

Reagent prediction with a transformer and its benefits for reaction product prediction

Mikhail Andronov, The Swiss AI Lab IDSIA

*Advanced Machine Learning for Innovative Drug Discovery (AIDD) –
Horizon 2020 Marie Skłodowska-Curie Innovative Training Network*

Cite this: DOI: 00.0000/xxxxxxxxxx

Reagent Prediction with a Molecular Transformer Improves Reaction Data Quality[†]

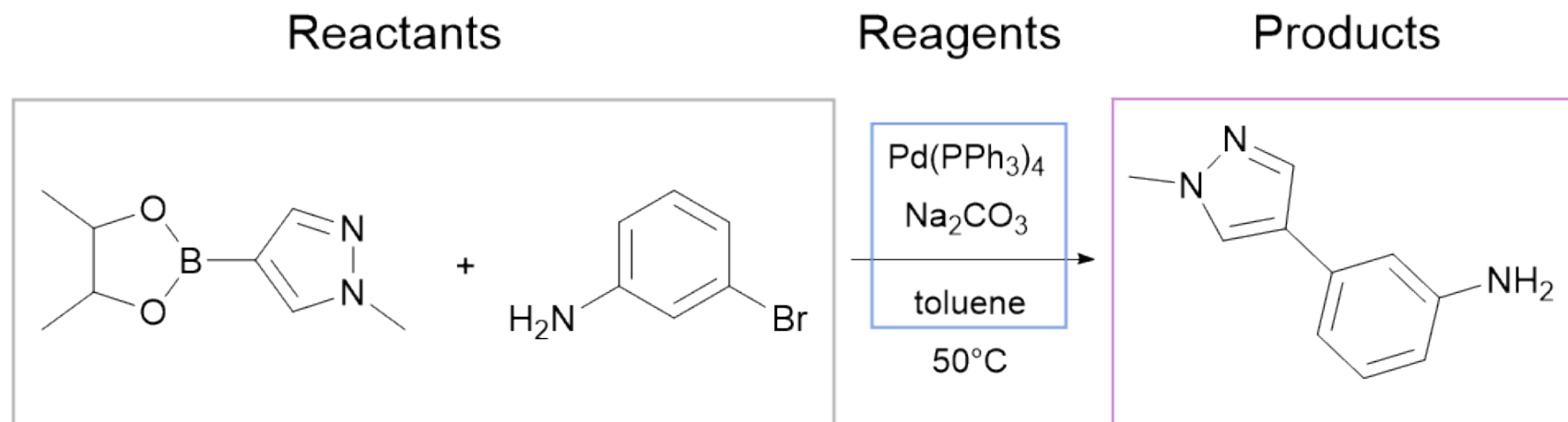
Mikhail Andronov,^{a,e} Varvara Voinarovska,^b Natalia Andronova,^c Michael Wand,^{a,d} Djork-Arné Clevert,^e and Jürgen Schmidhuber^{a,f}

Presented at 23rd EuroQSAR (Heidelberg, Germany)

Uploaded to *chemRxiv* in November

Under review to *Chemical Science*

Chemical reactions



A reaction type is defined by the reaction center and reagents.

With different reagents, reactants can turn into different products.

Any part of a reaction can be predicted.

