09.05.2022



Machine Learning for Synthesis Planning









Grew up in Fribourg, Switzerland

- French

nature

- Swiss German / German



nanotechnology



Materials Science & Engineering

Virtual screening & simulation workflows Prof Nicola Marzari

BSc ('14)/MSc ('16)



Lab work on ternary polymer blends for organic solar cells Prof Frank Nüesch

NCCR

Catalysis

TT Assistant Prof

since Feb 2022

in Digital Chemistry



твм

Fribourg • UniBe

EPFL

Collaboration with synthetic chemists

^b UNIVERSITÄT BERN

h

[],

PhD in Chemistry and Molecular Sciences ('21) Prof Jean-Louis Reymond

MPhil in Physics ('19) Dr Alpha Lee







Machine learning for chemical synthesis Intern/PhD/Postdoc Dr Teodoro Laino

RoboRXN

IBM **Research** Europe

Computational quest for 2D materials

Outline



CC(C)S.CN(C)C=O.Fc1cccnc1 F.O=C([O-])[O-].[K+].[K+] >>CC(C)Sc1ncccc1F











Data Representation Models Real-world applications

Explore Chemical Reaction Space What chemistry do the models learn?





What molecule to make?

Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules

Rafael Gómez-Bombarelli‡e (b). Jennifer N. Wel‡e (b). David Duvenaud^e#, José Miguel Hernández-Lobato⁵#, Benjamín Sánchez-Lengeling[‡], Dennis Sheberla† (b), Jorge Aguilera-Iparraguirre[†], Timothy D. Hirzel[†], Ryan P. Adams[®], and Alán Aspuru-Guzik^{‡ ±}. (b)

[Submitted on 5 Jan 2017]

Generating Focussed Molecule Libraries for Drug Discovery with Recurrent Neural Networks

Marwin H.S. Segler, Thierry Kogej, Christian Tyrchan, Mark P. Waller



Design Make Test

How to make it?

Reaction prediction

Synthesis planning

Experimental validation





Machine Intelligence for Chemical Reaction Space. Schwaller et al. 2022 https://wires.onlinelibrary.wiley.com/doi/full/10.1002/wcms.1604



US20030166932A1: General Procedure H A solution of trifluoromethanesulfonic acid 3,5,8,8-tetramethyl-7,8dihydronaphthalen-2-yl ester (Compound 35, 0.41 g, 1.2 mmol), Pd(OAc)2 (0.027 g, 0.12 mmol), BINAP (0.11 g, 0.18 mmol), Cs2CO3 (0.56 g, 1.72 mmol), ethyl 4aminobenzoate (0.25 g, 1.5 mmol) and 5 mL of toluene was flushed with argon for 10 min, then stirred at 100° C, in a sealed tube for 48 h. After the reaction mixture had been cooled to room temperature, the solvent was removed, and the residue was purified by flash column (hexane:ethyl acetate=4:1) to give 0.34 g (80%) of the title compound as a yellowish solid.

Patents (broad, accessible)



Experiments ELN/HTE (narrow)



(e.g. ORD, Kearnes et al.)





Simulations (narrow)

EPFL Chemical reaction data

US Patents Text-mining (Lowe 2012/17)

Millions of reactions

```
BrC(Br)(Br)Br.CC...>>...

CO.Nc1cccc([N+]...>>...

CC(=O)O[BH-]...>>...

(OC(C)=O)OC(C)=O..>>...

...

precursors>>products

Benchmark sets

USPTO_MIT

USPTO_STEREO
```

Reaction SMILES

CC(C)S.CN(C)C=O.Fc1cccnc1F.O=C([O-])[O-].[K+].[K+]>>CC(C)Sc1ncccc1F



EPFL Reaction representations



Reaction SMILES (text-based reaction representation, precursors>>products)

Atom-mapping (e.g. RXNMapper)

Atom-mapped reaction (required for reaction template, centre and bond change extraction)



