

IBM RXN for Chemistry: A Chemical Reaction Prediction Platform

Philippe Schwaller ([@pschwillr](https://twitter.com/pschwillr))

<http://dx.doi.org/10.26434/chemrxiv.7297379>

IBM Research AI

Theophile Gaudin, Teodoro Laino &
Costas Bekas
IBM Research Zurich, Switzerland



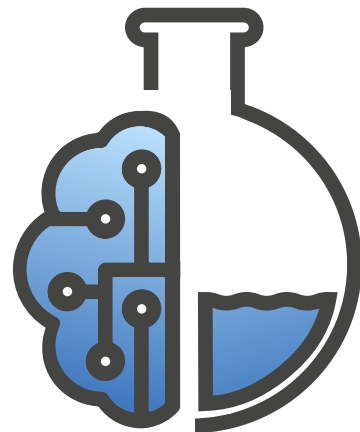
UNIVERSITY OF
CAMBRIDGE



THE WINTON PROGRAMME FOR THE

Physics of Sustainability

Peter Bolgar & Alpha Lee
University of Cambridge



Exploring the nearly endless chemical and materials space



In silico molecules / materials
VAE / GAN / RL

What to make?

Design

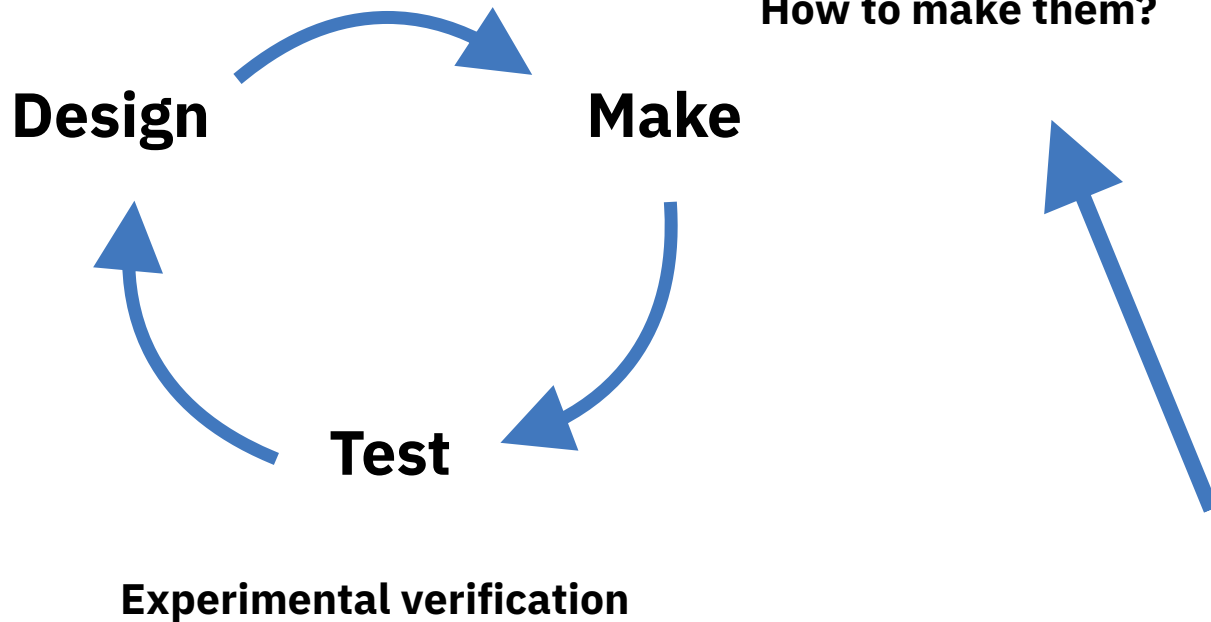
Make

Test

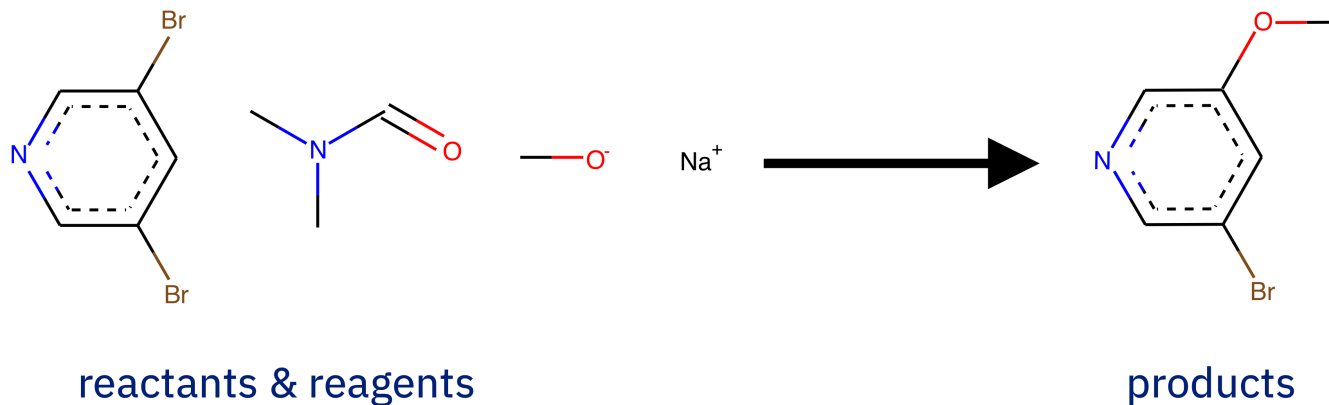
Experimental verification

Synthetic route planning
Reaction outcome prediction

How to make them?

