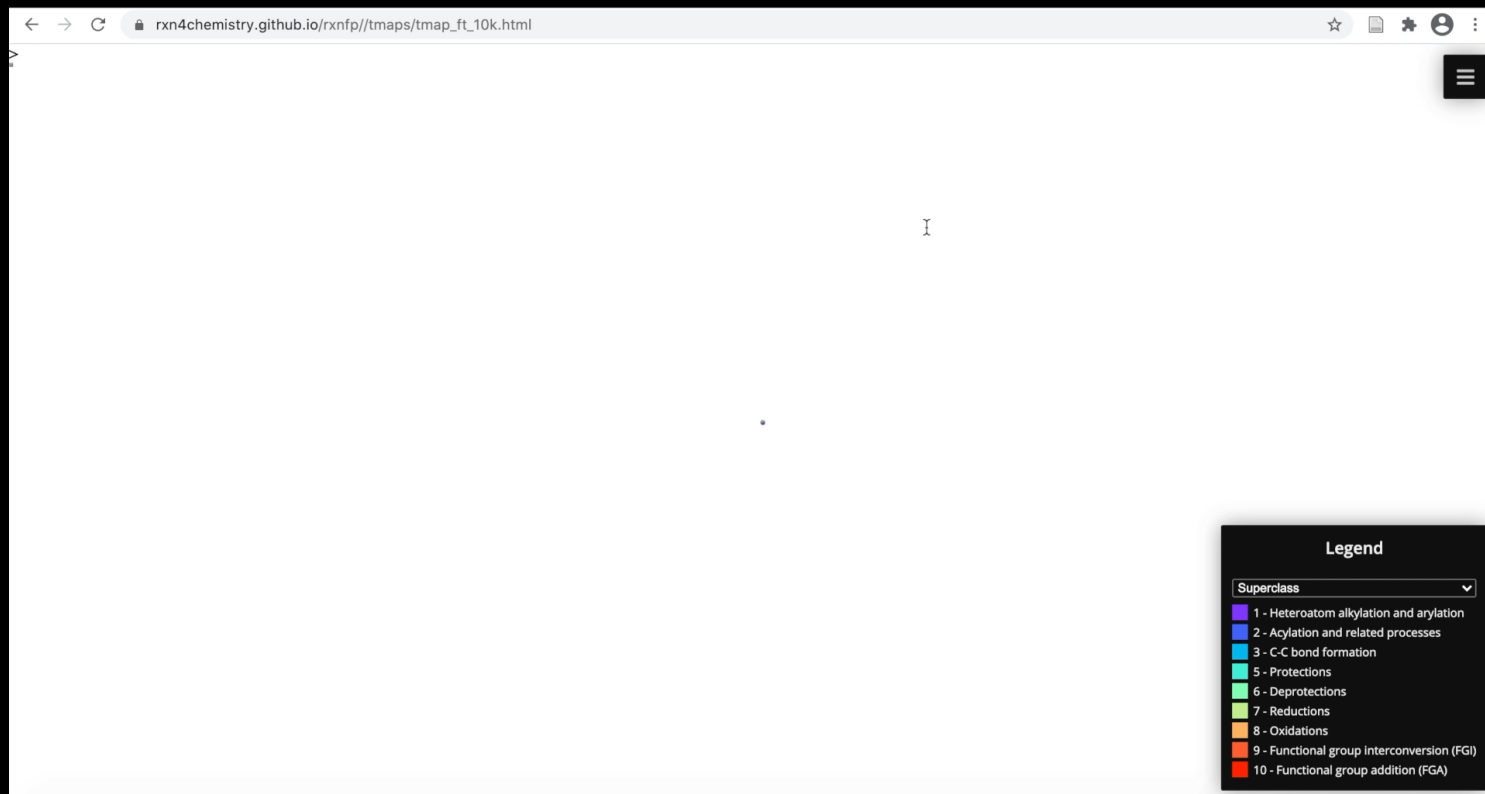


Mapping and classifying chemical reactions



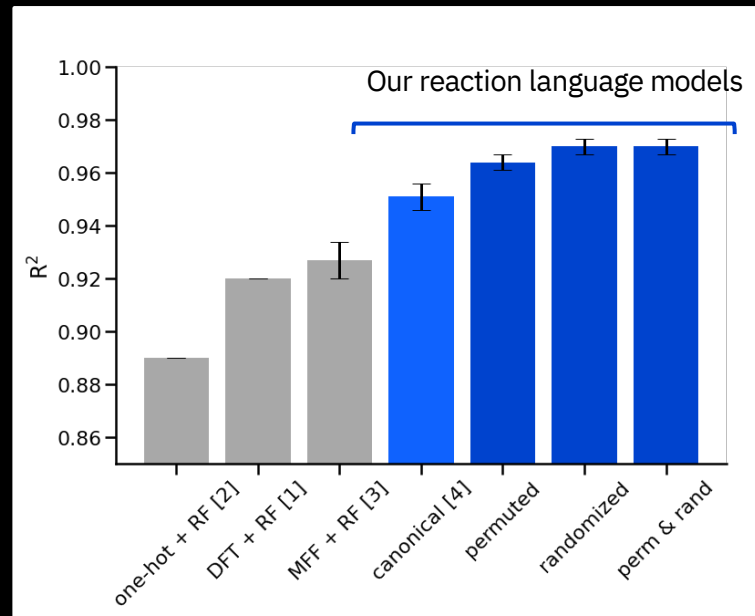
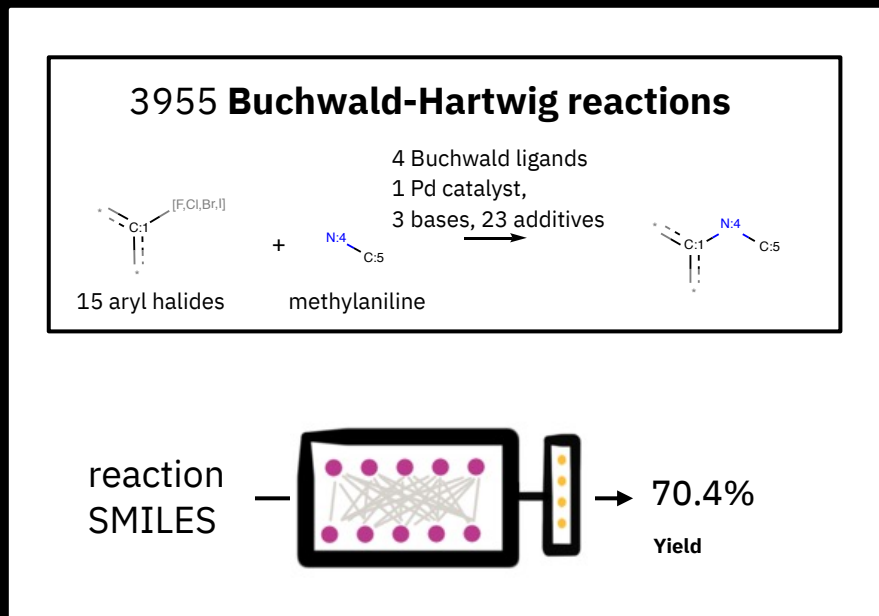
Schwaller, P.; Probst, D.; Vaucher, A. C.; Nair, V. H.; Kreutter, D.; Laino, T. & Reymond, J.-L., *Nat. Mach. Intell.*, **2021**, 3, 144-152.

Mapping and classifying chemical reactions



Schwaller, P.; Probst, D.; Vaucher, A. C.; Nair, V. H.; Kreutter, D.; Laino, T. & Reymond, J.-L., *Nat. Mach. Intell.*, **2021**, 3, 144-152.

Prediction of chemical reaction yields



- [1] Ahneman, D. T., Estrada, J. G., Lin, S., Dreher, S. D. & Doyle, A. G. Predicting reaction performance in C–N cross-coupling using machine learning. *Science* **360**, 186–190 (2018).
- [2] Chuang, K. V. & Keiser, M. J. Comment on “Predicting reaction performance in C–N cross-coupling using machine learning”. *Science* **362** (2018).

- [3] Sandfort, F., Strieth-Kalthoff, F., Kühnemund, M., Beecks, C. & Glorius, F. A structure-based platform for predicting chemical reactivity. *Chem* (2020).
- [4] Schwaller, P., Vaucher, A. C., Laino, T. & Reymond, J.-L. Prediction of chemical reaction yields using deep learning. *ChemRxiv preprint* doi:10.26434/chemrxiv.12758474 (2020).

Schwaller, P.; Vaucher, A. C.; Laino, T. & Reymond, J.-L., *Mach. Learn.: Sci. Technol.*, **2021**, *2*, 015016.
Schwaller, P.; Vaucher, A. C.; Laino, T. & Reymond, J.-L., *Chemrxiv.13286741*, **2020**.

