

# Human-in-the-loop for a disconnection aware retrosynthesis

Andrea Byekwaso, Alain Vaucher, Philippe Schwaller,  
Alessandra Toniato, Teodoro Laino



@acvaucher

**IBM Research**

*ACS Spring  
March 23, 2022*

# IBM and the data science / chemistry ecosystem?

- Accelerated discovery – **IBM RXN**
- Two sides in our relationship with the ecosystem:
  - **User**: existing data/tools to develop new AI models
  - **Provider**: make technology available & usable via platform
- Interesting interplay between **research** and **platform development!**
  
- Example (today's talk): disconnection aware retrosynthesis

# OUTLINE

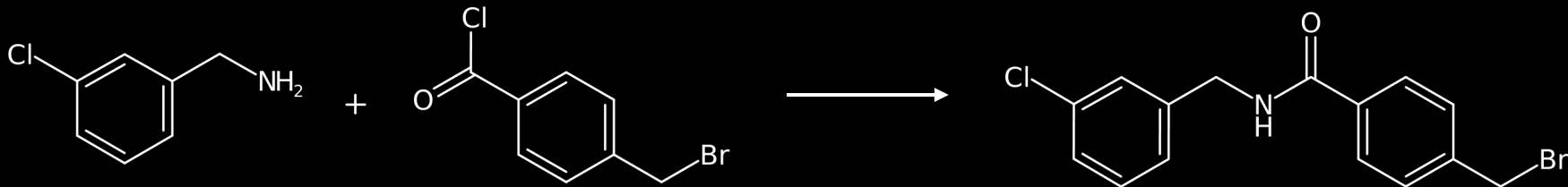
1. IBM RXN research
2. IBM RXN platform
3. Data-driven retrosynthetic models in practice
4. Disconnection aware retrosynthesis

# OUTLINE

1. IBM RXN research
2. IBM RXN platform
3. Data-driven retrosynthetic models in practice
4. Disconnection aware retrosynthesis

# Reaction prediction

Background: 1/7



## Textual representation (SMILES)

NCc1cccc(Cl)c1

O=C(Cl)c1ccc(CBr)cc1

O=C(NCc1cccc(Cl)c1)c1ccc(CBr)cc1

## “Sentence of atoms”

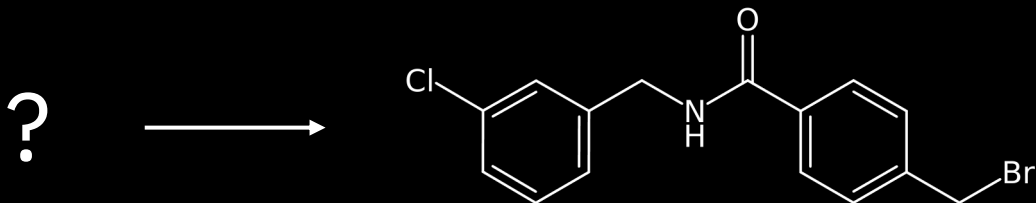
NCc1cccc(Cl)c1.O=C(Cl)c1ccc(CBr)cc1

“Translation”