 $\begin{array}{l} CC(=0)[0-].CC(=0)[0-].Cc1ccccc1.0=C([0-])[0-].O=S(=0)(0[c:11]1[cH:12][c:13]2[c:14]([cH:15][c:16]1[CH3:17])[C:18]([CH3:19])=[CH:20]\\ [CH2:21][C:22]2([CH3:23])[CH3:24])C(F)(F)F,[CH3:1][CH2:2][0:3][C:4](=[0:5])[c:6]1[cH:7][cH:8][c:9]([NH2:10])[cH:25][cH:26]1.[Cs+].\\ [Pd+2].c1ccc(P(c2cccc2)c2ccc3ccccc3c2-c2c(P(c3cccc3)c3ccccc3)ccc3cccc23)cc1>>[CH3:1][CH2:2][0:3][C:4](=[0:5])[c:6]1[cH:7][cH:8]\\ [c:9]([NH:10][c:11]2[cH:12][c:13]3[c:14]([cH:15][c:16]2[CH3:17])[C:18]([CH3:19])=[CH:20][CH2:21][C:22]3([CH3:23])[CH3:24])[cH2:2][cH2:2][cH2:2]]\\ \end{array}$

Reaction template

Condensed Graph of Reaction



Chemical reactions from US patents (1976-Sep2016)

Dataset posted on 13.06.2017, 18:49 by Daniel Lowe

Reactions extracted by text-mining from United States patents published between 1976 and September 2016. The reactions are available as CML or reaction SMILES. Note that the reactions SMILES are derived from the CML. The files can be unzipped using a program like 7-Zip.



8727

downloads

12748

views

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

https://figshare.com/articles/Chemical_reactions_from_US_patents_1976-Sep2016_/5104873

EPFL Chemical reaction data

US Patents



Reaction SMILES

Millions of reactions

BrC(Br)(Br)Br.CC...>>... CO.Nc1cccc([N+]...>>... CC(=O)O[BH-]...>>... (OC(C)=O)OC(C)=O..>>... ... precursors>>products Benchmark sets USPTO_MIT USPTO_STEREO

CC(C)S.CN(C)C=O.Fc1cccnc1F.O=C([O-])[O-].[K+].[K+]>>CC(C)Sc1ncccc1F





SMILES, a Chemical Language and Information System. 1. Introduction to Methodology and Encoding Rules

DAVID WEININGER Medicinal Chemistry Project, Pomona College, Claremont, California 91711

Received June 17, 1987



Krenn, Mario, et al. "SELFIES and the future of molecular string representations." *arXiv preprint arXiv:2204.00056* (2022).

Action and the second of the s



Split -> "sequences of atoms" called tokens

CC(C)S.CN(C)C = O.Fc1cccnc1F.O = C([O-])[O-].[K+].[K+] >> CC(C)Sc1ncccc1F

\rightarrow Borrow methods developed for human languages

Nam & Kim, arXiv:1612.09529; Liu et al., ACS Centr. Sci. 2017; Schwaller et al., Chem. Sci, 2018

EPFL Sequence-2-sequence models

French: Le chat est noir.

German: Die Katze ist schwarz.



Sutskever et al., Sequence to Sequence Learning with Neural Networks. NeurIPS, 2014.

Cho et al., Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. EMNLP, 2014.

Sequence-2-sequence models with attention

One state per input



Attention = ability to focus on most important features

Bahdanau et al., Neural Machine Translation by Jointly Learning to Align and Translate. ICLR, 2015 Luong et al., Effective approaches to attention-based neural machine translation, EMNLP, 2015

EPFL

Transformer architecture



- Stacks of attention layers
- Multi-head attention



- No rules integrated / no chemical knowledge
- Accurate predictions on unseen reactions
- Better than rule and graph-based approaches

Schwaller et al., Molecular Transformer – A Model for Uncertainty-Calibrated Chemical Reaction Prediction. ACS Central Science, 2019





USPTO-MIT benchmark (no stereochemistry)



Separated vs mixed setting



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**

[HTML] A graph-convolutional neural network model for the prediction of chemical reactivity

<u>CW Coley</u>, <u>W Jin</u>, <u>L Rogers</u>, <u>TF Jamison</u>... - Chemical ..., 2019 - pubs.rsc.org We present a supervised learning approach to predict the products of organic reactions given their reactants, reagents, and solvent (s). The prediction task is factored into two stages comparable to manual expert approaches: considering possible sites of reactivity ...