OUTLINE

1. AI models for chemistry
2. Synthesis automation – RoboRXN
3. AI for biocatalysis

Data and chemical reactions

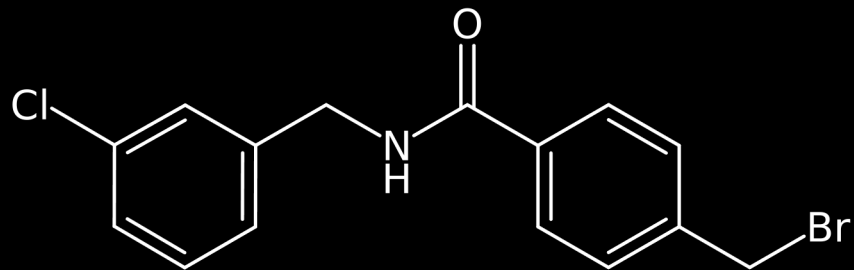
- Chemists have been doing reactions in roughly the same way for **decades**
- Set of **standard lab operations**
- **Millions** of reactions reported in the literature



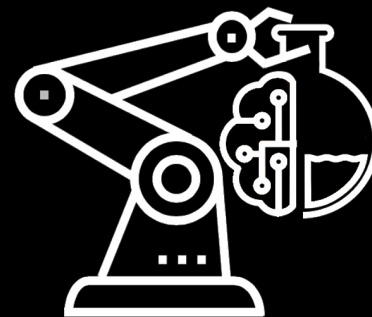
How can exploit this **data** to **accelerate discovery**?

- Assist chemists in synthesis planning
- ... and run the syntheses for them!

Automated execution on robot

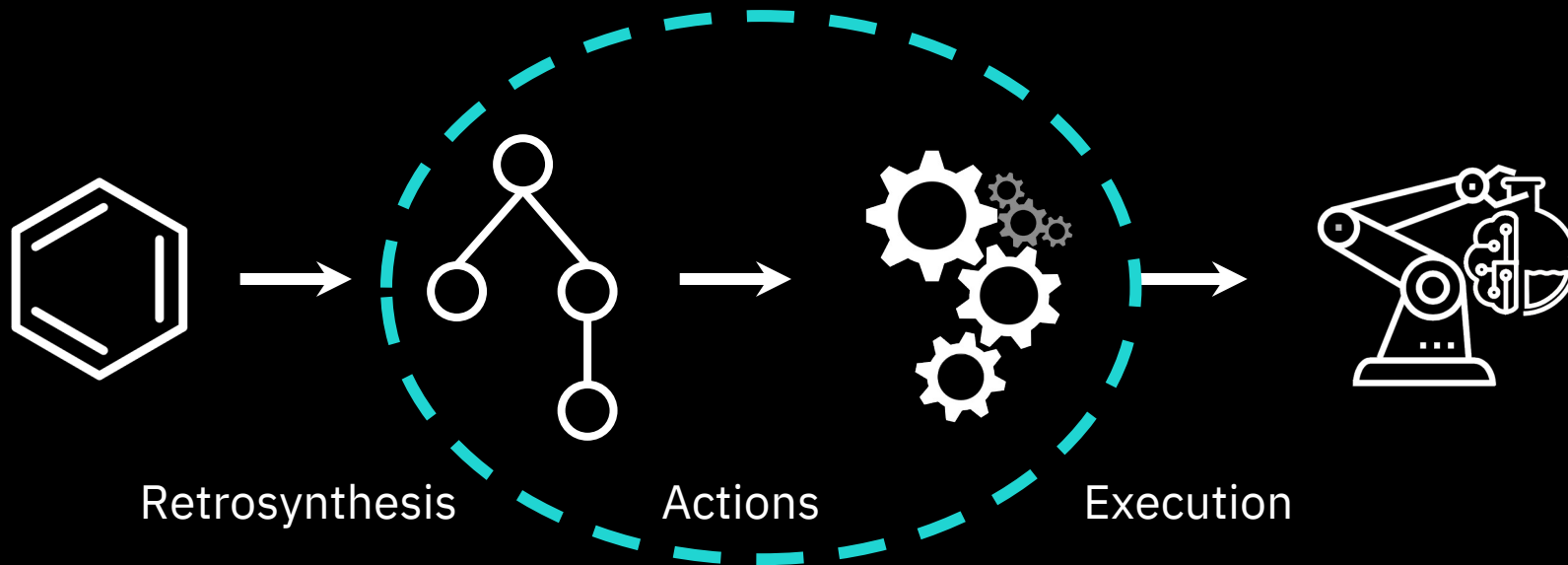


Target molecule



Synthesis execution

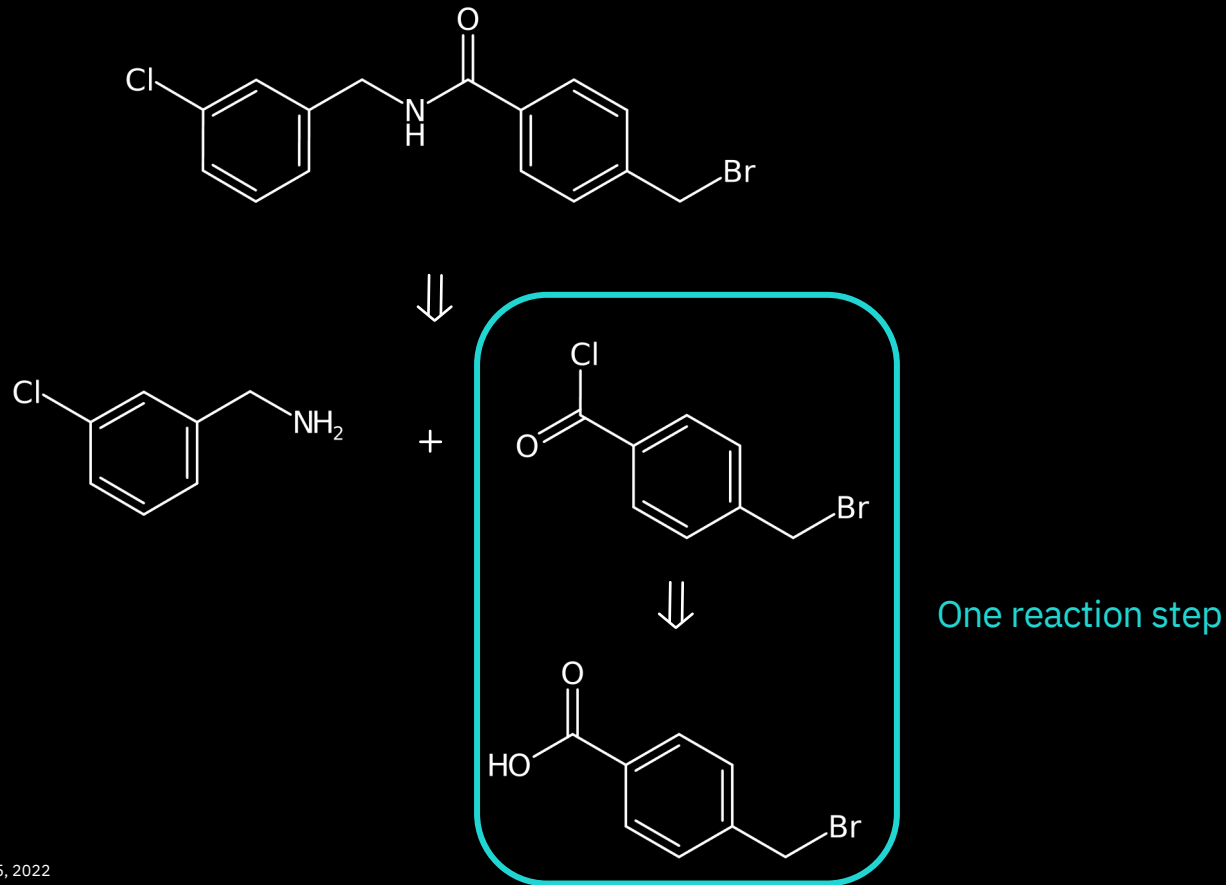
Automated execution on robot



OUTLINE

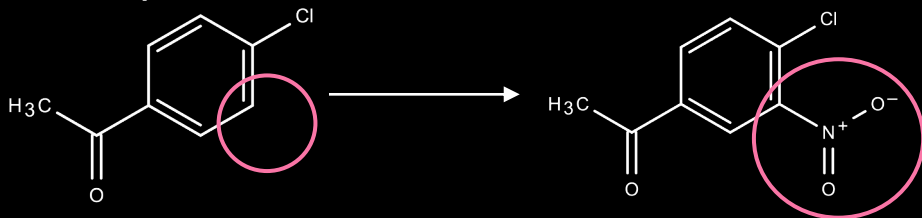
1. AI models for chemistry
2. Synthesis automation – RoboRXN
 - A. Prediction of synthesis actions
 - B. Execution on robot
3. AI for biocatalysis

Synthesis actions



Synthesis actions

Example:

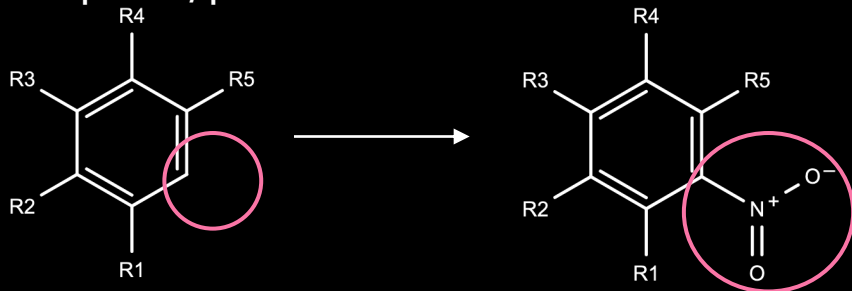


– Same template but different synthesis actions!