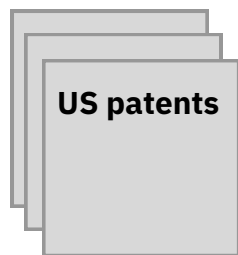


Data-Driven Chemical Reaction Prediction



Data > Representation > Model > Results

Data



Lowe (2012,
2017)



text-mining

Reaction SMILES
- reactants >> products
CML files

filtered



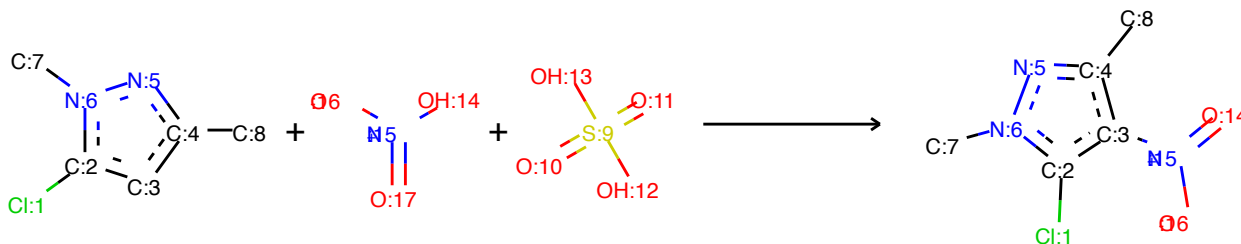
Benchmark sets

Jin et al. **USPTO_MIT**
NeurIPS 2017

Schwaller et al. **USPTO_STEREO**
NeurIPS 2017

Bradshaw et al. **USPTO_LEF**
ArXiv 2018

[Cl:1][c:2]1[cH:3][c:4]([CH3:8])[n:5][n:6]1[CH3:7].[OH:14][N+:15]([O-:16])=[O:17].[S:9](=[O:10])
(=[O:11])([OH:12])[OH:13]>>[Cl:1][c:2]1[c:3]([N+:15])(=[O:14])[O-:16])[c:4]([CH3:8])[n:5][n:6]1[CH3:7]

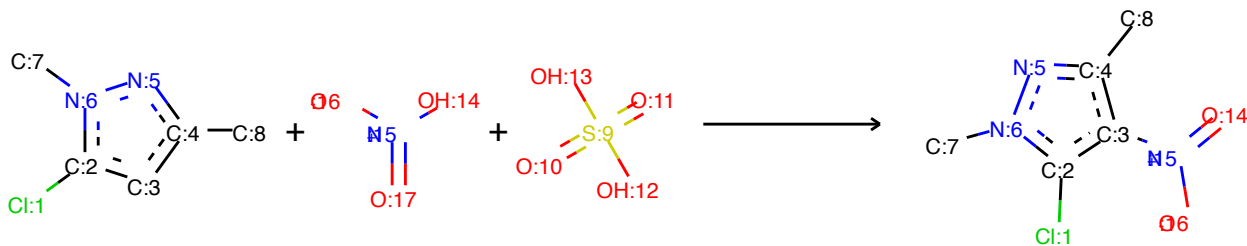


Atom-mapping from Indigo TK

- USPTO dataset description:

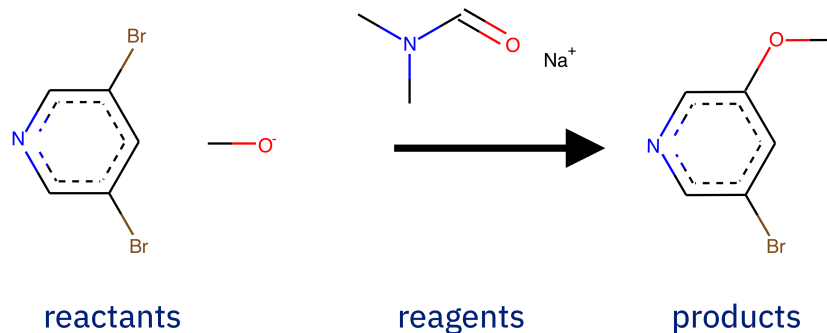
While typically correct, the **atom-maps are wrong in many cases** and hence **should not be entirely relied on**.

[Cl:1][c:2]1[cH:3][c:4]([CH3:8])[n:5][n:6]1[CH3:7].[OH:14][N+:15]([O-:16])=[O:17].[S:9](=[O:10])(=[O:11])([OH:12])[OH:13]>>[Cl:1][c:2]1[c:3]([N+:15](=[O:14])[O-:16])[c:4]([CH3:8])[n:5][n:6]1[CH3:7]