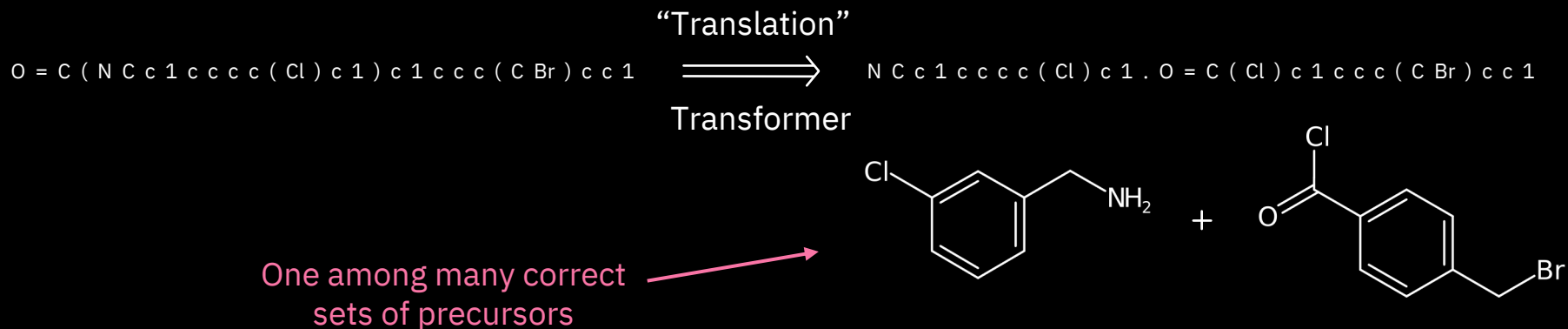
O=C(NCc1cccc(Cl)c1)c1ccc(CBr)cc1

Molecular Transformer

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



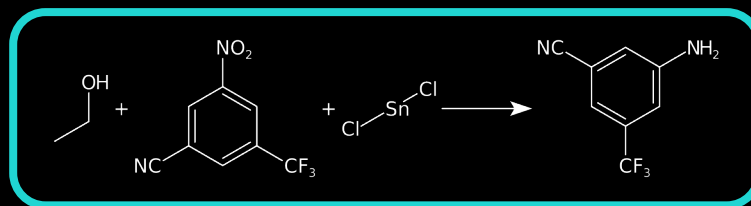
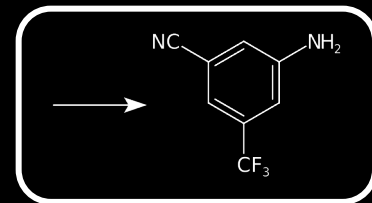
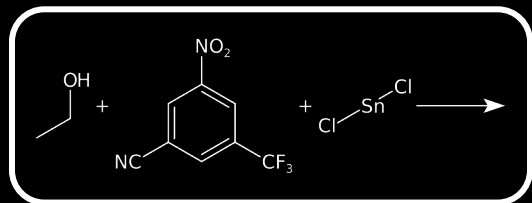
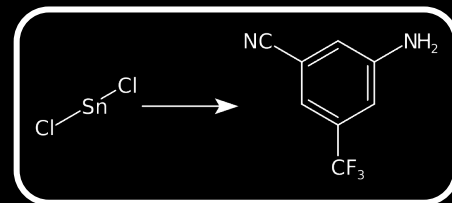
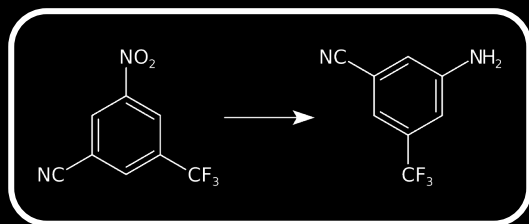
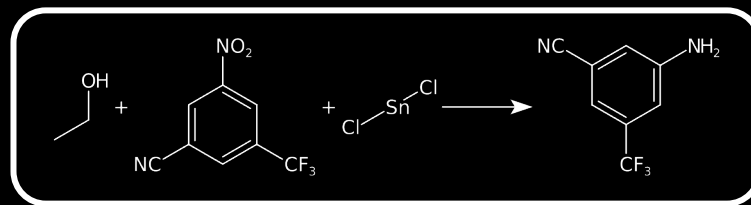
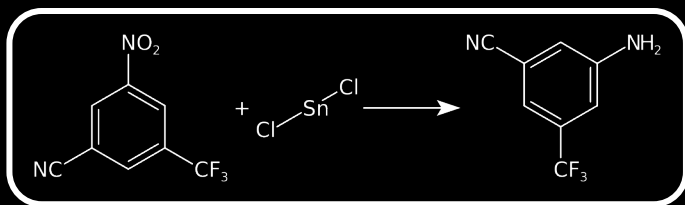
Similar approach, both sides switched



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.