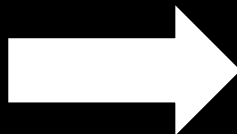
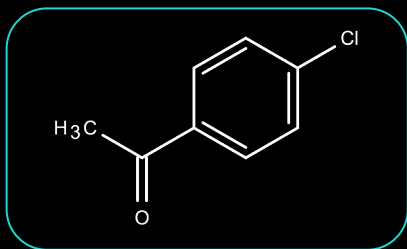
– Hard to predict

– Ideally: **ML model!**

Template/pattern:



Synthesis actions



Operation 1

Operation 2

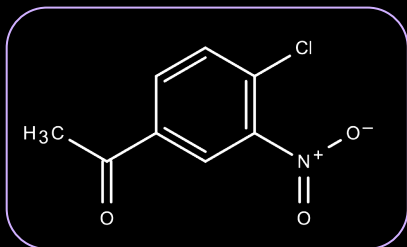
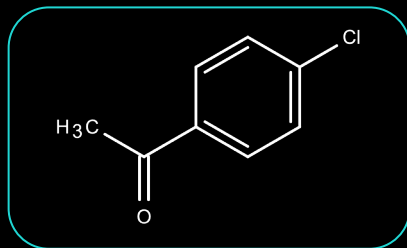
Operation 3

Operation 4

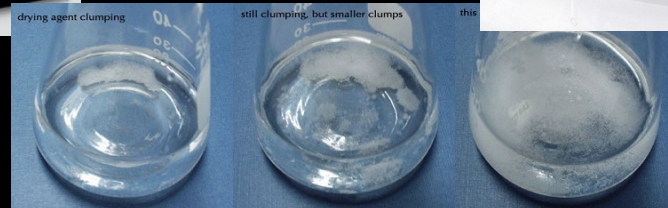
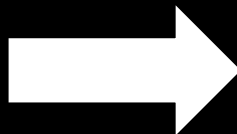
...

```
C1=CC(C(=O)C)=CC=C1Cl>>C1=CC(C(=O)C)=CC([N+]([O-])=O)=C1Cl
```

Synthesis actions



C1=CC(C(=O)C)=CC=C1Cl>>C1=CC(C(=O)C)=CC([N+]([O-])=O)=C1Cl



Synthesis actions

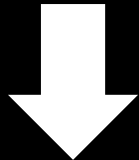
- No dataset!
- Information is available **indirectly**
- First: extract actions from text

Example procedure from a patent

A mixture of 1-(4-isopropyl-phenyl)-5-oxo-pyrrolidine-3-carboxylic acid ethyl ester obtained in step 2 (0.7 g, 2.65 mmol) and ethanol were cooled to 10-15° C. Sodium borohydride (0.25 g, 6.6 mmol) was added portion wise over a period of 20 min and the reaction mixture was stirred for 3.5 hrs at 20-25° C. The organic volatiles were evaporated and the residue was taken into brine solution (15 ml). The aqueous layer was extracted with ethyl acetate, dried over Na₂SO₄ and evaporated to obtain 4-hydroxymethyl-1-(4-isopropyl-phenyl)-pyrrolidin-2-one as an off white solid (0.5 g, 81%).

Models for Paragraph-to-actions

... Sodium borohydride (0.25 g, 6.6 mmol) was added portion wise over a period of 20 min and the reaction mixture was stirred for 3.5 hrs at 20-25° C ...



```
Add(name='Sodium borohydride',  
      quantity=['0.25 g', '6.6 mmol'],  
      duration='20 min')
```

```
Stir(temperature='20-25°C',  
      duration='3.5 hrs')
```

What kind of model?

–Rule-based model?

–Fully data-driven model?

Both!

Hand-annotated data

Action ID	Type and properties	Edit properties	Delete action
35210	PHASESEPARATION		
35211	COLLECTLAYER organic		
35212	WASH with water (500 ml)		

Initial actions from rule-based model

Sentence to annotate

>1700 annotated sentences

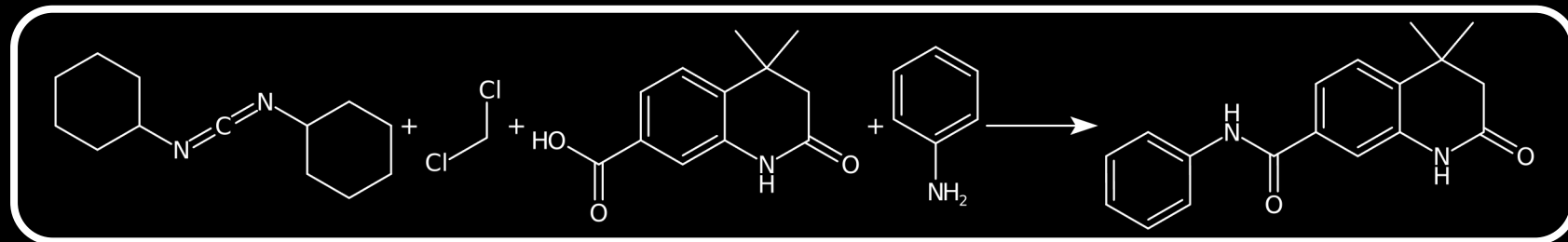
Delete property	Type	Text
	material	brine

Results

Model	100% accuracy
Combined rule-based model	21.9
Pretrained translation model	24.7
Model without pretraining	37.8
Refined translation model	60.8

Vaucher, A. C.; Zipoli, F.; Geluykens, J.; Nair, V. H.; Schwaller, P.; Laino, T., *Nat. Commun.* **2020**, *11*, 3601.

SMILES-to-actions



```
C(=NC1CCCCC1)=NC1CCCCC1 . ClCCl . CC1(C)CC(=O)Nc2cc(C(=O)O)ccc21 . Nc1ccccc1 >> CC1(C)CC(=O)Nc2cc(C(=O)Nc3ccccc3)ccc21
```

2.7 g (12.3 mmol) 4,4-Dimethyl-1,2,3,4-tetrahydro-2-oxo-7-quinolinecarboxylic acid were added to a solution of 3.8 g (18.5 mmol) N,N'-dicyclohexylcarbodiimide and 1.1 ml (12.3 mmol) aniline in 80 ml dichloromethane. The reaction mixture was stirred for 4 hours at ambient temperature and the precipitate was filtered off with suction and recrystallised from ethanol. There was obtained 1.2 g of the title compound; m.p. 249-251° C.

ML model

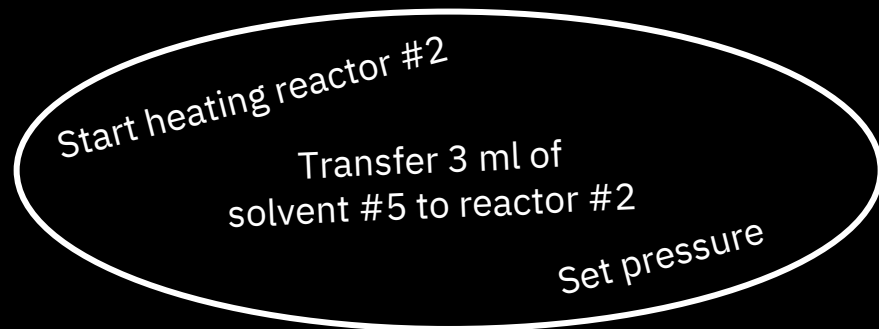
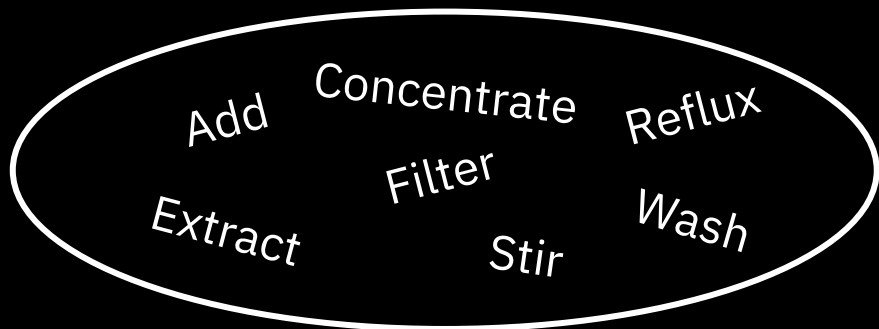
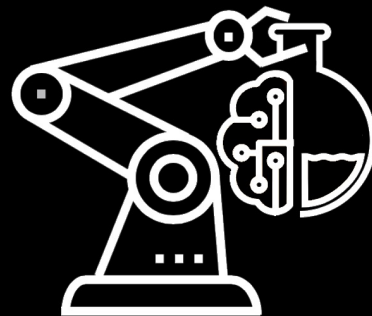
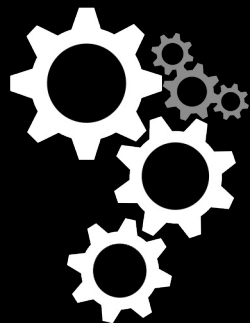
1. MAKESOLUTION with N,N'-dicyclohexylcarbodiimide (3.8 g, 18.5 mmol) and aniline (1.1 ml, 12.3 mmol) and dichloromethane (80 ml)
2. STIR for 4 hours at ambient temperature
3. ADD 4,4-Dimethyl-1,2,3,4-tetrahydro-2-oxo-7-quinolinecarboxylic acid (2.7 g, 12.3 mmol)
4. STIR for 4 hours at ambient temperature
5. FILTER keep precipitate
6. RECRYSTALLIZE from ethanol
7. YIELD title compound (1.2 g)