☆ 55 Cited by 132 Related articles All 8 versions

Graph-edit-based, atom-mapping dependent

Methods		Top- n acc	uracy (%)	
-	1	3	5	10
USPTO_480	k_mixed			
MEGAN (Sacha et al., 2021)	86.3	92.4	94.0	95.4
Molecular Transformer (Schwaller et al., 2019)	88.6	93.5	94.2	94.9
Graph2SMILES (D-GCN) (ours)	90.3	94.0	94.6	95.2
Graph2SMILES (D-GAT) (ours)	90.3	94.0	94.8	95.3
Augmented Transformer (Tetko et al., 2020)	90.6	-	96.1	-
Chemformer (Irwin et al., 2021)	91.3	-	93.7	94.0

Chemformer: a pre-trained transformer for computational chemistry

Ross Irwin¹, Spyridon Dimitriadis^{1,2}, Jiazhen He¹ and Esben Jannik Bjerrum^{3,1} ib

State-of-the-art augmented NLP transformer models for direct and single-step retrosynthesis

Igor V. Tetko ⊠, Pavel Karpov, Ruud Van Deursen & Guillaume Godin ⊠

Augmented Transformer

PERMUTATION INVARIANT GRAPH-TO-SEQUENCE MODEL FOR TEMPLATE-FREE RETROSYNTHESIS AND REACTION PREDICTION Graph2SMILES

Zhengkai Tu 1,2 and Connor W. Coley 1,3

Molecule Edit Graph Attention Network: Modeling Chemical Reactions as Sequences of Graph Edits

Mikołaj Sacha, Mikołaj Błaż, Piotr Byrski, Paweł Dąbrowski-Turnański, Mikołaj Chromiński, Rafał Loska, Paweł Włodarczyk-Pruszyński, and Stanisław Jastrzębski*



Extensive data augmentations



 c1c(N)ccc(C)c1
 Cc1ccc(N)cc1

 c1cc(C)ccc1N
 c1(N)ccc(C)cc1

 c1c(C)ccc(N)c1
 Nc1ccc(C)cc1

 c1(C)ccc(N)cc1
 c1cc(N)ccc1C

Molecule SMILES randomizations

{aryl_halide}.{methylaniline}.{pd_catalyst}.{ligand}.{base}.{additive}>>{product}
{ligand}.{base}.{methylaniline}.{additive}.{pd_catalyst}.{aryl_halide}>>{product}
{base}.{methylaniline}.{pd_catalyst}.{aryl_halide}.{ligand}>>{product}
{additive}.{base}.{aryl_halide}.{ligand}.{methylaniline}.{pd_catalyst}>>{product}
{aryl_halide}.{pd_catalyst}.{base}.{ligand}.{methylaniline}.{additive}>>{product}
}

Molecule permutations

Bigger model (6 layers of size 512) compared to Molecular Transformer (4 layers of size 256) and 100x test-time augmentation -> much slower at inference time

State-of-the-art augmented NLP transformer models for direct and single-step retrosynthesis

Igor V. Tetko 🖂, Pavel Karpov, Ruud Van Deursen & Guillaume Godin 🖂

Large-scale pretraining



Chemformer: a pre-trained transformer for computational chemistry

Ross Irwin¹, Spyridon Dimitriadis^{1,2}, Jiazhen He¹ and Esben Jannik Bjerrum^{3,1}

Graph2SMILES -> Graph encoder with a SMILES decoder



PERMUTATION INVARIANT GRAPH-TO-SEQUENCE MODEL FOR TEMPLATE-FREE RETROSYNTHESIS AND REACTION PREDICTION

Zhengkai Tu^{1,2} and Connor W. Coley^{1,3}

EPFL Improving Few- and Zero-Shot Reaction Template Prediction Using Modern Hopfield Networks



Philipp Seidl, Philipp Renz, Natalia Dyubankova, Paulo Neves, Jonas Verhoeven, Jörg K. Wegner, Marwin Segler, Sepp Hochreiter, and Günter Klambauer*

-> Probably more on that on Wednesday

EPFL Stereochemistry & experimental validation

- 14-step synthesis of a lipid-linked oligosaccharide
- >40% accuracy increase with Carbo Transformer
- Similar performance gains on JACS/CARBO test sets



IBM Resear



Transfer learning enables the molecular transformer to predict regio- and u^* stereoselective reactions on carbohydrates

Pesciullesi*, Schwaller* et al., Nature Communications, 2020

EPFL What is transfer learning?