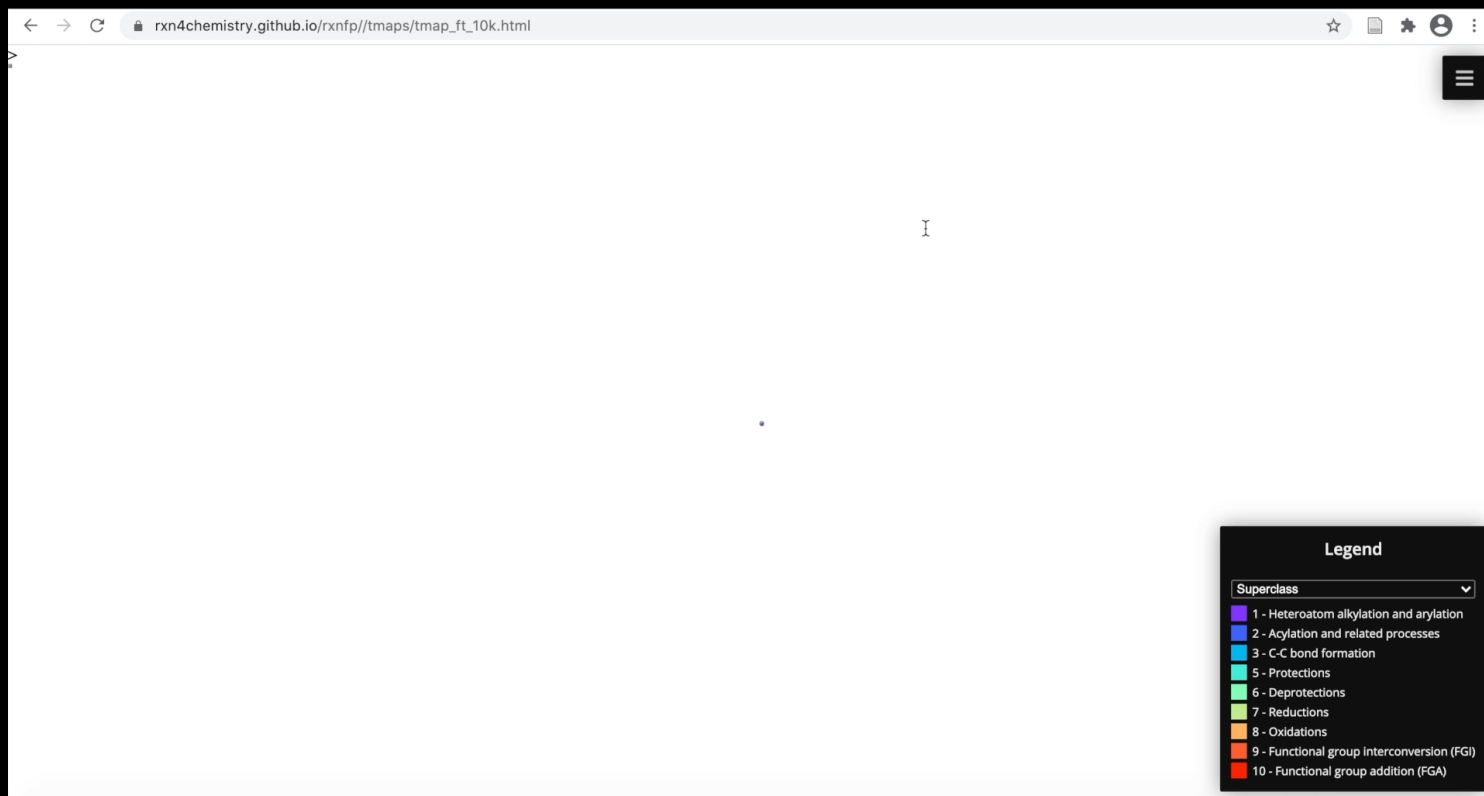
# Completing partial chemical equations

Background: 3/7



# Classifying and mapping reactions

Background: 4/7

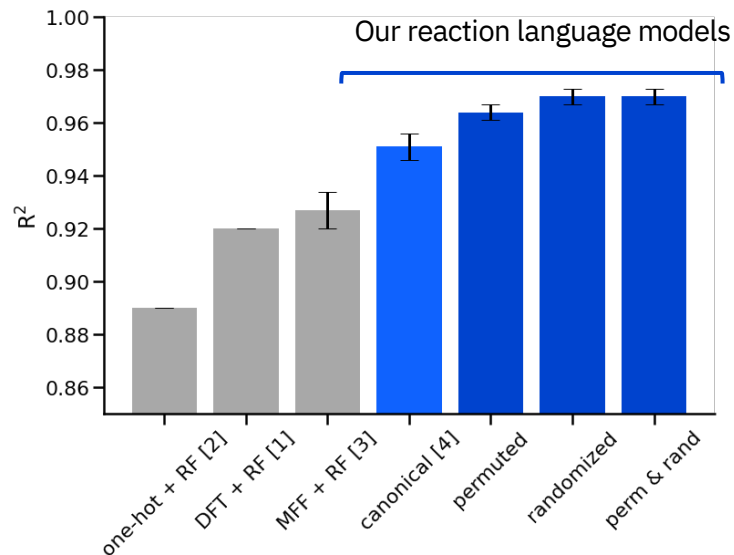
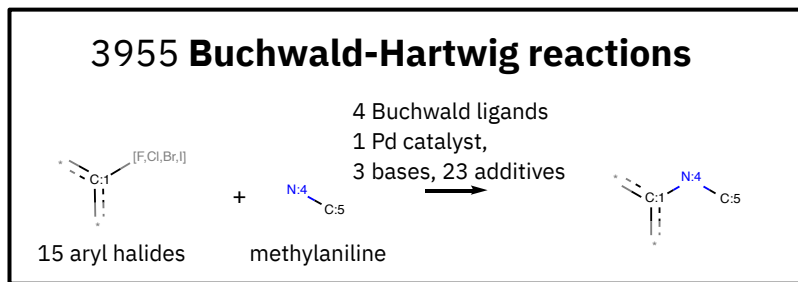


Schwaller, P.; Probst, D.; Vaucher, A. C.; Nair, V. H.; Kreutter, D.; Laino, T. & Reymond, J.-L., *Nat. Mach. Intell.*, **2021**, 3, 144-152.



# Prediction of chemical reaction yields

Background: 5/7



[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).

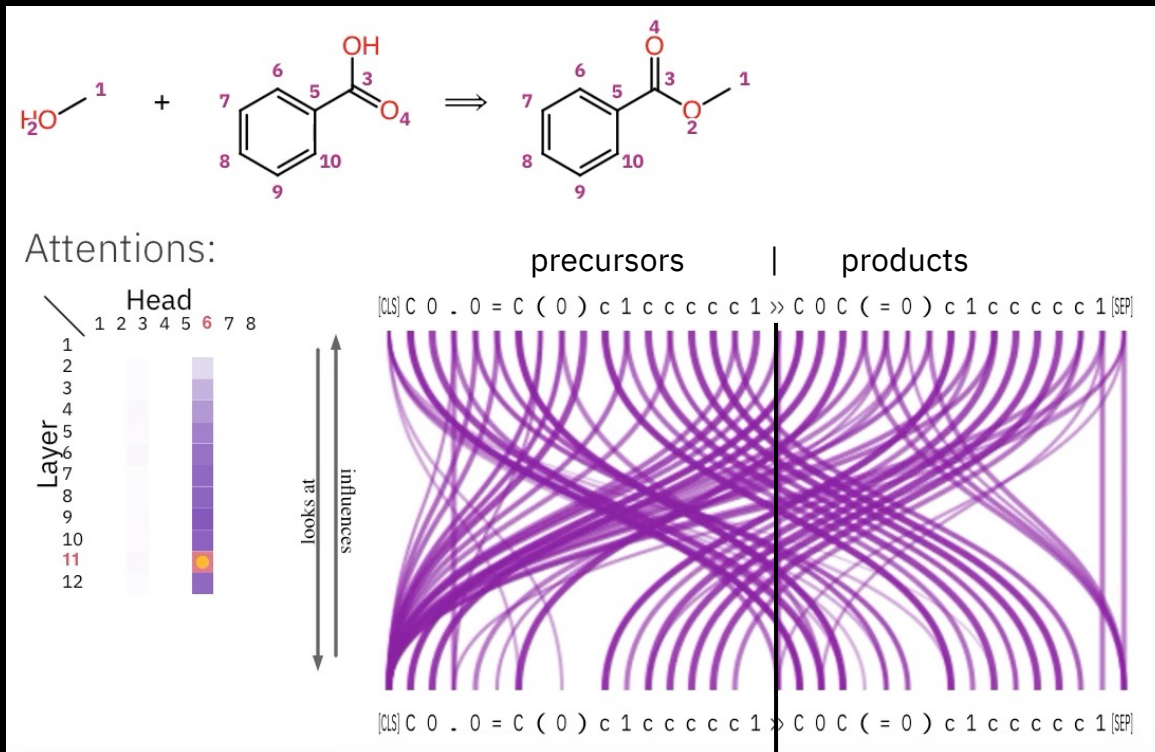
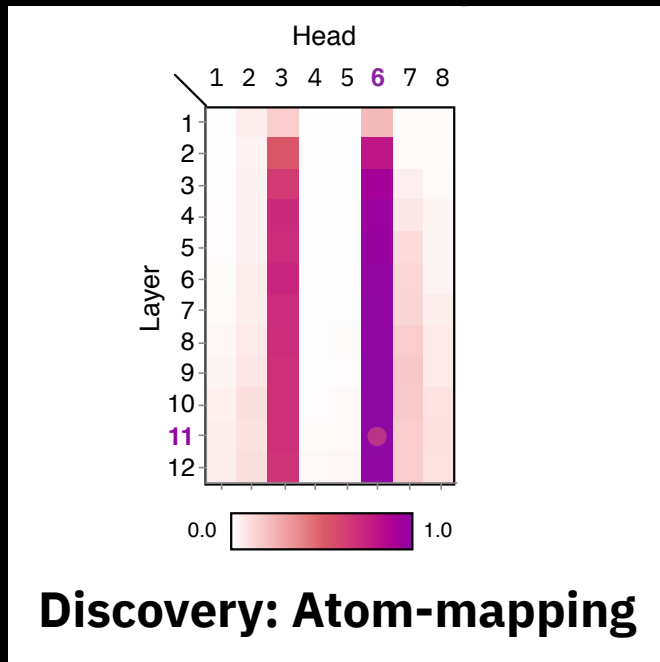
[3] Sandfort, F., Strieth-Kalthoff, F., Kühnemund, M., Beecks, C. & Glorius, F. A structure-based platform for predicting chemical reactivity. *Chem* (2020).

[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.