1. ADD \$1\$
2. ADD \$4\$
3. ADD \$2\$
4. ADD \$3\$
5. STIR for @3@ at #4#
6. FILTER keep precipitate
7. RECRYSTALLIZE from ethanol
8. YIELD \$-1\$

OUTLINE

1. AI models for chemistry
2. Synthesis automation – RoboRXN
 - A. Prediction of synthesis actions
 - B. Execution on robot
3. AI for biocatalysis

Execution on chemical robot

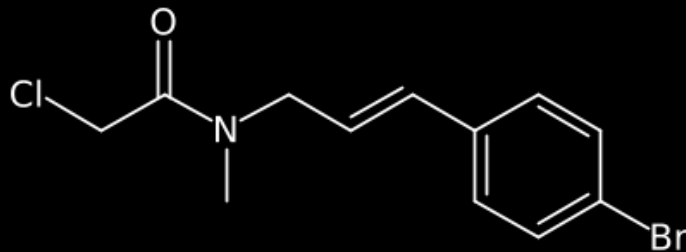


Execution on chemical robot

Cloud-based setup for
autonomous synthesis



DEMO



CN(C/C=C/c1ccc(Br)cc1)C(=O)CCl

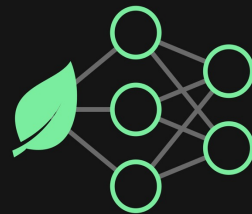
OUTLINE

1. AI models for chemistry
2. Synthesis automation – RoboRXN
3. AI for biocatalysis

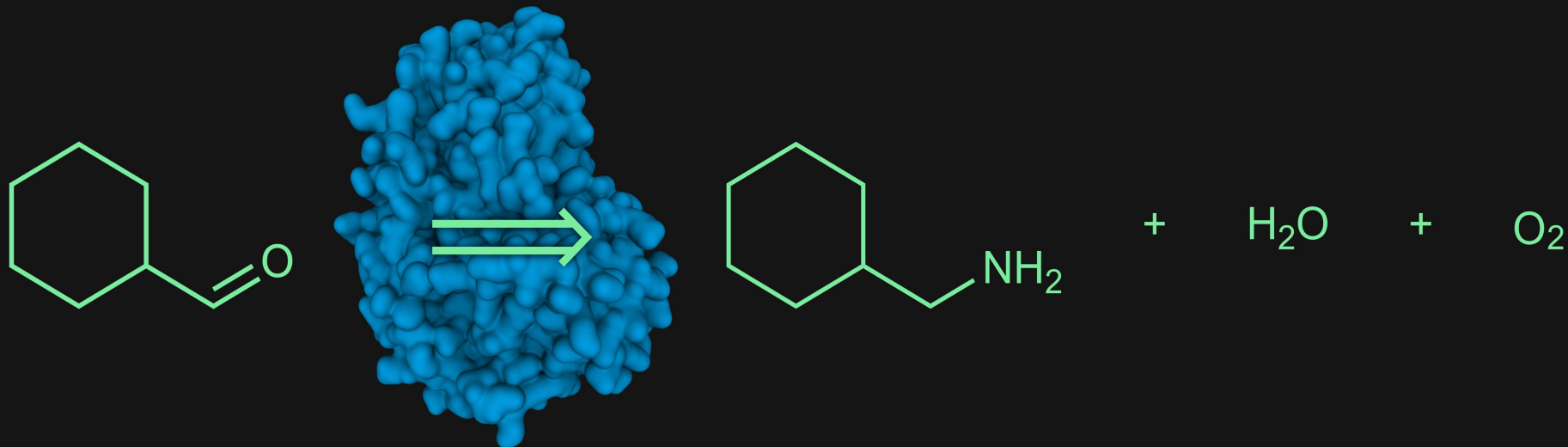
OUTLINE

1. AI models for chemistry
2. Synthesis automation – RoboRXN
3. AI for biocatalysis
 - A. Biocatalyzed synthesis planning
 - B. Identification of active sites

Enzymatic catalyst



GreenCatRXN

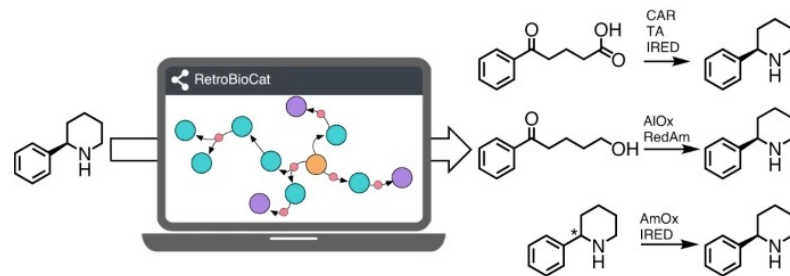


EC 1.4.3.-

~~TEMPO~~ ~~NaClO~~ ~~DEM~~

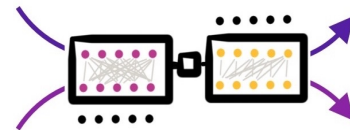
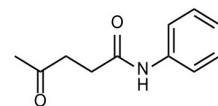
Enzymes in synthesis planning

- CASP (computer-assisted synthetic planning) using enzymes currently relying on rule-based approaches (Finnigan et al., Nat. Catal., 2021)
- Kreutter, Schwaller, and Reymond have shown that including enzyme name information yields good results



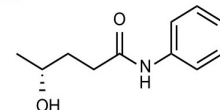
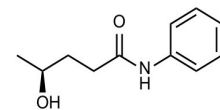
Finnigan et al, Nat. Catal., 2021

ketoreductase p3-b03



ketoreductase p1-a04

Enzymatic Transformer

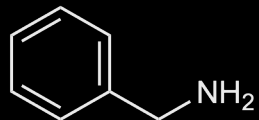


Kreutter et al., Chem. Sci., 2021

Our intuition

cyclohexanemethylamine + water + oxygen

= ammonia + cyclohexanecarbaldehyde + hydrogen peroxide



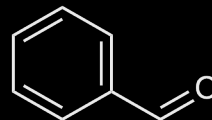
+



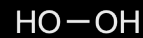
+



+

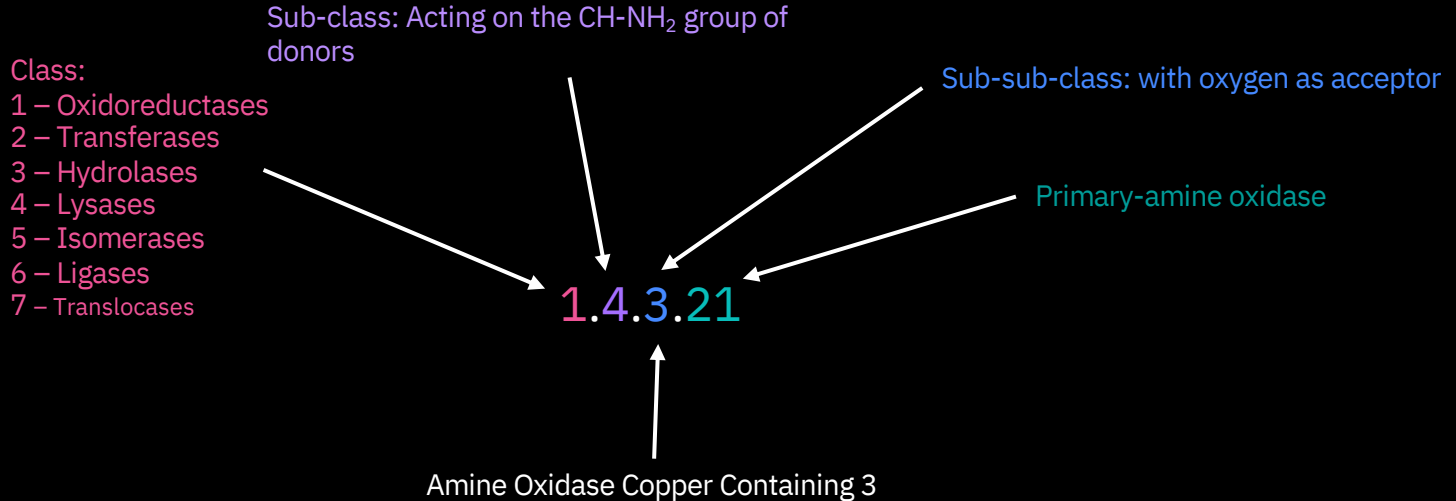


+



Our intuition

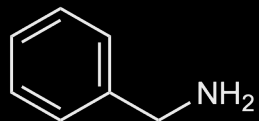
cyclohexanemethylamine + water + oxygen = ammonia + cyclohexanecarbaldehyde + hydrogen peroxide