Transfer learning applicable to any reaction subspace of interest!

EPFL *Retrosynthesis* (Corey, Nobel prize, 1990)



Lego analogy: Amol Thakkar

AIDD Summer School / ML for Synthesis Planning Steps to target

EPFL

Single-step retrosynthesis using language models (2017)

Introduction of USPTO-50k benchmark dataset $i = \frac{1}{N}$

Figure 1. Phenylalanine synthetic scheme.

Retrosynthetic Reaction Prediction Using Neural Sequence-to-Sequence Models

Bowen Liu[†][ip], Bharath Ramsundar[‡][ip], Prasad Kawthekar[‡], Jade Shi[†], Joseph Gomes[†], Quang Luu Nguyen[†], Stephen Ho[†], Jack Sloane[†], Paul Wender^{†§}[ip], and Vijay Pande^{*†‡1}



EPFL

Single-step retrosynthesis (USPTO 50k)

- no reaction class information given

Model	Top-1	Top-5	Top-10
SMILES-based			
SCROP [<u>51]</u>	43.7	65.2	68.7
Two-way transformer [52]	47.1	73.1	76.3
Aug transformer [<u>5</u>]	48.3	73.4	77.4
Chemformer	53.6	61.1	61.7
Chemformer-Large	54.3	62.3	63.0

Multi-step synthesis planning



Single-step prediction approach: template-based / graph-based / SMILES-based

Search algorithm: MCTS / RL / A* / beam search

EPFL

Multi-step synthesis planning

Molecular Transformer Retro

IBM Research





Molecular Transformer Forward

Predicting retrosynthetic pathways using a combined linguistic model and hyper-graph exploration strategy

Philippe Schwaller, Riccardo Petraglia, Valerio Zullo, Vishnu H Nair, Rico Andreas Haeuselmann, Riccardo Pisoni, Costas Bekas, Anna Iuliano, Teodoro Laino *Chem. Sci.*, 2020, Advance Article <u>https://doi.org/10.1039/C9SC05704H</u>

EPFL



Our approach: Precursors prediction (reactants, reagents, catalysts, solvents, ..)



Reactants prediction (atom-mapping dependent)



Rule-based approaches:

- Segler et al. Nature, 555, 604–61
- Coley et al. Science 365 (6453)
- Genheden et al. J Cheminf. 12 (1), 1-9
- Thakkar et al. Chem. Sci. 11 (1), 154-168

and other SMILES-based approaches:

- Liu et al. ACS Cent. Sci. 3 (10), 1103-1113
- Tetko et al. Nature Comm. 11 (1), 1-11

EPFL Single step -> multi-step performance?



Template Library Size

Top-1 accuracy (single step) Multi-step performance

New metrics/benchmarks are required!! (collaborations?)

Datasets and their influence on the development of computer assisted synthesis planning tools in the pharmaceutical domain $\underline{}^{\pm}$

Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy $\!\!\!^{+}_{-}$

Philippe Schwaller 🕲 *ª, Riccardo Petraglia ª, Valerio Zullo º, Vishnu H. Nair ª, Rico Andreas Haeuselmann ª, Riccardo Pisoni ª, Costas Bekas ª, Anna Iuliano 💿 ^b and Teodoro Laino ^a

Model-based metrics



Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy $\!\!\!\!^{\pm}$

Philippe Schwaller 💿 **, Riccardo Petraglia *, Valerio Zullo ^b, Vishnu H. Nair *, Rico Andreas Haeuselmann *, Riccardo Pisoni *, Costas Bekas *, Anna Iuliano 🗓 ^b and Teodoro Laino *




How to *facilitate adoption*?

IBM **RXN for Chemistry Free platform:** rxn.res.ibm.com and API access.





Get back the product and a confidence score

EPFL *Retrosynthesis* interface



EPFL

Retrosynthesis user feedback



Alessandro Castrogiovanni

Sparr Research Group University of Basel, 2019

42

What's next?

IBM Research

What's next? #RoboRXN **IBM** Research

RXN

or Chemist





EPFL



Vaucher, A.C., Schwaller, P., Geluykens, J. *et al.* Inferring experimental procedures from text-based representations of chemical reactions. *Nat Commun* **12**, 2573 (2021). Vaucher, A.C., Zipoli, F., Geluykens, J. et al. Automated extraction of chemical synthesis actions from experimental procedures. Nat Commun **11**, 3601 (2020).