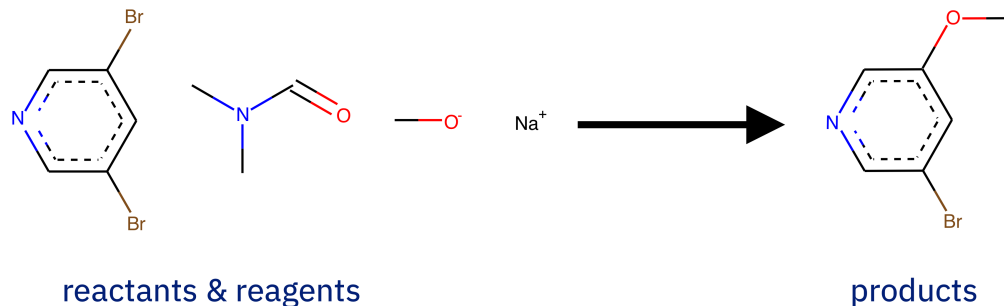


Separated vs mixed reagents

Separated



Mixed



No distinction between reactants and reagents

Representing molecules for ML

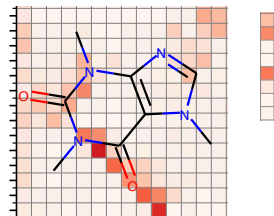
Molecular fingerprints

000010000....0100

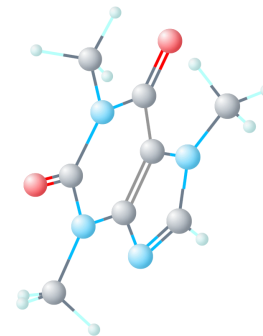
Text-based representations

- SMILES / INCHI
- CN1C=NC2=C1C(=O)N(C(=O)N2C)C

Graph-based



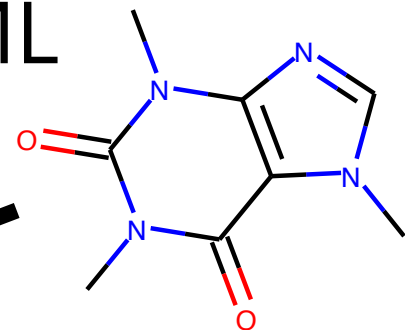
3D structure & surface



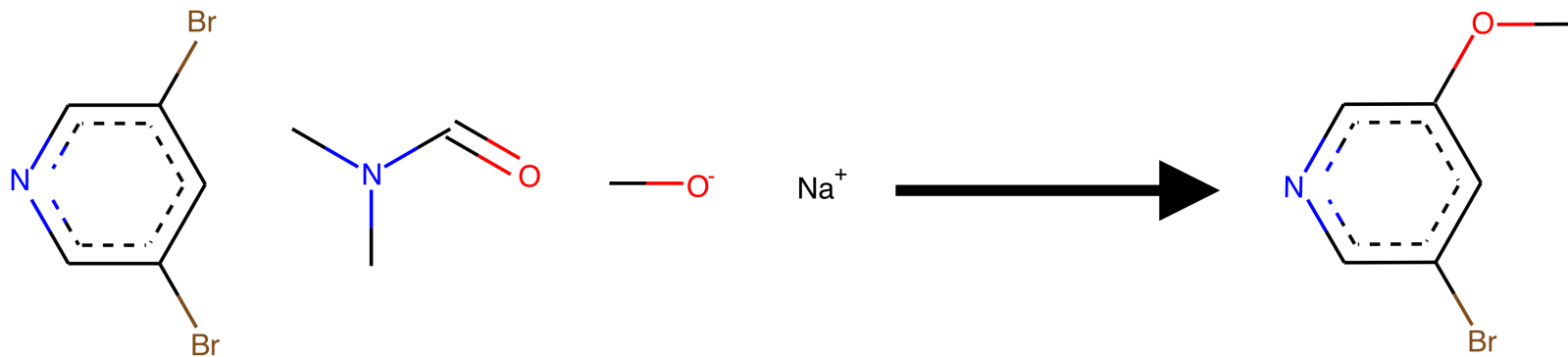
1D

2D

3D



Atoms as letters, molecules as words

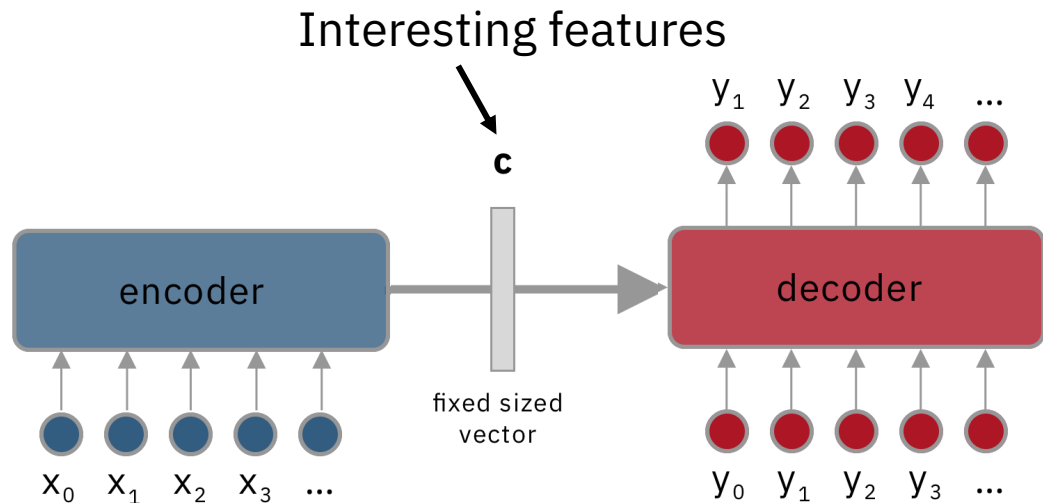


Br c1 c n c c (Br) c 1 . C N (C) C = O . C [O-] . [Na+]

C O c 1 c n c c (Br) c 1

SMILES to SMILES prediction with sequence-2-sequence models

Model: Standard Seq-2-Seq



Problem: fixed size

INPUTS = reactants + reagents

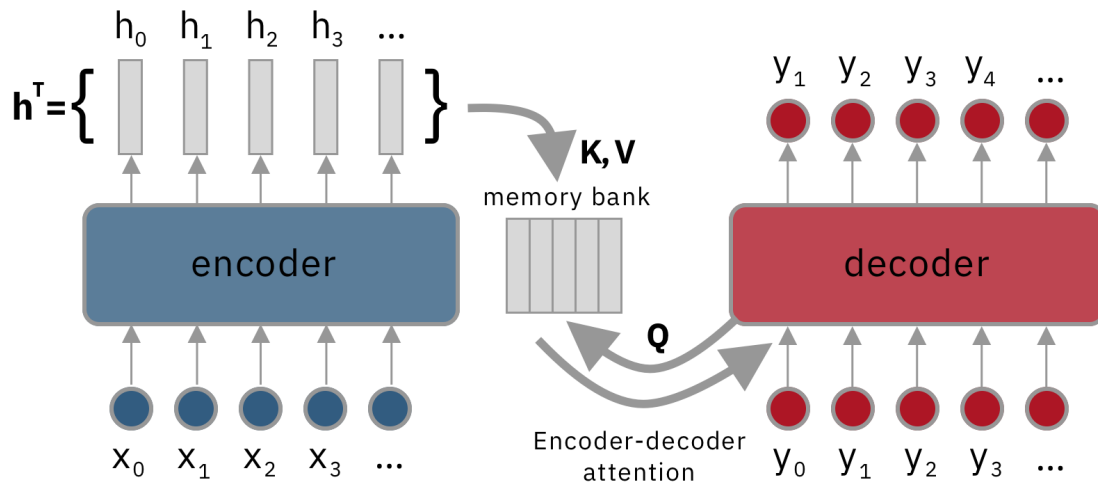
Br c 1 c n c c (Br) c 1 . C N (C) C = O . C [O -] . [Na +]

OUTPUTS = products

C O c 1 c n c c (Br) c 1

Seq-2-Seq with Attention

One state per input



Attention = ability to selectively concentrate on one aspect of context

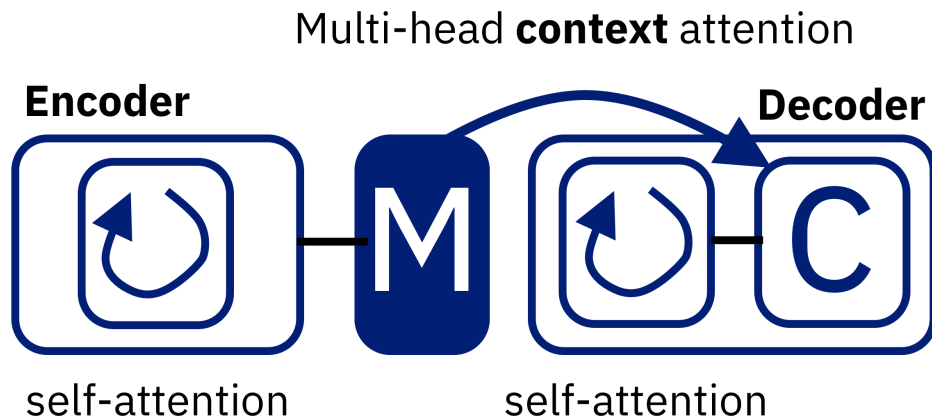
Transformer Architecture

INPUTS = reactants + reagents

Br c 1 c n c c (Br) c 1 . C N (C) C = O . C [O -] . [Na +]