Our intuition

cyclohexanemethylamine + water + oxygen

= ammonia + cyclohexanecarbaldehyde + hydrogen peroxide



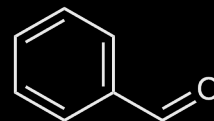
+



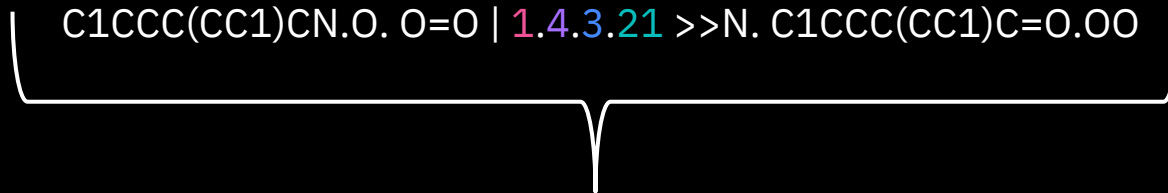
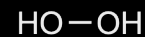
+



+



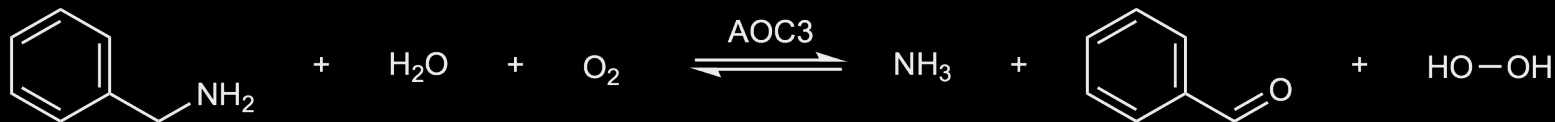
+



Enzymatic reaction SMILES

Our intuition

cyclohexanemethylamine + water + oxygen = ammonia + cyclohexanecarbaldehyde + hydrogen peroxide



Enzymatic reaction SMILES:

```
C1CCC(CC1)CN.O.O=O | 1.4.3.21 >>N.C1CCC(CC1)C=O.OO
```

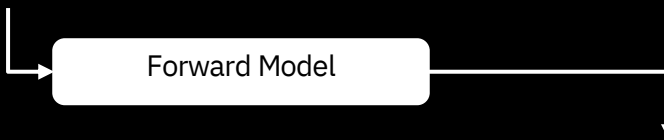
Tokenized reaction

```
C 1 C C C ( C C 1 ) C N . O . O = O | [v1] [u4] [t3] [q21] >>N.C 1 C C C ( C C 1 ) C = O . O O
```

Our approach for biocatalyzed synthesis planning

Forward prediction

C1CCC(CC1)CN.O.O=O | [v1] [u4] [t3] [q21] >> N.C1CCC(CC1)C=O.OO

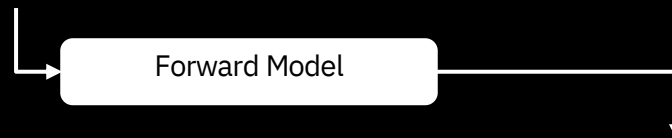


C1CCC(CC1)CN.O.O=O | [v1] [u4] [t3] [q21] >> N.C1CCC(CC1)C=O.OO

Our approach for biocatalyzed synthesis planning

Forward prediction

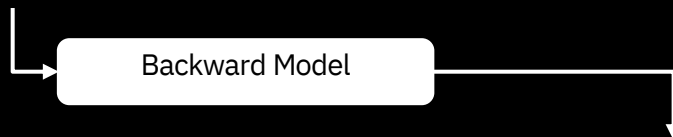
C1CCC(CC1)CN.O.O=O | [v1] [u4] [t3] [q21] >> N.C1CCC(CC1)C=O.OO



C1CCC(CC1)CN.O.O=O | [v1] [u4] [t3] [q21] >> N.C1CCC(CC1)C=O.OO

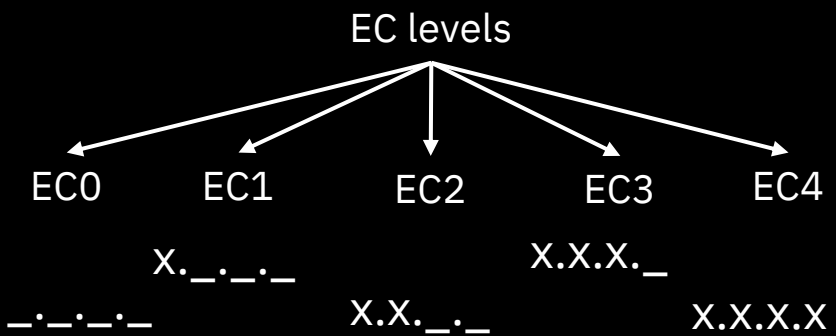
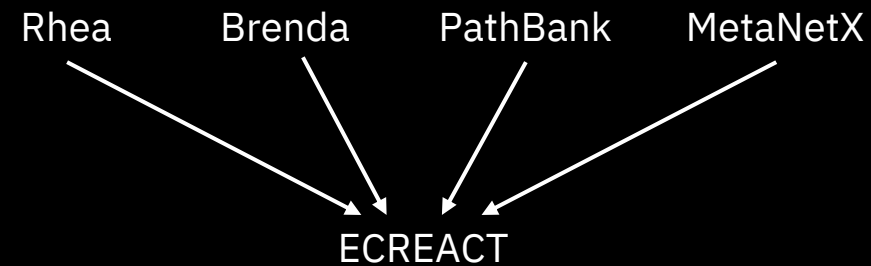
Backward prediction

C1CCC(CC1)CN.O.O=O | [v1] [u4] [t3] [q21] >> N.C1CCC(CC1)C=O.OO

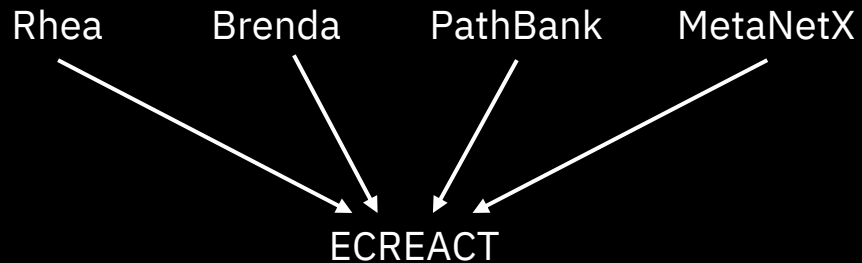


C1CCC(CC1)CN.O.O=O | [v1] [u4] [t3] [q21] >> N.C1CCC(CC1)C=O.OO

Dataset



Dataset



Total number of reactions: 63.403

Most represented reaction type: Transferases (EC 2.x.x.x)

Least represented reaction type: Translocases (EC 7.x.x.x)

Probst D. et al., ChemRxiv, 2021

Rafael A. et al., Nucleic Acids Res., 2012;

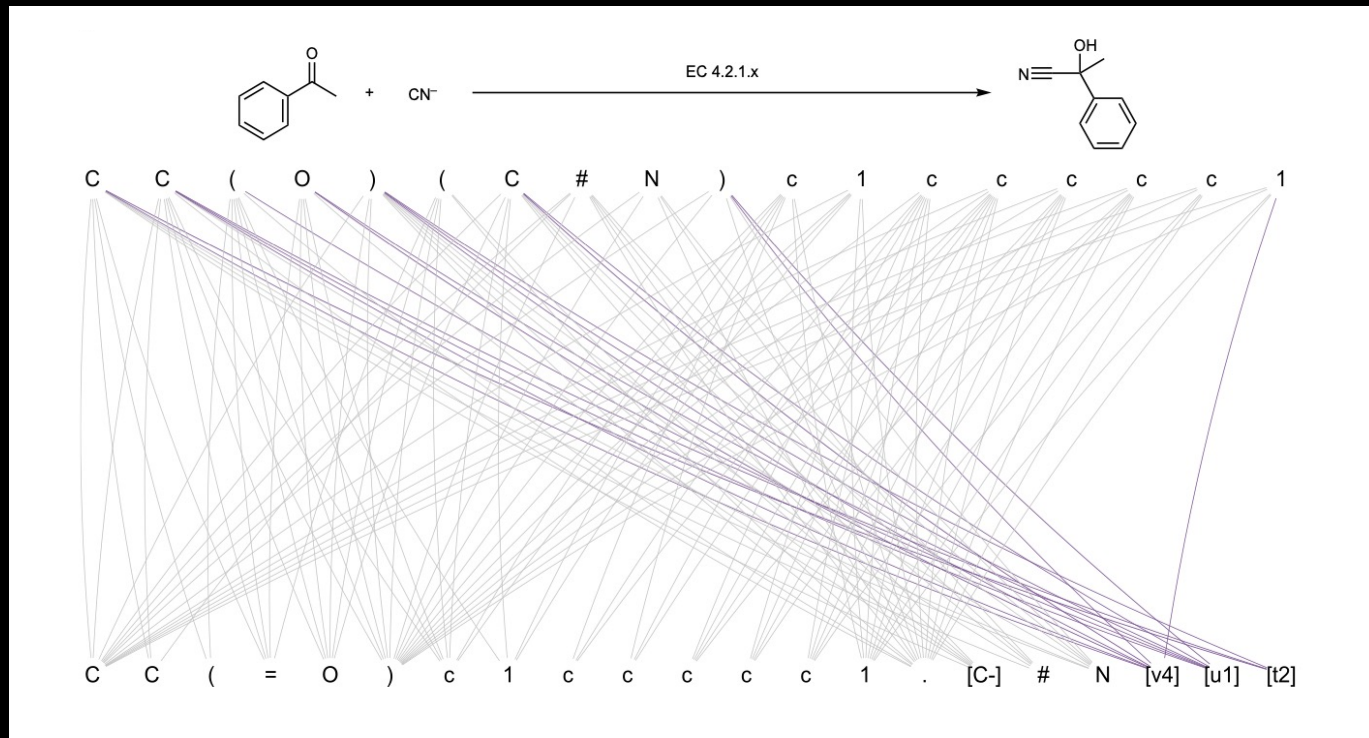
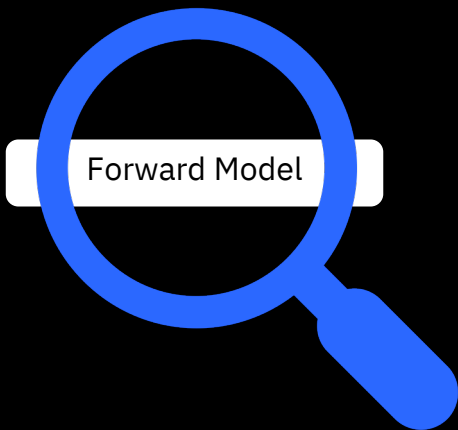
Chang A., et al. Nucleic Acids Res., 2021;