# Atom mapping: RXNMapper

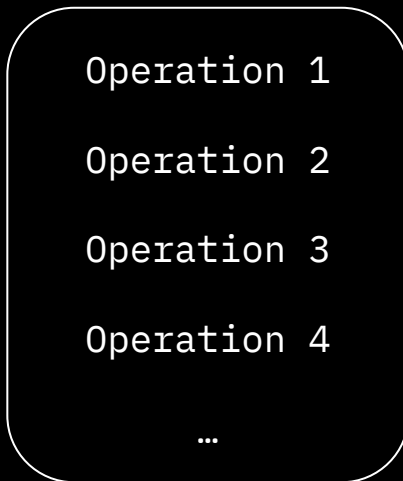
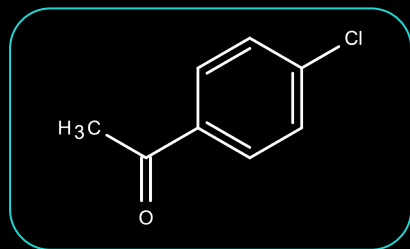
Background: 6/7



Schwaller, P.; Hoover, B.; Reymond, J.-L.; Strobel, H. & Laino, T., *Sci. Adv.*, **2021**, *7*, eabe4166.

# Synthesis actions & synthesis automation

Background: 7/7



C1=CC(C(=O)C)=CC=C1Cl1>>C1=CC(C(=O)C)=CC([N+](=[O-])=O)=C1Cl1

Vaucher, A. C.; Schwaller, P.; Gelykens, J.; Nair, V. H.; Iuliano, A.; Laino, T., *Nat. Commun.*, **2021**, *12*, 2573.

Alain Vaucher / IBM Research Europe / March 23, 2022

# Synthesis actions & synthesis automation

Background: 7/7




# OUTLINE

1. IBM RXN research
2. IBM RXN platform
3. Data-driven retrosynthetic models in practice
4. Disconnection aware retrosynthesis

IBM RXN



Similar reactions list



Status: 0/2 processed


Reaction 1 

Score **1** Reaction class N/A


### Model tuner 1

Name	Creation date	AI model	Status
 data-retro	2022-02-18	Retrosynthetic route prediction	 Ready to run

Data reaction: \_50k.txt 

File preview: 

Items per page: 10 1-1 of 1 item

Reaction 3 

Score **0.999** Reaction class N/A

Freely available  
platform:  
[rxn.res.ibm.com](https://rxn.res.ibm.com)

# API wrapper on GitHub: github.com/rxn4chemistry/rxn4chemistry

Launcher rxn4chemistry\_tour.ipynb Python 3 (ipykernel)

running a reaction prediction is as simple as:

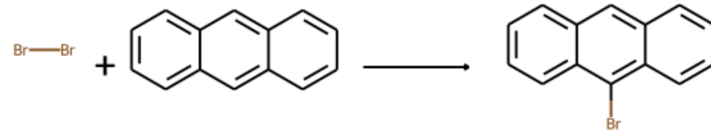
```
[5]: predict_reaction_response = rxn4chemistry_wrapper.predict_reaction(
      'BrBr.c1ccc2cc3ccccc3cc2c1'
    )
```

**NOTE:** we have set limitations on the number of calls per second and per minute in the public version of RXN for Chemistry. These limits can be tweaked or removed in on-premise deployments. Those limitations are currently set to 5 calls per minute, in most cases this is not a problematic limitation.

```
[6]: predict_reaction_results = rxn4chemistry_wrapper.get_predict_reaction_results(
      predict_reaction_response['prediction_id']
    )
```

```
[7]: get_reaction_from_smiles(predict_reaction_results['response']['payload']['attempts'][0]['smiles'])
```

```
[7]:
```



It is possible to run reaction prediction in batches (not storing the information in any project) to use the service in a highthroughput fashion:

```
[8]: predict_rection_batch_response = rxn4chemistry_wrapper.predict_reaction_batch(
```

# OUTLINE

1. IBM RXN research
2. IBM RXN platform
3. Data-driven retrosynthetic models in practice
4. Disconnection aware retrosynthesis



# Retrosynthesis: numerous data-driven models!

TABLE 1 Different retrosynthetic prediction approaches

Input (product)	Output	Single-step approach	Multi-step algorithm
<i>Sequence-based</i>			
SMILES	Reactants SMILES	Seq-2-seq LSTM	None: 56
SMILES	Largest reactant SMILES	Transformer	None: 83
SMILES	Reactants SMILES	Transformer	None: 83,95–99, MCTS: 112
SMILES	1: Synthons prediction, 2: Synthons completion	Transformer	None: 105
MACCS keys	Reactants MACCS keys	Transformer	None: 100
SMILES	Precursors SMILES	Transformer	Beam search: 28
<i>Fingerprint-based</i>			
Fingerprint	Reaction template	Similarity	None: 93
Fingerprint	Reaction template	Feed-forward NN	None: 78,106,107, MCTS: 23,29,109,110, RL: 114, A*: 115,116
Fingerprint	Reaction template	Modern Hopfield network	None: 91
<i>Graph-based</i>			
Molecular graph	Reaction template	Graph neural network	None: 52,101, SMC: 117
Molecular graph	1: Synthons prediction, 2: Synthons completion	Graph neural network	None: 102,103,118
Molecular graph	Sequence of graph-edits	Graph neural network	None: 30


Only a selection!

Schwaller, P. et al., Machine intelligence for chemical reaction space, *WIREs Comput. Mol. Sci.*, **2022**, e1604.

# Retrosynthesis: recent models

PAPER • OPEN ACCESS

## Chemformer: a pre-trained transformer for computational chemistry

Ross Irwin<sup>1</sup>, Spyridon Dimitriadis<sup>1,2</sup>, Jiazhen He<sup>1</sup> and Esben Jannik Bjerrum<sup>3,1</sup> 

 > [physics](#) > P

Physics > Chemical P

[Submitted on 29 Jan 2022]

**Retrofor**  
**Retrosyn**

Yue Wan, Be

## AI-Driven Synthetic Route Design Incorporated with Retrosynthesis Knowledge

Shoichi Ishida, Kei Terayama, Ryosuke Kojima, Kiyosei Takasu, and Yasushi Okuno\*

[Cite this article](#) | [Chem Inf Model](#) 2022, XXXX, XXX | [Article Views](#) | [Altmetric](#) | [Citations](#)

Research article | [Open Access](#) | [Published: 15 March 2022](#)