OUTPUTS = products

C O c 1 c n c c (Br) c 1



- NO recurrent neural networks
- Stacks of **attention layers**
- **Multi-head** attention

Attention Is All You Need – Vaswani et al. NeurIPS 2017

Attention types in Transformer

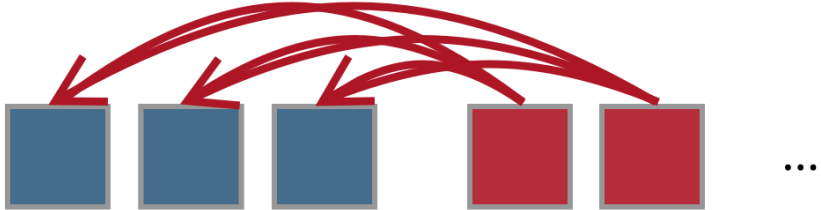
a) Self-attention



b) Masked self-attention



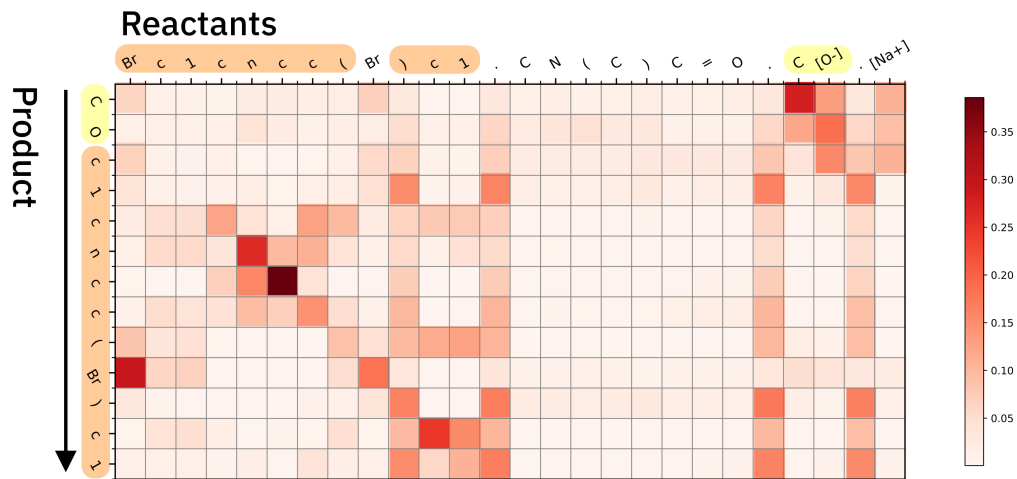
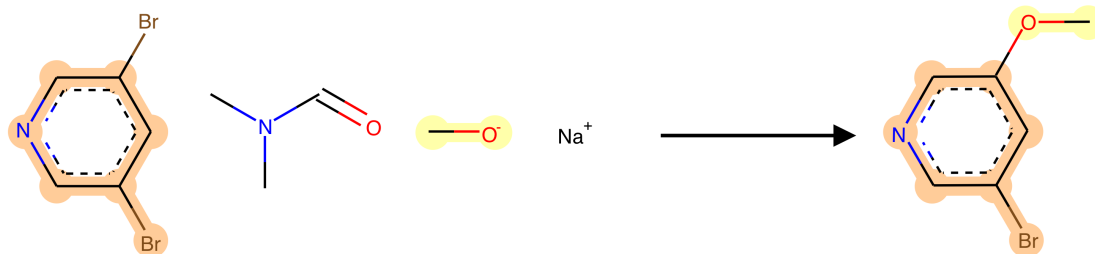
c) Encoder-decoder attention



Encoder states

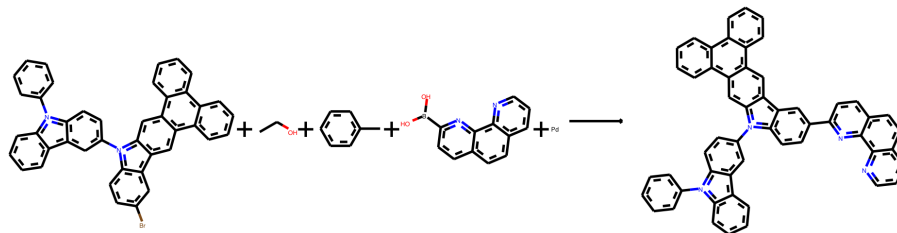
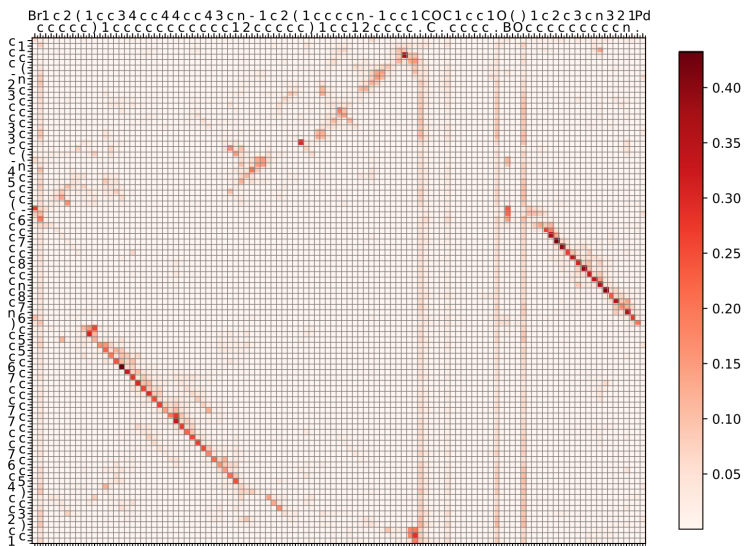
Decoder states

Visualizing encoder-decoder attention



Plotting the attention weights at every decoder time step.