David S., Nucleic Acids Res., 2020;

Thomas N, ACS, 2019



Attention Analysis



Try GreenCatRXN on RXN for Chemistry: <https://rxn.res.ibm.com>

Have you ever thought about comparing
Traditional Organic Synthesis
VS
Enzymatic Synthesis ?



GitHub: <https://github.com/rxn4chemistry/rxn4chemistry>

OUTLINE

1. AI models for chemistry
2. Synthesis automation – RoboRXN
3. AI for biocatalysis
 - A. Biocatalyzed synthesis planning
 - B. Identification of active sites

Active site in proteins

Proteins' activity is directly related to the structure of the active site

Biological function annotation of proteins usually rely on 3D structural models ([Yousaf et al., 2021, Kozlovskii and Popov, 2020 and 2021])

Some methods identify active sites via sequence similarity (Pfam [Mistry et al., 2020] and PSI-BLAST [Altschul et al., 1997])

PROTEIN STRUCTURE

Scaffold to support and position active site

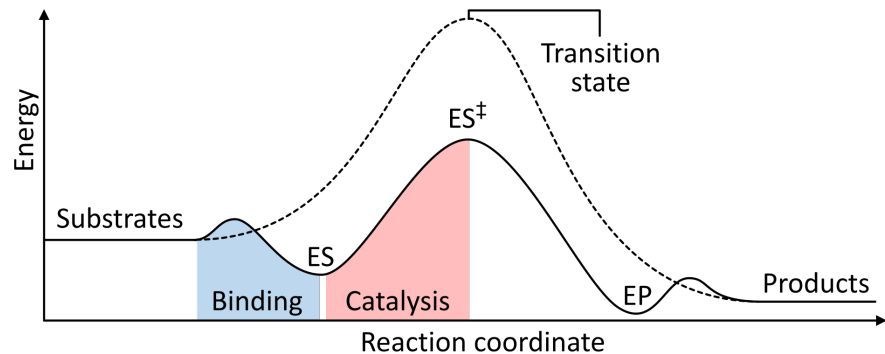
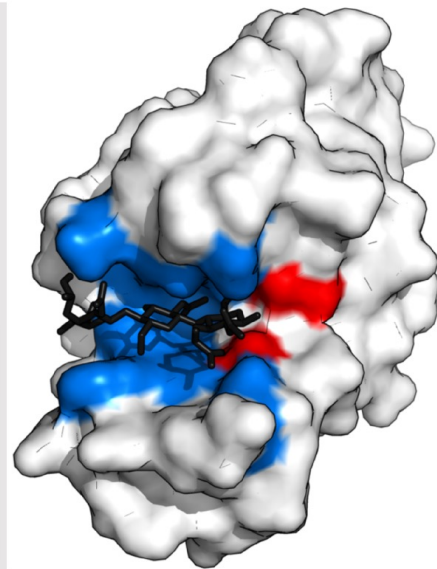
ACTIVE SITE

BINDING SITES

Bind and orient substrate(s)

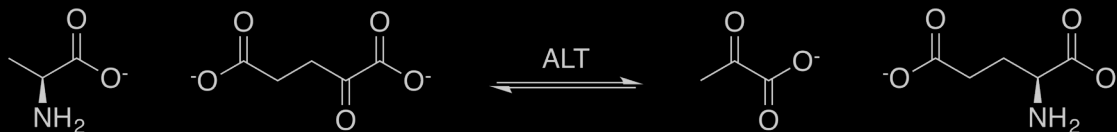
CATALYTIC SITE

Reduce chemical activation energy



Biocatalysis as a language

L-alanine + 2-oxoglutarate = pyruvate + L-glutamate



Reaction SMILES

C[C@H](N)C([O-])=O . [O-]C(=O)CCC(=O)C([O-])=O >> CC(=O)C([O-])=O . N[C@@H](CCC([O-])=O)C([O-])=O

Introducing AA sequences

```
10 20 30 40 50
MASSTGDRSQ AVRHGLRAKV LTLDGMNPRV RRVEYAVRGP IVQRALELEQ
60 70 80 90 100
ELRQGVKKPF TEVIRANIGD AQAMGQRPIT FLRQVLALCV NPDLLSSPNF
110 120 130 140 150
PDDAKKRAER ILQACGGHSL GAYSVSSGIQ LIREDVARYI ERRDGGIPAD
160 170 180 190 200
PNNVFLSTGA SDAIVTVLKL LVAGEGHTRT GVLIPQPYP LYSATLAELG
210 220 230 240 250
AVQVDYYLDE ERAWALDVAE LHRALGQARD HCRPRALCVI NPGNPTGQVQ
260 270 280 290 300
TRECIEAVIR FAFEERLFL L ADEVYQDNVY AAGSQFHSFK KVLMEMGPPY
310 320 330 340 350
AGQQELASFH STSKGYMGEC GFRGGYVEV NMDAAVQQQM LKLM SVRLCP
360 370 380 390 400
PVPGQALLDL VVSPAPTDP SFAQFQAEKQ AVLAELA AKA KLTEQVFNEA
410 420 430 440 450
PGISCNPVQG AMYSFPRVQL PPRaveraQE LGLAPDMFFC LRLLEETGIC
460 470 480 490
VWPGSGFGQR EGTYHFRMTI LPPLEKLRL L LEKLSRFHAK FTLEYS
```



Reaction SMILES

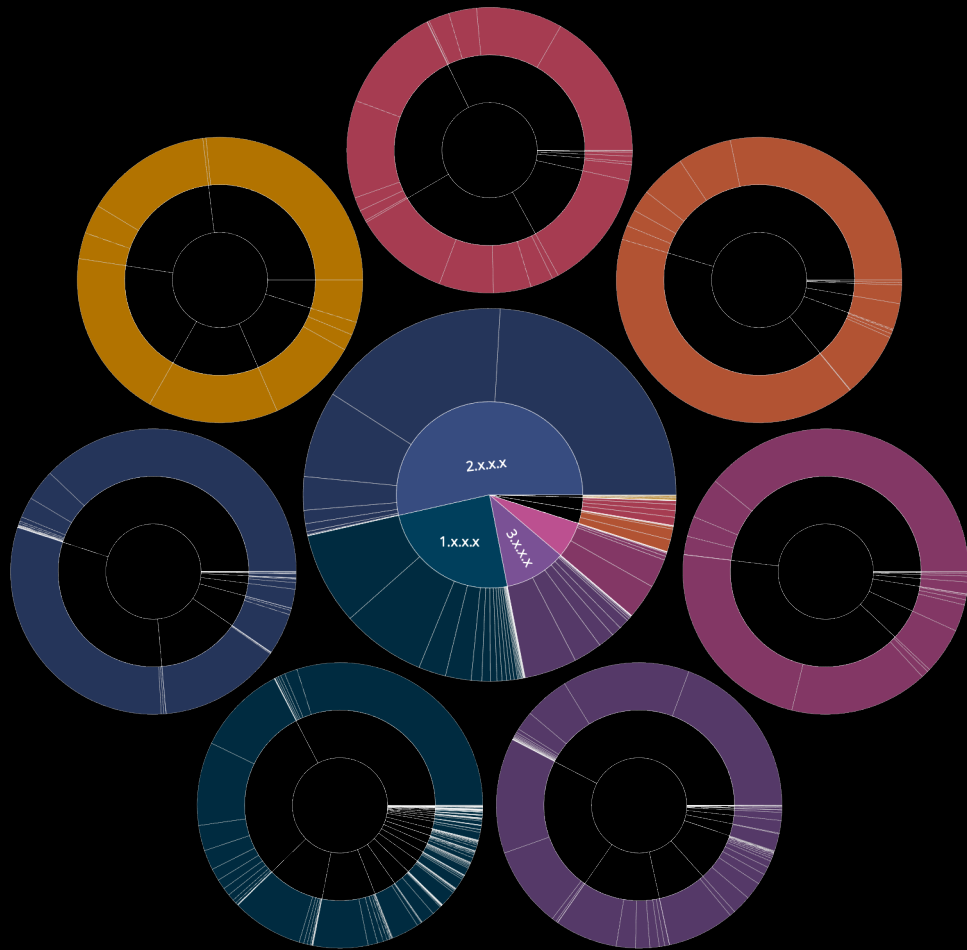
C[C@H](N)C([O-])=O . [O-]C(=O)CCC(=O)C([O-])=O | **MAS...LEYS** >> CC(=O)C([O-])=O . N[C@@H](CCC([O-])=O)C([O-])=O