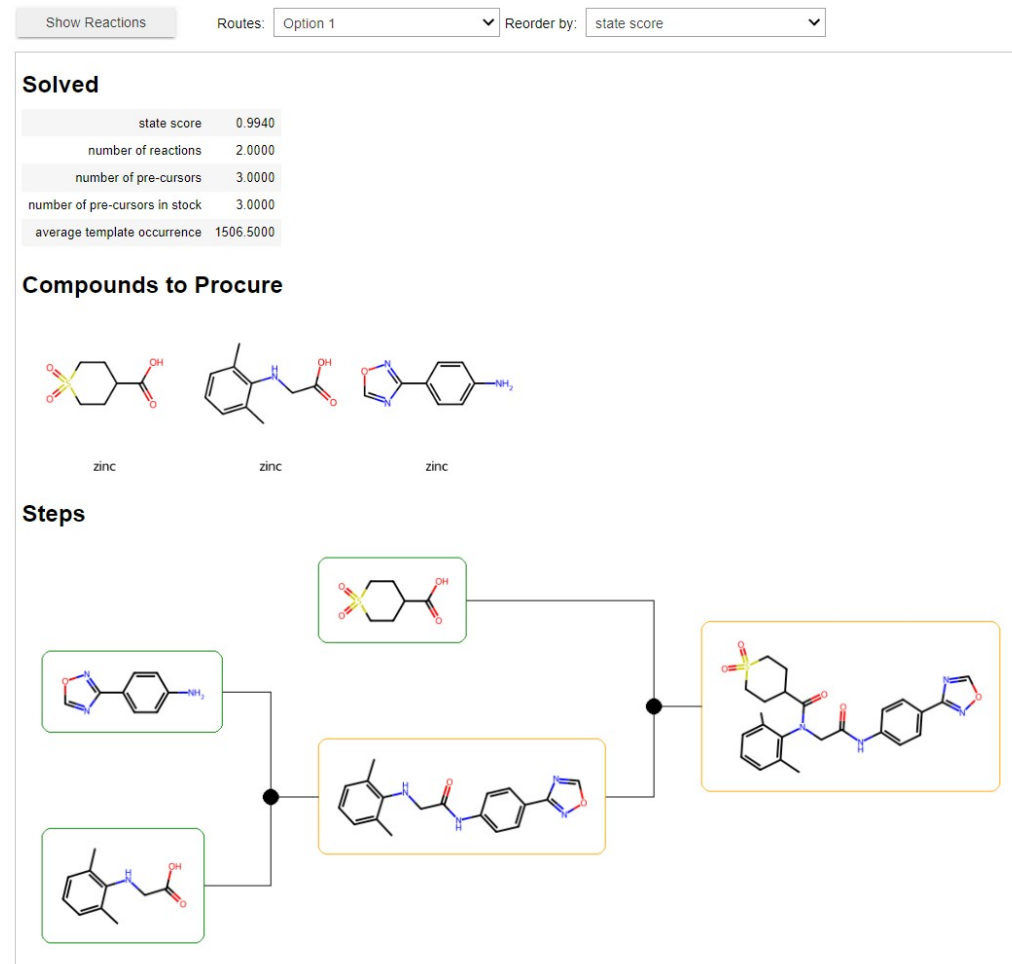
Why predict reagents?

1). To help CASP

Aizynthfinder generates routes without reagents.