Bromo-Suzuki coupling reaction



INPUT:

Brc1ccc2c(c1)c1cc3c4ccccc4c4ccccc4c3cc1n2-c1ccc2c(c1)c1cccc1n2-c1cccc1.CCO.Cc1cccc1.OB(O)c1ccc2ccc3cccnc3c2n1.[Pd]

OUTPUT:

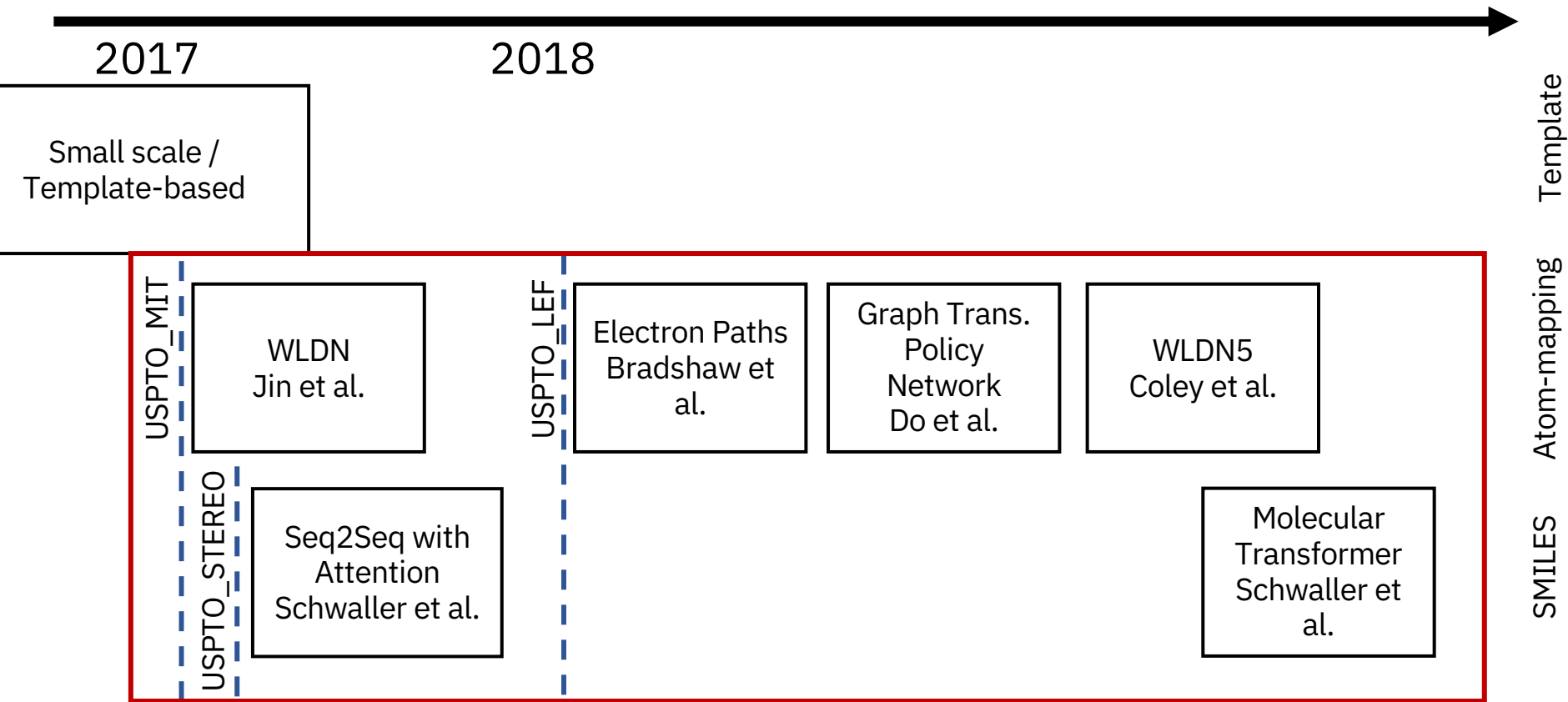
c1ccc(-n2c3ccccc3c3cc(-n4c5ccc(c6ccc7ccc8cccnc8c7n6)cc5c5cc6c7ccccc7c7cccc7c6cc54)ccc32)cc1



Features of Smiles-2-Smiles Approach

- Trained end-2-end
- Fully data-driven
- Rule / template-free
- Physics agnostic model
- Attention weights & confidence score

Recent Data-Driven Reaction Prediction on open-source USPTO dataset



Top-1 accuracy [%] comparison

		require atom-mapping in training set					
USPTO*		S2S [1]	WLDN[2]	ELEC [3]	GTPN [4]	WLDN5 [5]	MT
<u>_MIT_sep</u>	500k	80.3	79.6		82.4	85.6	90.4
<u>_MIT_mixed</u>			74.0				88.6
<u>_LEF_sep</u>	350k		84.0	87.0	87.4	88.3	92.0
<u>_LEF_mixed</u>							90.3

[1] Schwaller et al.: Chem. Sci., 2018, 9, 6091-6098

[2] Jin et al.: NIPS, 2017, 30, 2607-2616

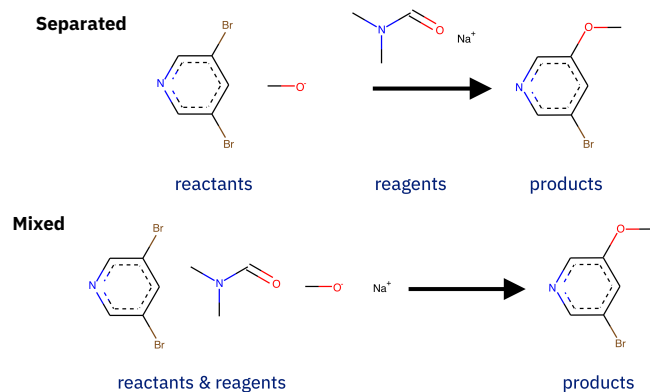
[3] Bradshaw et al.: [arXiv:1805.10970](https://arxiv.org/abs/1805.10970)