Alain Vaucher

Reaction Transformer models

From a sequence to a sequence.



encoder

input

decoder

output

From a sequence to a single value/label.



Chemical reactions can be represented as text.



Reaction classification

Schwaller et al., Nature Mach. Int., 2021

Yields predictions

Schwaller et al., MLST, 2021 Schwaller et al., chemrxiv.13286741

encoder-only



Self-supervised training

Pretraining



Masked Language Modelling

- Predict mask given context
- Unlimited training data

EPFL Reaction *classification*

Reaction SMILES Reaction class Fischer-Speier Classification model >98% accurate **Esterification** CO.O=C(O)c1ccc([N+](=O)[O-])cc1F.O=S(=O)2.6.3 (0)0>>COC(=0)c1ccc([N+](=0)[0-])cc1F Ketone reductive amination CCOC(=0)C(C)=0.Nc1cc(Cl)cc(Cl)c1>>CCOC(=0)1.2.5 C(C)Nc1cc(Cl)cc(Cl)c1 **Fingerprints Reaction encoder** = Encoded reaction properties °C0.0=C(0)c1ccc([N+](=0)[0-])cc1F.0=S(=0) (0)0>>COC(=0)c1ccc([N+](=0)[0-])cc1F [0.14, 0.25, ..., 0.9] **[0.22, 0.83, ..., -0.12]** Mapping the space of chemical reactions using nature attention-based neural networks CCOC(=0)C(C)=0.Nc1cc(Cl)cc(Cl)c1>>CCOC(=0)machine intelligence Philippe Schwaller 1227, Daniel Probst², Alain C. Vaucher 1, Vishnu H. Nair¹, David Kreutter², C(C)Nc1cc(Cl)cc(Cl)c1

Teodoro Laino¹ and Jean-Louis Reymond²



Chemical Reaction Atlas

Powered by TMAP



Unrecognised
 Heteroatom alkylation and arylation
 Acylation and related processes
 C-C bond formation
 Heterocycle formation
 Protections
 Deprotections
 Reductions
 Oxidations
 Functional group interconvertions
 Functional group additions
 Resolutions

EPFL Chemical reaction space exploration



EPFL Search for *similar reactions*

IBM RXN								Visu	ual ec	litor	Sm	Smiles string editor									Run forward reaction prediction			\rightarrow		⊗
		\square		ß	ଯ	X			¢	°,	100%	\$	<u>*</u> .	χ.	٨	$\langle \bar{A} \rangle$	V	۲	Ь	.).	J.		¢	0	0	
R																										н
0																										С
/																										Ν
\sim																										0
A°																										S
A																										Р
0																										F
[]																										CI
Ø																	k									Br
→																										Т
→R1																										₽T
	0	$\hat{\Box}$	0	\bigcirc	\bigtriangleup		0	\bigcirc		۵														Chira	1	

 \rightarrow Access to reaction conditions, procedures, articles

EPFL Open-source code

https://github.com/rxn4chemistry/rxnfp | pip install rxnfp

machine intelligence

MATTERS ARISING

https://doi.org/10.1038/s42256-021-00367-2



Reusability report: Learning the language of synthetic methods used in medicinal chemistry

Jon Paul Janet 💿, Anna Tomberg and Jonas Boström 💿 🖂

ARISING FROM P. Schwaller et al. Nat. Mach. Intell. https://doi.org/10.1038/s42256-020-00284-w (2020).

Medicinal Chemistry, Research and Early Development, Cardiovascular, Renal and Metabolism (CVRM), BioPharmaceuticals R&D, AstraZeneca, Gothenburg, Sweden. [™]e-mail: jonas.bostrom@astrazeneca.com

EPFL Prediction of *patent reaction yields*



Long story short \rightarrow I didn't work

EPFL Prediction of *HTE reaction yields*

3955 Buchwald-Hartwig reactions with measured yield



Yield prediction using DFT descriptors + random forest (RF)

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).

Multi fingerprint features (MFF) + RF

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

Learning to predict reaction yields

 $\label{eq:started_st$



Reaction Transformer

77.6 %

• output

input encoder-onlv 20.0 %

56.3 %

HILE data from: 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).