## Improving the performance of models for one-step retrosynthesis through re-ranking

[Min Htoo Lin](#), [Zhengkai Tu](#) & [Connor W. Coley](#) 

[Journal of Cheminformatics](#) **14**, Article number: 15 (2022) | [Cite this article](#)

**269** Accesses | **9** Altmetric | [Metrics](#)


Share Add to Export



# Retrosynthesis: metrics?

PAPER • OPEN ACCESS

## Chemformer: a pre-trained transformer for computational chemistry

Ross Irwin<sup>1</sup>, Spyridon Dimitriadis<sup>1,2</sup>, Jiazhen He<sup>1</sup> and Esben Jannik Bjerrum<sup>3,1</sup> 

arXiv > physics > P

Physics > Chemical P

[Submitted on 29 Jan 2022]

Retrofor  
Retrosyn

Yue Wan, Be

## AI-Driven Synthetic Route Design Incorporated with Retrosynthesis Knowledge

Shoichi Ishida, Kei Takayama, Daisuke Kojima, Kinsei Takasugi and Yasushi Okuno\*

# Metrics?

Research article | [Open Access](#) | [Published: 15 March 2022](#)

## Improving the performance of models for one-step retrosynthesis through re-ranking

[Min Htoo Lin](#), [Zhengkai Tu](#) & [Connor W. Coley](#) 

[Journal of Cheminformatics](#) **14**, Article number: 15 (2022) | [Cite this article](#)

**269** Accesses | **9** Altmetric | [Metrics](#)

Share Add to Export



# Retrosynthesis: metrics?

“On direct synthesis and retrosynthesis prediction benchmark datasets we publish state-of-the-art results for **top-1 accuracy**.”

AI-Driven Synthetic Route Design Incorporated with Retrosynthesis Knowledge

“We adopt the conventional **top-k accuracy** of the full reactants to evaluate the retrosynthesis performance.”

“Typically, these data-driven methods are evaluated in terms of **top-N accuracy**”

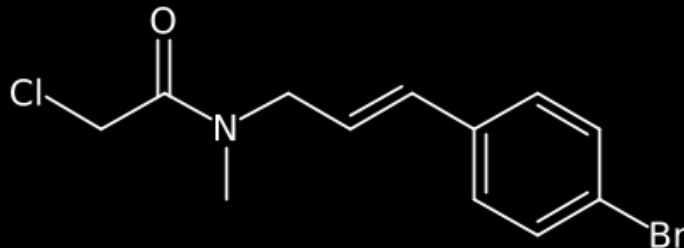
“We acknowledge that **top-N accuracy** alone does not paint a complete picture of a one-step model’s performance, as others have also argued.”

# Retrosynthesis: metrics?

- top-N accuracy is not fully satisfactory – a necessary evil
- Multi-step: hard to assess as well!
  - Top-N
  - Turing test
  - Percentage of solved molecules
- What do the chemists need?

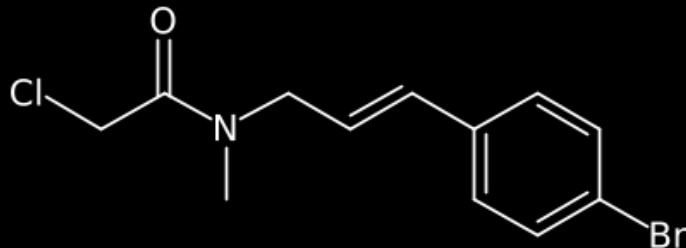
# Our experience with IBM RXN

- DEMO: CN(C/C=C/c1ccc(Br)cc1)C(=O)CCl



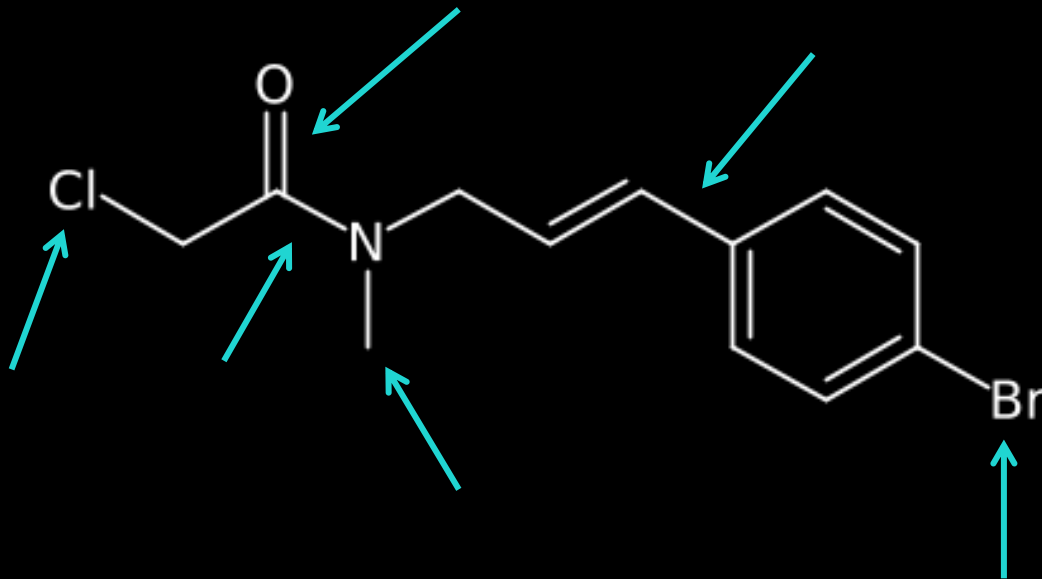
# Our experience with IBM RXN

- Many chemists prefer “interactive” mode
- DEMO: CN(C/C=C/c1ccc(Br)cc1)C(=O)CCl



# Even more interactive control?

- Let the chemists decide where to break the compound?