[4] Do et al.: arXiv:1812.09441

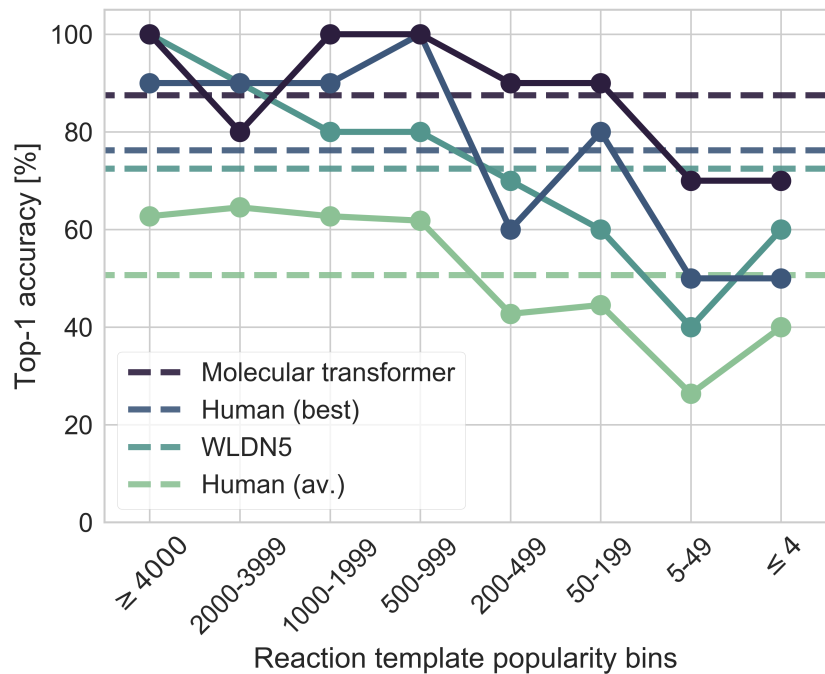
[5] Coley et al. : Chem. Sci., 2019, 10, 370-377

MT

**Molecular Transformer for Chemical Reaction
Prediction and Uncertainty Estimation**



Human prediction benchmark



87.5 % Molecular Transformer

76.5 % best human

72.5 % Coley et al. model

50.6 % average human

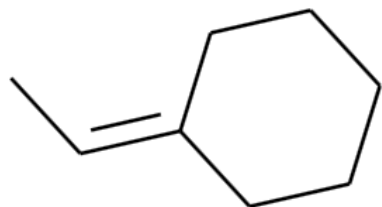
- **80 reactions** (10 reactions per bin)
- Given to **11 chemists**
- **Mixed** input setting for all

Original study performed by
Coley et al. 2018 – WLDN5

common ←————→ **rare**

What else can we do?

Reaction scoring



+ HCl

Markovnikov

Anti-Markovnikov

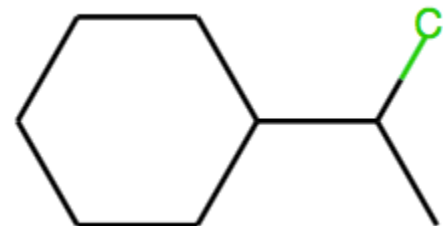
When we provide the model with **reactants**>>**products**