Results: 70/30 random split



- [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 crosscoupling 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).

Results: *reduced* training data





2] 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).

But why do *Transformers* work so well?

Self-supervised training

millions of **unlabeled** reaction SMILES

BrC(Br)(Br)Br.CC...>>... CO.Nc1cccc([N+]...>>... ... (OC(C)=O)OC(C)=O..>>... precursors>>products

Lowe, 2017



Work with Ben Hoover and Hendrik Strobelt





u^b IBM Research

Self-supervised training



Masked Language Modelling:

- Reaction corruption task
- Unlimited training data

But why do *Transformers* work so well?

Ben Hoover and Self-supervised training Visual inspection Hendrik Strobelt millions of unlabeled reaction SMILES BrC(Br)(Br)Br.CC...>>... CO.Nc1cccc([N+]...>>... ... (OC(C)=0)0C(C)=0..>>... precursors>>products input encoder-only Lowe, 2017 Learned patterns $u^{\scriptscriptstyle b}$

Work with

IBM Research

UNIVERSITÄT



Input: precursors >> product

EPFL Visual inspection of attention weights



[CH3:1][OH:2].O[C:3](=[O:4])[c:5]1[cH:6][cH:7][cH:8][cH:9][cH:10]1>>[CH3:1][O: 2][C:3](=[O:4])[c:5]1[cH:6][cH:7][cH:8][c H:9][cH:10]1



RXNMapper



Discovery: Atom-mapping

CO.O=C(O)c1ccccc1>>COC(=O)c1ccccc1

Independent *benchmark*

- 1851 curated reactions
- 5 popular tools (RXNMapper, ChemAxon, NameRXN, Indigo, and RDTools)

"**RXNMapper was ranked best** with an accuracy of **83.74%** by being the second fastest algorithm (0.05 seconds)."

Compared to **38.47%** accuracy for Indigo (open-source alternative)

Atom-to-atom Mapping: A Benchmarking Study of Popular Mapping Algorithms and Consensus Strategies

Arkadii Lin, Natalia Dyubankova, Timur I. Madzhidov, Ramil I. Nugmanov, Jonas Verhoeven, Timur R. Gimadiev, Valentina A. Afonina, Zarina Ibragimova, Assima Rakhimbekova ... See all authors 🗸

First published: 02 November 2021 | https://doi.org/10.1002/minf.202100138



Grammar of chemical reactions

EPFL Atom-mapping in *AI-assisted synthesis planning*

Template-based approaches Planning chemical syntheses with deep neural networks and symbolic AI

Marwin H. S. Segler^{1,2}, Mike Preuss³ & Mark P. Waller⁴

A robotic platform for flow synthesis of organic compounds informed by AI planning

😰 Connor W. Coley^{1,*}, 😳 Dale A. Thomas III^{1,2,*,†}, Justin A. M. Lummiss^{3,*,†}, 😳 Jonathan N. Jaworski^{3,‡}, 😳 Christopher P....

Datasets and their influence on the development of computer assisted synthesis planning tools in the pharmaceutical domain†

Amol Thakkar, 💿 *ab Thierry Kogej,^a Jean-Louis Reymond, 💿 ^b Ola Engkvist^a and Esben Jannik Bjerrum 💿 *a

Graph NN-based approaches (bond changes / graph edits)

A graph-convolutional neural network model for the prediction of chemical reactivity $\!\!\!\!\!\!\!^{\dagger}$

Connor W. Coley 😳², Wengong Jin ⁵, Luke Rogers ², Timothy F. Jamison 🙆², Tommi S. Jaakkola ⁵, William H. Green ⁶, Regina Barzilay ^{+ 5} and Klavs F. Jensen ⁶

Learning Graph Models for Retrosynthesis Prediction VR Somnath, C Bunne, CW Coley, A Krause, R Barzilay Neural Information Processing Systems (NeurIPS) 35

SMILES-to-SMILES approaches

Molecular Transformer: A Model for Uncertainty-Calibrated Chemical Reaction Prediction

Philippe Schwaller*, Teodoro Laino, Théophile Gaudin, Peter Bolgar, Christopher A. Hunter, Costas Bekas, and Alpha A. Lee*

Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy $\!\!\!^{\dagger}$

Philippe Schwaller ()*, Riccardo Petraglia ^a, Valerio Zullo ^b, Vishnu H. Nair ^a, Rico Andreas Haeuselmann ^a, Riccardo Pisoni ^a, Costas Bekas ^a, Anna Iuliano ()^b and Teodoro Laino ^a

Atom-mapping dependent -> training data

atom-maps -> templates atom-maps -> graph edits

All benefit from better atom-maps!