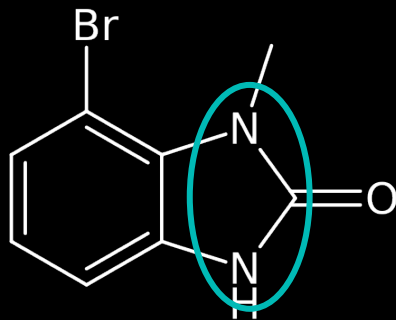
# OUTLINE

1. IBM RXN research
2. IBM RXN platform
3. Data-driven retrosynthetic models in practice
4. Disconnection aware retrosynthesis

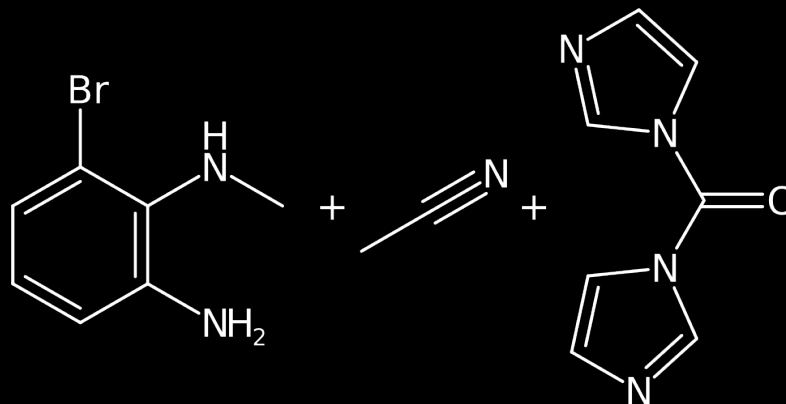
# Disconnection-aware retrosynthesis

Goal:

Input (target compound)



Output (precursors)

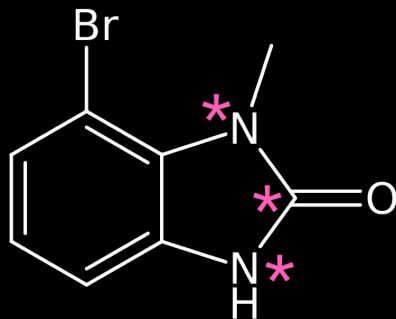


# Machine learning model

## Inspired by the **Molecular Transformer**

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

Input



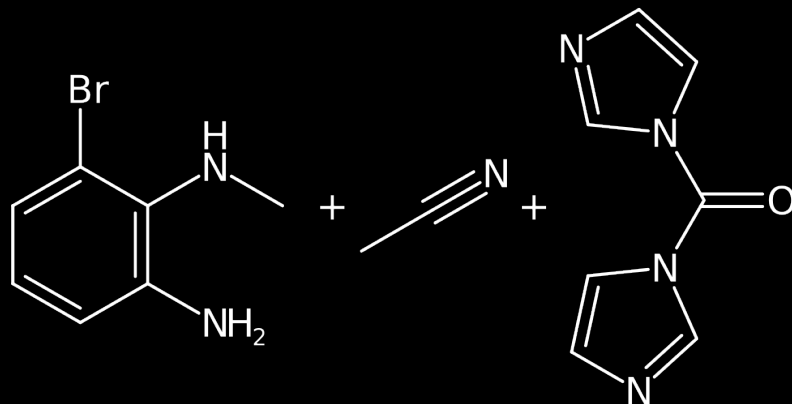
Cn1c2c(Br)cccc2[nH]c1=O

C[n:1]1c2c(Br)cccc2[nH:1][c:1]1=O

“Translation”



Output

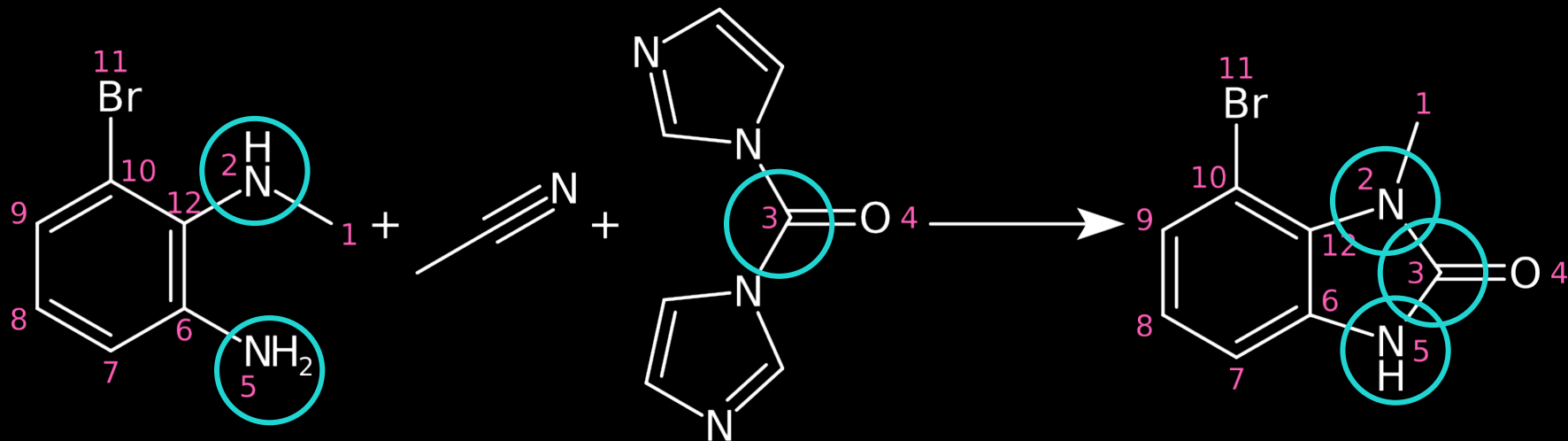


CNc1c(N)cccc1Br.CC#N.O=C(n1ccnc1)n1ccnc1

\*Not shown here for readability: the Model uses tokenized SMILES strings: “C[n:1]1c2c(Br)...” → “C [n:1] 1 c 2 c ( Br ) ...”