CC=C1CCCCC1.Cl>>CCC1(Cl)CCCCC1



Score: 0.99

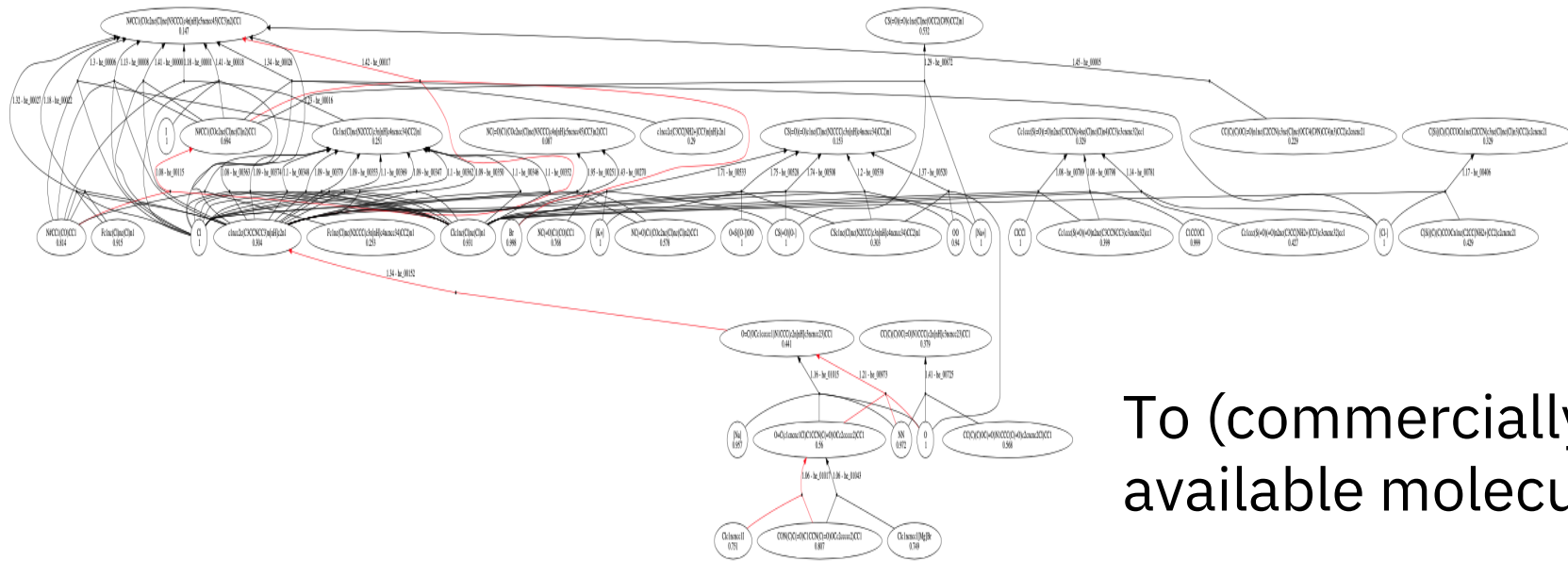
CC=C1CCCCC1.Cl>>CC(Cl)C1CCCCC1



Score: 0.001

Single step validation in retrosynthesis

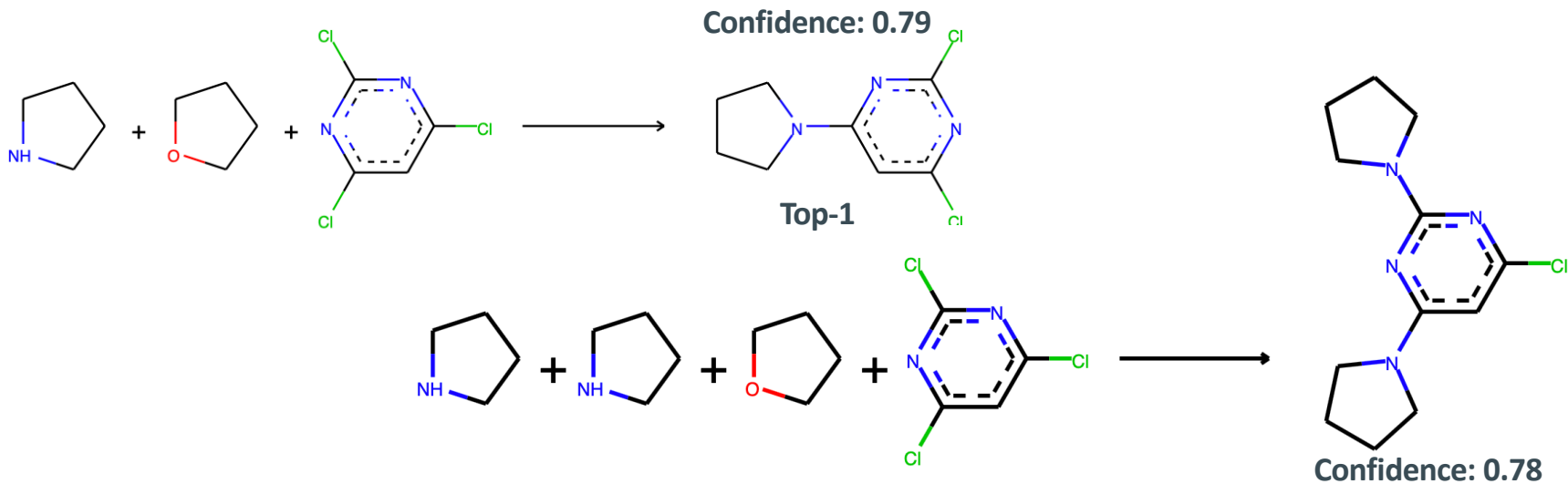
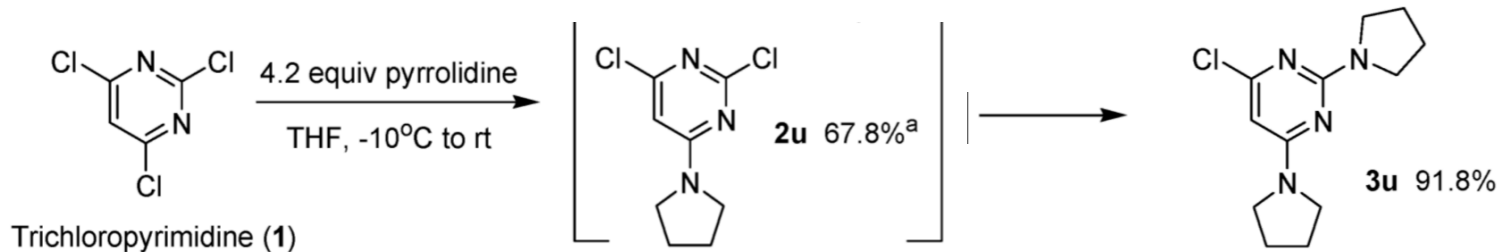
Target



To (commercially) available molecules

Regioselectivity – Halogenated Pyrimidines

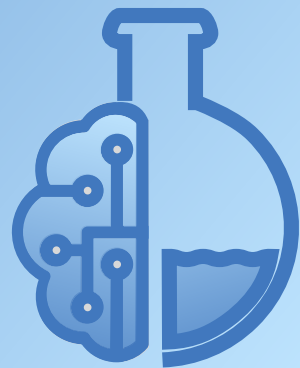
Org. Process Res. Dev., **2006**, *10* (5), pp 921–926 DOI: 10.1021/op060093q