Atom-mapping independent

-> Atom rearrangement not tracked

EPFL

Open-source code

https://github.com/rxn4chemistry/rxnmapper | pip install rxnmapper

```
from rxnmapper import RXNMapper
rxn_mapper = RXNMapper()
rxns = ['CC(C)S.CN(C)C=0.Fc1cccnc1F.0=C([0-])[0-].[K+].[K+]>>CC(C)Sc1ncccc1F', 'C1C0CC
results = rxn_mapper.get_attention_guided_atom_maps(rxns)
```

```
[{'mapped_rxn': 'CN(C)C=0.F[c:5]1[n:6][cH:7][cH:8][cH:9][c:10]1[F:11].0=C([0-])[0-].[CH3:1][CH
    'confidence': 0.9565619900376546},
    {'mapped_rxn': 'C1C0CC01.CC(C)(C)[0:3][C:2](=[0:1])[CH2:4][0:5][NH:6][C:7](=[0:8])[NH:9][CH2:
    'confidence': 0.9704424331552834}]
```

EPFL

Demo on RXNMapper.ai

Mapping:



Model sizes – opposite of trend





Reaction prediction: Schwaller et al. *ACS Cent. Sci.* 2019 Pesciullesi et al. *Nature Comm.* 11, 4874 (2020)



Synthesis planning: Schwaller et a *Chem. Sci., 2020,11, 3316-3325* **Classification/fingerprints:** Schwaller et al. *Nature Mach. Intell., 2021*

Reaction yields: Schwaller et al. *Mach. Learn.: Sci. Technol.*, 2021; Schwaller et al. NeurIPS 2020 Ml4Mol workshop

Atom-mapping: Schwaller et al. *Science Advances, 2021*





EPFL Conclusions

- Apply ML/NLP methods to synthesis
- **RXN for Chemistry** platform / API access
 - reaction prediction / synthesis planning
 - paragraph-2-actions / RoboRXN simulator
- Data-driven reaction fingerprints
 - Reaction atlas / nearest neighbour search
 - Reaction yields
- Transformers capture the grammar of chemical reactions
- Fast, high-quality atom-maps even on strongly imbalanced reactions.









rxnmapper.ai
EPFL Review & Outlook



Machine intelligence for chemical reaction space

https://wires.onlinelibrary.wiley.com/ Co doi/full/10.1002/wcms.1604

Philippe Schwaller 🔀, Alain C. Vaucher, Ruben Laplaza, Charlotte Bunne, Andreas Krause, Clemence Corminboeuf, Teodoro Laino

First published: 07 March 2022 | https://doi.org/10.1002/wcms.1604



IBM Research AI

7 | [

^b UNIVERSITÄT BERN





Bojana Rankovic, Oliver Schilter (with IBM)

Teodoro Laino, Alain Vaucher, Alessandra Toniato, Theophile Gaudin, Costas Bekas, Daniel Probst, Matteo Manica and the RoboRXN team, Ben Hoover, Hendrik Strobelt

Jean-Louis Reymond, Giorgio Pesciullesi, David Kreutter, Alice Capecchi, Amol Thakkar, and the Reymond group

Alpha Lee, Peter Bolgar and Ryan-Rhys Griffiths

Clemence Corminboeuf, Andreas Krause, Ruben Laplaza, Charlotte Bunne







O PyTorch

Transformers

EPFL Many thanks for your attention!

@pschwllr @SchwallerGroup



https://github.com/pschwllr https://github.com/rxn4chemistry



phs@zurich.ibm.com- philippe.schwaller@epfl.ch

EPFL Many thanks for your attention!

@pschwllr @SchwallerGroup



https://github.com/pschwllr https://github.com/rxn4chemistry



phs@zurich.ibm.com- philippe.schwaller@epfl.ch



Colab exercise: https://github.com/schwallergroup/dmd s_language_models_for_reactions