# Dataset generation

- Start from atom-mapped reaction:



# Data and model

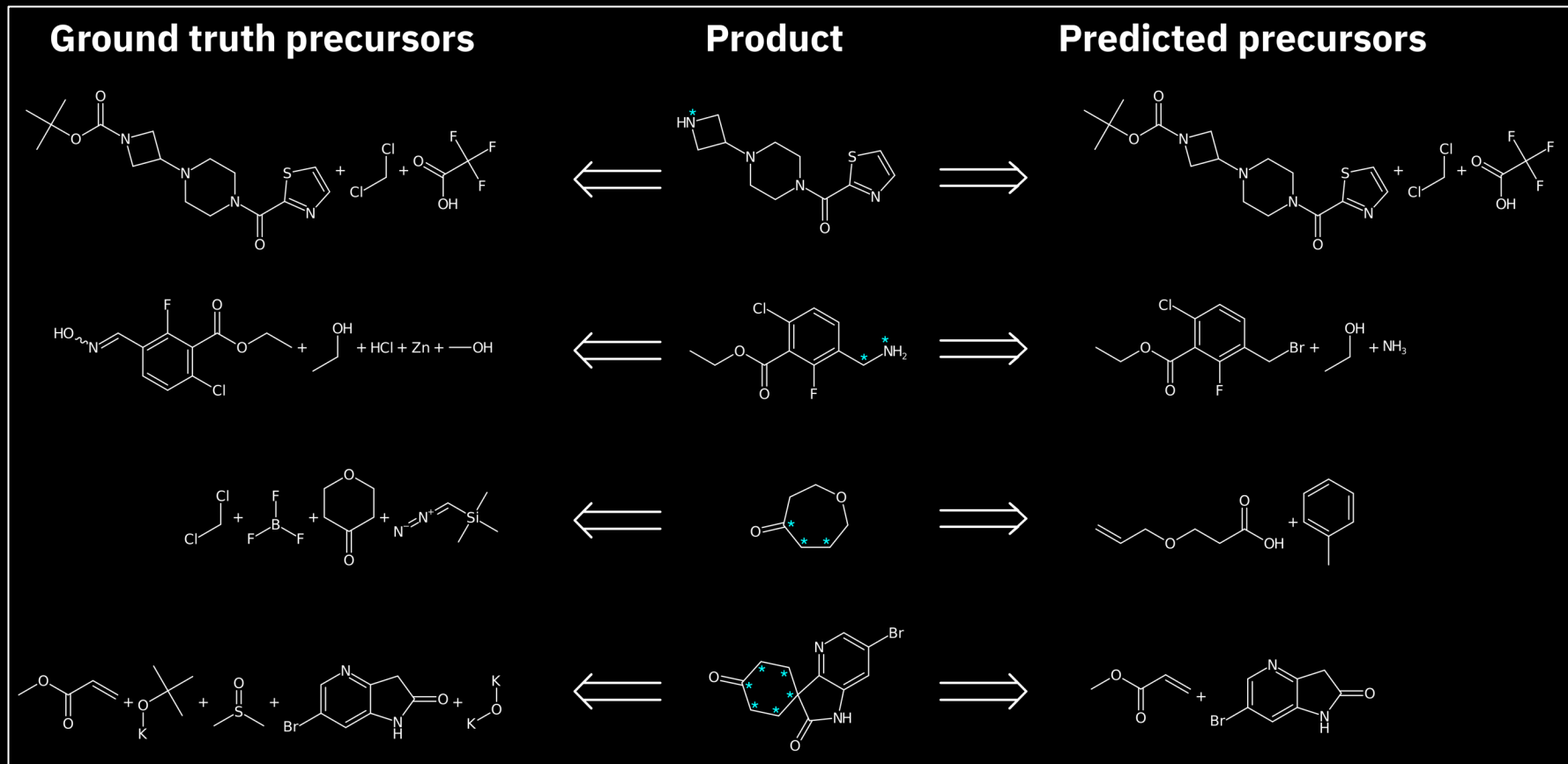
- >2M **patent reactions** from Pistachio [1]
- Atom-mapped with **RXNMapper** [2]
- Training, validation, and test sets of sizes **2.27M**, **10.0k**, and **126k**.
- **Transformer-based seq-2-seq** model implemented with OpenNMT [3]

[1] Nextmove Software Pistachio, <http://www.nextmovesoftware.com/pistachio> (Accessed Sep 23, 2021).

[2] Schwaller, P.; Hoover, B.; Reymond, J.-L.; Strobelt, H. & Laino, T., *Sci. Adv.*, **2021**, 7, eabe4166..

[3] Klein, G. et al, OpenNMT: Open-Source Toolkit for Neural Machine Translation. ACL 2017.