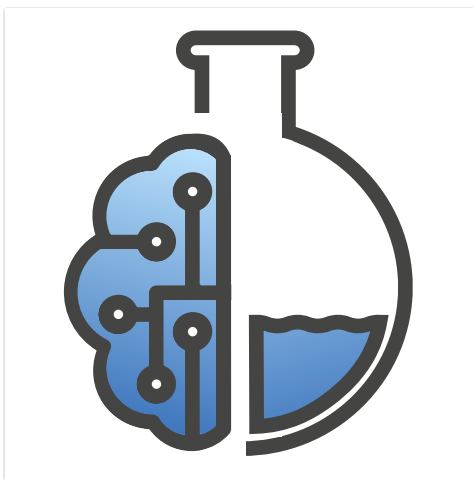


IBM RXN for Chemistry

Freely available on:
rxn.res.ibm.com

#RXNFORCHEMISTRY

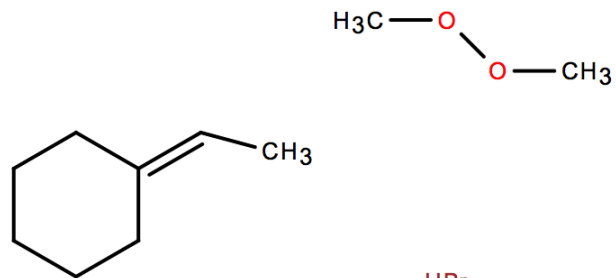




Live DEMO



Use The Smiles String Editor

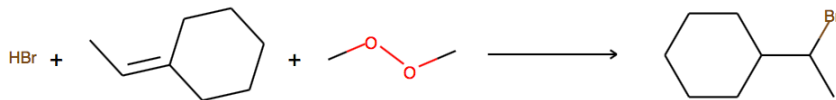
H
C
N
O
S
P
F
Cl
Br
I

Simply **draw reactants** & run the prediction

easy_reactions_20180813_11:59:18.856

Add to my molecules list

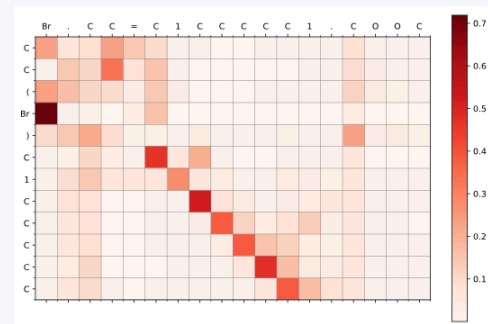
Add Reaction to Collection



STRING

Br.CC=C1CCCC1.COOC>>CC(Br)C1CCCC1

Attention Weight



Confidence: 1.00

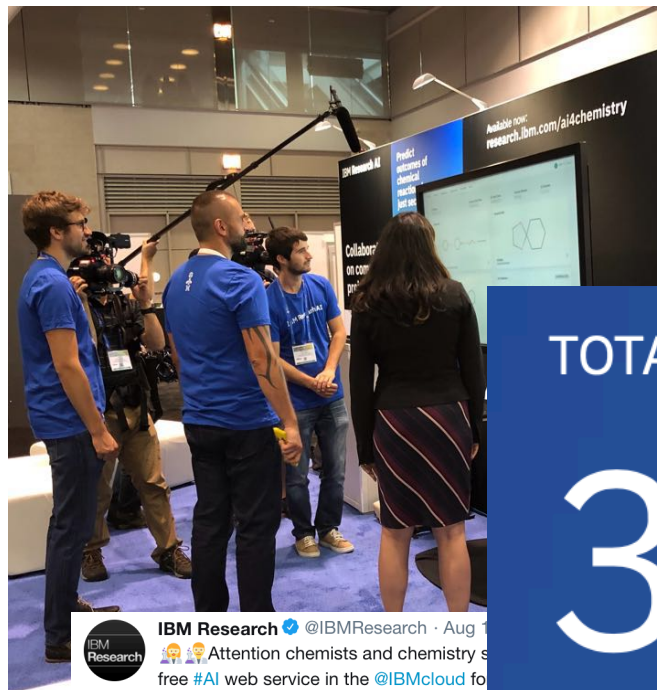
Help us improve! Send us your feedback.

What do you think about this result?

It's correct!

It's not so good!

Get back the **product**, the **attention** weights and the **confidence**



↑
452
↓

Posted by u/ibmzrl 2 days ago

IBM launches the first, FREE AI web service for predicting chemical reactions

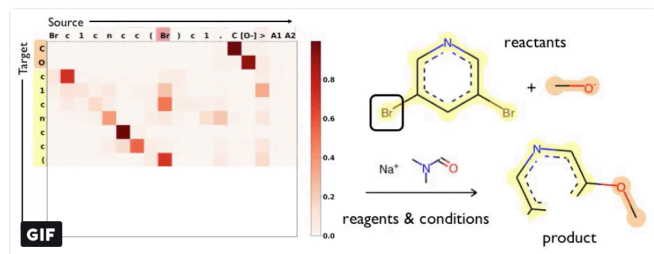
rxn.res.ibm.com/

79 Comments Share Save ...

TOTAL PREDICTED REACTIONS
37442



IBM Research @IBMRsearch · Aug 1
Attention chemists and chemistry students! Introducing a free #AI web service in the @IBMcloud for now rxn.res.ibm.com @AmerChemSociety #IBMBoston



21 656 942



Human vs AI challenge

Human Total Time
02:40:58.3

AI Total Time
00:02:36.0

Human Correct Answers
35

AI Correct Answers
64

Total Challenges
86

Your challenges results

Challenge

Chloro N-arylation

Human Total Time

00:02:30.2

AI Total Time

00:00:04.3

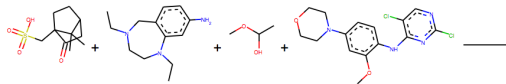
Human Answer

Correct

AI Answer

Correct

Reactants

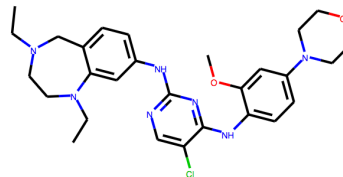


STRING

```
CC1(C)C2CCC1(CS(=O)(=O)O)C(=O)C2.CCN1CCN(CC)c2cc(N)ccc2C  
1.COC(C)O.COc1cc(N2CCOCC2)ccc1Nc1nc(Cl)ncc1Cl
```



Ground Truth



STRING

```
CCN1CCN(CC)c2cc(Nc3ncc(Cl)c(Nc4ccc(N5CCOCC5)cc4O)n3)ccc  
2C1
```



API access & open-source code

IBM RXN API - <https://rxn.res.ibm.com>

```
curl --data '{ "reactants" :  
  "C(O)1=C(O)C=C(C2=C(O)C=C([H])C([H])=C2C2=C([H])C=C([H])C=C2)C=C1O.[H].[H].  
  [H]", "mol":"" }' --header "Content-Type: application/json" --header "  
  Authorization:apk-have-here-your-own-key" -X POST  
https://rxn.res.ibm.com/rxn/api/api/v1/predictions/pr?projectId=5c532f56d6cb7  
600019ea342
```

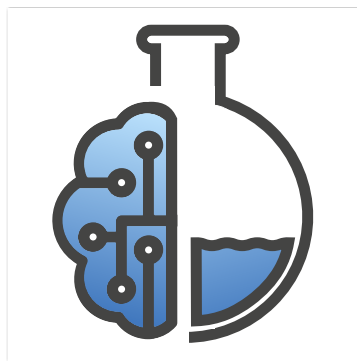
Code, trained models, data:

<https://github.com/pschwillr/MolecularTransformer>

Key Messages

- **SMILES** (and other line notations) might be an alternative representation → **Natural Language Processing** methods
- The **Molecular Transformer** is **universally applicable** across existing reaction datasets
 - No distinction between reactants-reagents required
 - Predictions in few hundreds of ms
 - Current state of the art in data-driven reaction prediction
- Have a look at **IBM RXN for Chemistry (it's free!)**, give us feedback and help us improve

Thank you for your attention!



Philippe Schwaller

phs@zurich.ibm.com / Twitter: [@pschwillr](https://twitter.com/pschwillr)

- “Found in Translation”: predicting outcomes of complex organic chemistry reactions using neural sequence-to-sequence models

(<https://doi.org/10.1039/C8SC02339E>)

- Molecular Transformer for Chemical Reaction Prediction and Uncertainty Estimation (<https://dx.doi.org/10.26434/chemrxiv.7297379>)

IBM **Research** AI



UNIVERSITY OF
CAMBRIDGE



THE WINTON PROGRAMME FOR THE

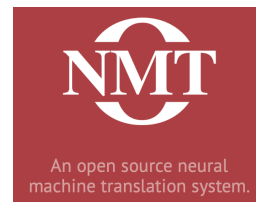
Physics of Sustainability



docker



Open-Source Cheminformatics
and Machine Learning



 PyTorch