Dataset preparation

ECREACT: Source from Brenda, MetaNetX, PathBank, and Rhea

Annotated reaction SMILES using AA sequences

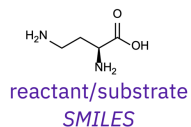
Combined with USPTO [Lowe, 2012] (~1M reactions)



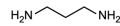
Language modelling for biocatalysed reactions

Training

Enzymatic reaction representation



amino acid sequence



`NCC[C@H](N)C(=O)O|KTYF...SQTSQIHKKDNHIRGQARFCP...YVLK>>NCCCN`

Tokenization



`"[CLS]" "N_" ... "O_" "KTY" ... "YVLK" "N_" ... "N_" "[SEP]"`

Model & training tasks

BERT model

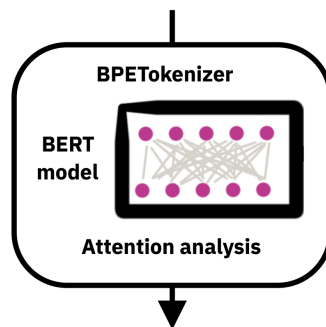


**MLM +
n-gram MLM
tasks**

Learning via MTL on organic and enzymatic reactions

Inference

`NCC[C@H](N)C(=O)O|KTYF...SQTSQIHKKDNHIRGQARFCP...YVLK>>NCCCN`

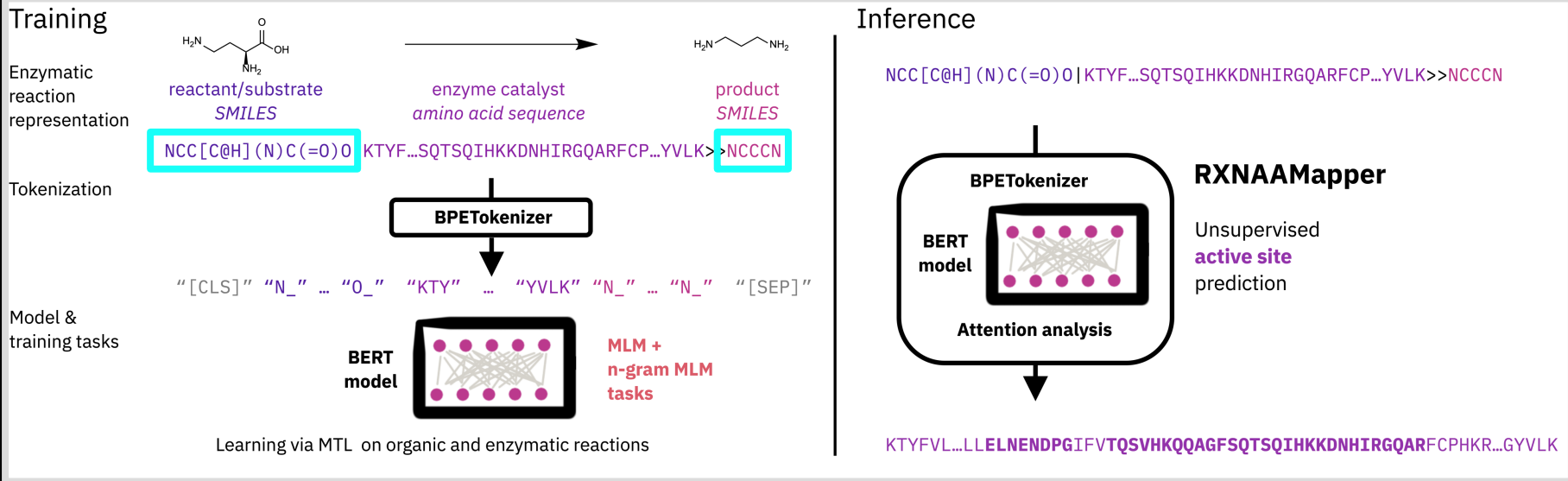


RXNAAMapper

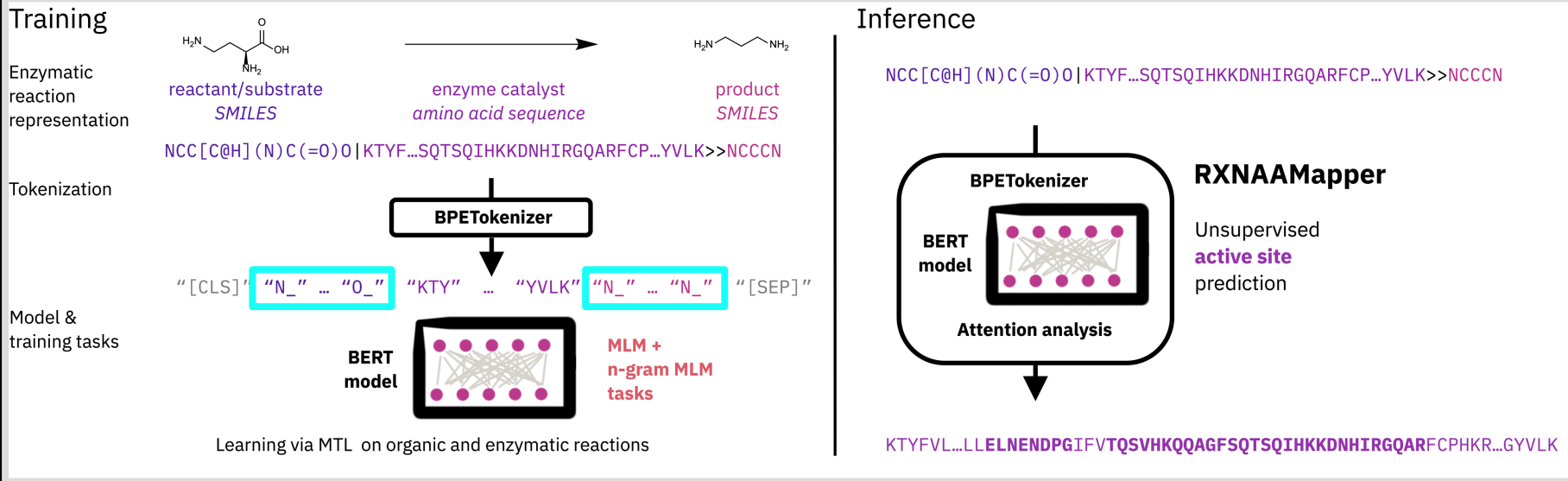
Unsupervised
active site
prediction

`KTYFVL...LLELNENDPGIFVTQSVHKQAGFSQTSQIHKKDNHIRGQARFCPHKR...GYVLK`

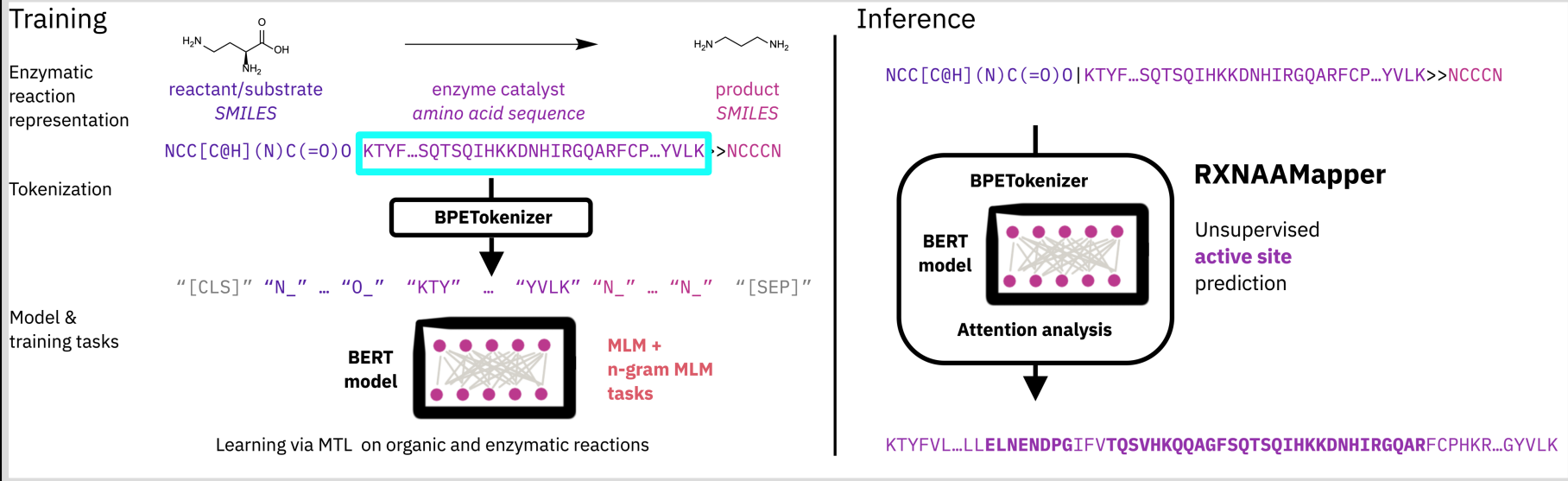
Language modelling for biocatalysed reactions