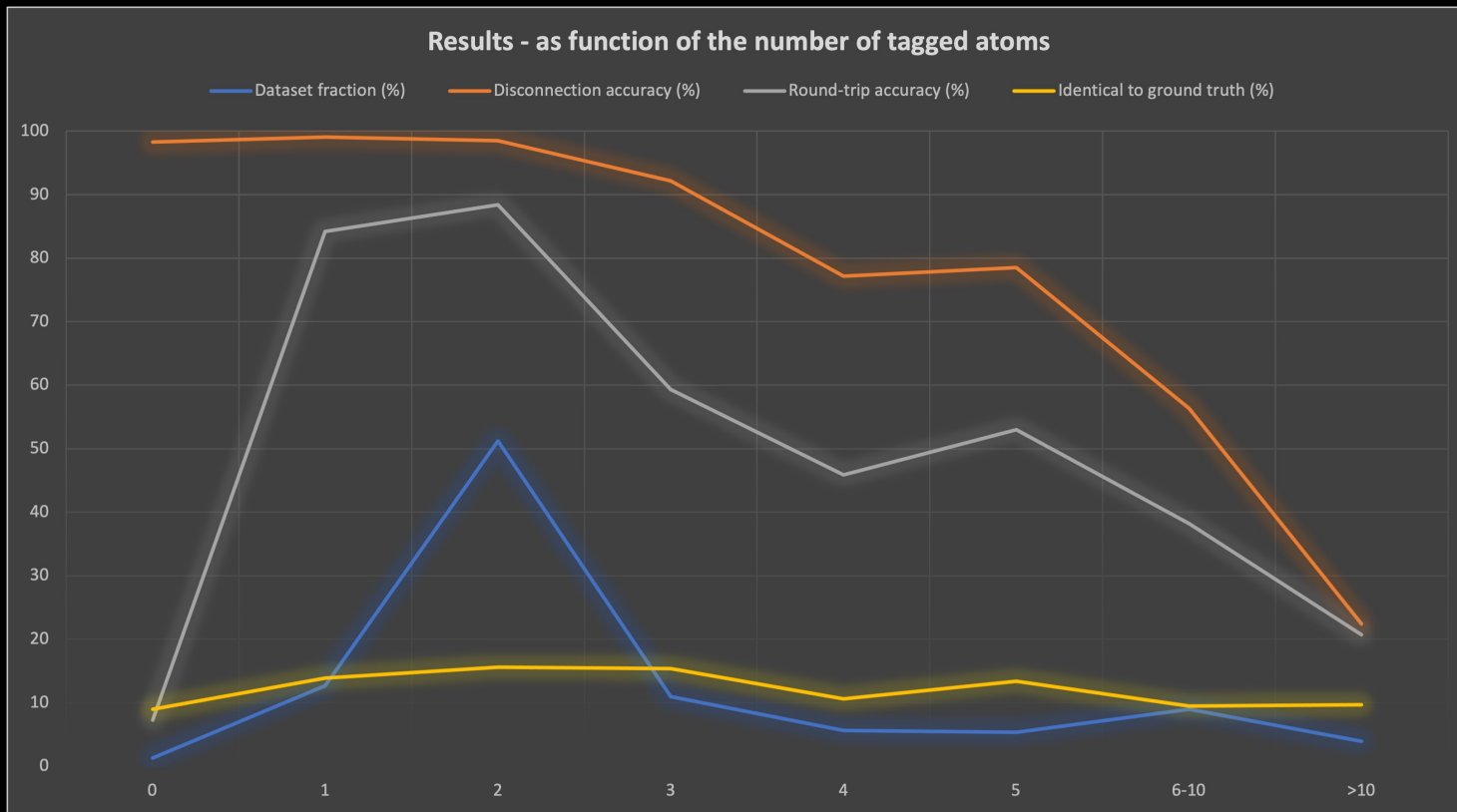
# Results: examples



# Results

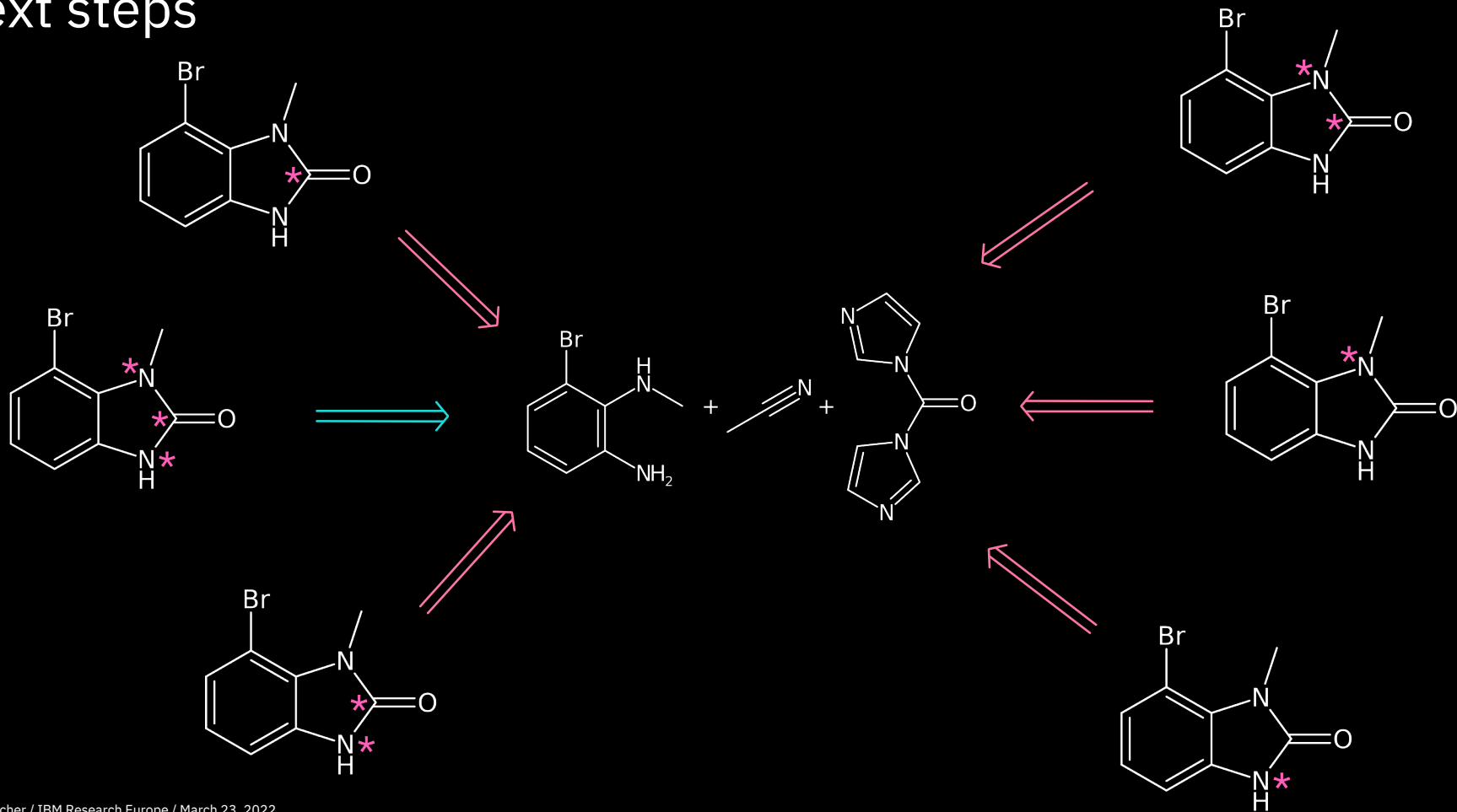
<b>Number of tagged atoms</b>	<b>Dataset fraction (%)</b>	<b>Disconnection accuracy (%)</b>	<b>Round-trip accuracy (%)</b>	<b>Identical to ground truth (%)</b>
0	1.29	98.3	7.3	9.0
1	12.65	99.1	84.2	13.9
2	51.22	98.5	88.4	15.6
3	10.99	92.2	59.3	15.4
4	5.64	77.2	45.9	10.6
5	5.34	78.5	53.0	13.4
6-10	8.96	56.3	38.2	9.5
>10	3.91	22.4	20.7	9.7
<b>Overall</b>	<b>100.0</b>	<b>88.9</b>	<b>76.0</b>	<b>14.1</b>

# Results





# Next steps



# Thank you for your attention!

## Questions or comments

*E-mail:* [ava@zurich.ibm.com](mailto:ava@zurich.ibm.com)

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

Preprint with initial results: [ibm.biz/disconnection-aware-retro](https://ibm.biz/disconnection-aware-retro)