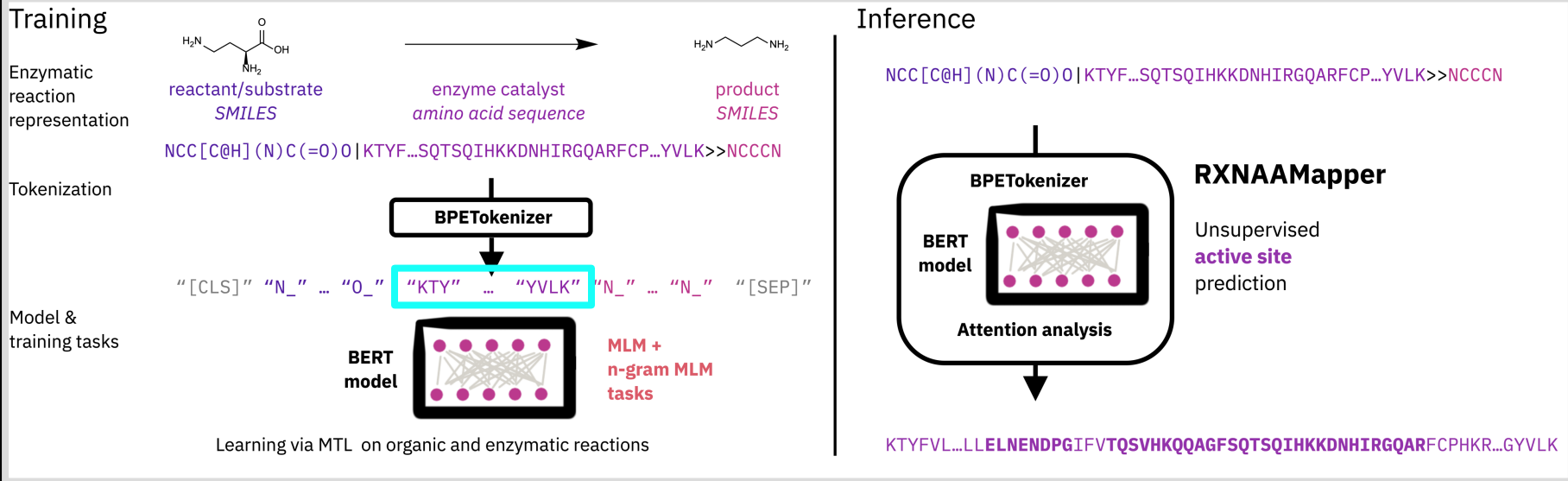
Language modelling for biocatalysed reactions



Language modelling for biocatalysed reactions



Language modelling for biocatalysed reactions



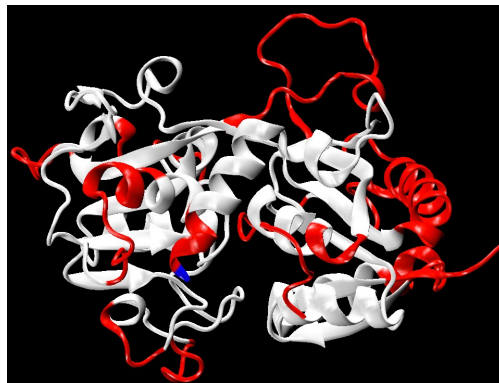
Results

	Overlap Score	False Positive Rate
Random Model	4.98%	84.20%
Pfam	24.01%	78.01%
BERT-base	28.98%	75.56%
RXNAAMapper (ours)	31.51%	66.63%

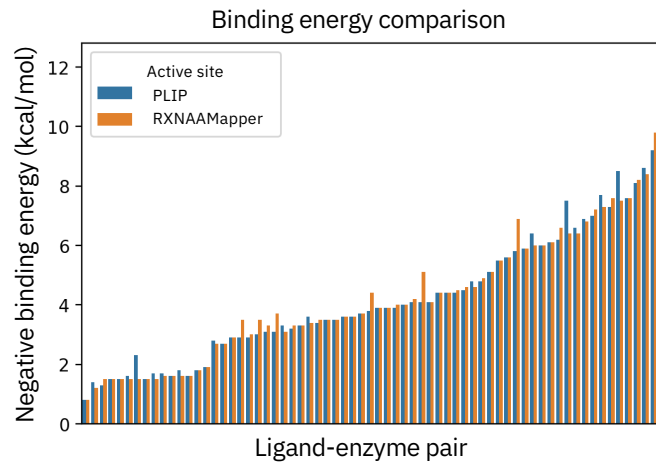
Results

(a)

- Background
- Prediction
- Overlap
- Ground-truth



(b)



Availability

rxn-aa-mapper

Reactions SMILES-AA sequence mapping

setup

```
conda env create -f conda.yml
conda activate rxn_aa_mapper
```

In the following we consider on [example](#)

predict active site

The trained model can be used to map reactant atoms to AA sequence locations that potentially represent the active site.

```
from rxn_aa_mapper.aa_mapper import RXNAAMapper

config_mapper = {
    "vocabulary_file": "./examples/vocabulary_token_75K_min_600_max_750_500K.txt",
    "aa_sequence_tokenizer_filepath": "./examples/token_75K_min_600_max_750_500K.json",
    "model_path": "/tmp/rxnaamapper-pretrained-model",
    "head": 3,
    "layers": [11],
    "top_k": 1,
}

mapper = RXNAAMapper(config=config_mapper)
mapper.get_reactant_aa_sequence_attention_guided_maps(["NC(=O)c1ccc[n+]([C@@H]2O[C@H](COP(
```

<https://github.com/rxn4chemistry/rxnaamapper>

OUTLINE

1. AI models for chemistry
2. Synthesis automation – RoboRXN
3. AI for biocatalysis

IBM RXN for Chemistry
The free AI Tool in the Cloud for Digital Chemistry

Language Models for converting Experimental Procedures, predicting Chemical Reactions or Retrosynthesis Pathways and automating Chemical Synthesis.

[Start your Project Now](#)

[Tweet #rxnforchemistry](#)

We designed a unique tool for digital chemistry, based on the Molecular Transformer, using a simple Ketcher drawing interface and made it available on IBM Cloud to perform a variety of tasks: converting Experimental Procedures into Action Sequences, predicting Chemical Reactions, Retrosynthesis Pathways, Experimental Procedures and automating the compilation and execution of Chemical Syntheses.

Synthesis process for predicting Chemical Reactions

The magic behind the app is a set of Language Models based on Transformers that can predict the most likely outcome of a Chemical Reaction or understand natural language description of Chemical

User interface, freely available on:
rxn.res.ibm.com

rxn4chemistry / rxn4chemistry Public

Python wrapper for the IBM RXN for Chemistry API

93 stars
10 watching
22 forks

Commits

Commit	Author	Date	Commits
druglisberg feat: bumping version 1.6.0	7716584	on 6 May	103
chore: switching from master to main		2 months ago	
chore: bumping package version.		8 months ago	
chore: bumping package version.		8 months ago	
improved handling of errors in API responses (#38)		3 months ago	
feat: bumping version 1.6.0.		2 months ago	
chore: adding sh scripts and fixing deploy.		2 years ago	
feat: initial open source release.		2 years ago	
feat: major release to align with RXN API update.		8 months ago	
feat: initial open source release.		2 years ago	
docs: updated README.md		2 months ago	
feat: handling better None responses.		2 months ago	
feat: major release to align with RXN API update.		8 months ago	
log: removing loguru and using NullHandler.		8 months ago	
log: removing loguru and using NullHandler.		8 months ago	

Python wrapper for the IBM RXN for Chemistry API

[Build and publish rxn4chemistry on PyPI](#) [pip package](#) [License MIT](#) [launch](#) [lander](#)

Access via API / Python wrapper:
github.com/rxn4chemistry/rxn4chemistry

Thank you for your attention!

If you have any questions:

E-mail: ava@zurich.ibm.com

Twitter: [@acvaucher](https://twitter.com/acvaucher)

Acknowledgments:

Antonio Cardinale	Philippe Schwaller
Alessandro Castrogiovanni	Aleksandros Sobczyk
Joppe Geluykens	Alessandra Toniato
Teodoro Laino	Heiko Wolf
Matteo Manica	Federico Zipoli
Vishnu H. Nair	Loïc Kwate Dassi
Yves Gaetan Nana Teukam	Daniel Probst