2). To address data flaws

Many reactions in USPTO don't always have well-specified reagents



Literature: conditions prediction

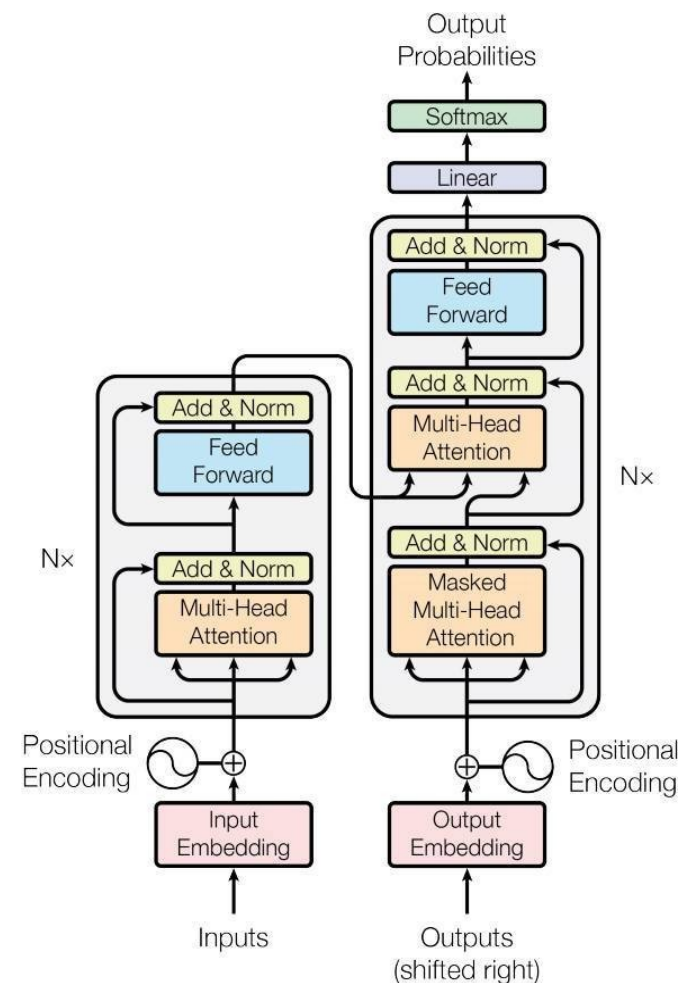
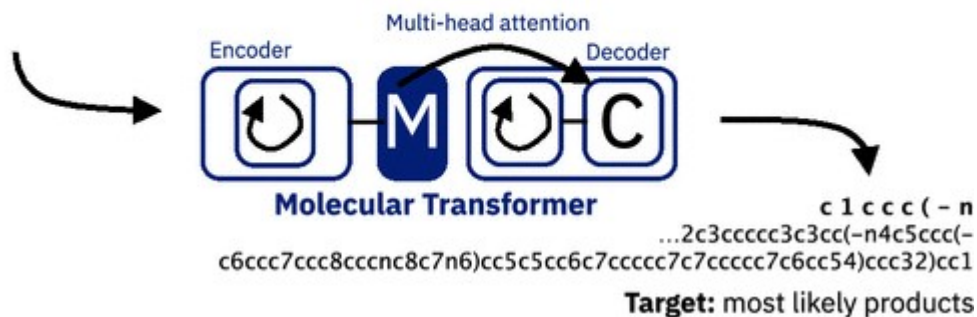
Paper	Reaction types	Goal of predictions	Dataset	Model	Data format
<i>Walker et al. 2019</i>	Five name reactions	Solvent	Reaxys	SVM	Molecular fingerprints (OpenBabel)
<i>Afonina et al. 2021</i>	Hydrogenation reactions	Catalyst, temperature, pressure	Reaxys	MLP	Molecular fingerprints (ISIDA Fragmentor 2017)
<i>Gao et al. 2018</i>	Broad range of reactions	Restricted set of reagents, temperature	Reaxys	MLP	Molecular fingerprints (RDKit)
<i>Maser et al. 2021</i>	Four name reactions	Restricted set of reagents, temperature	Reaxys	GBM, GNN	Molecular graphs

Transformer

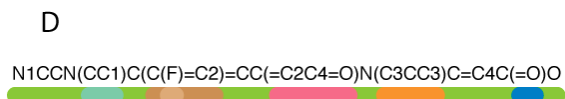
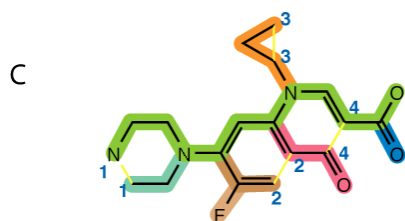
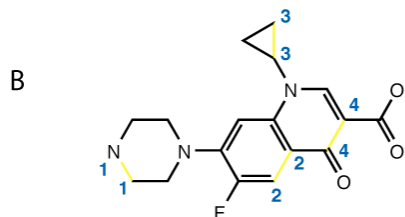
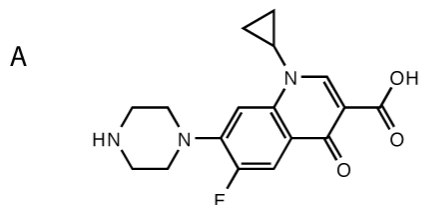
Today the standard base model for all kinds of NLP tasks.
Originally proposed for machine translation.

Input: reactants-reagents (atom-wise tokenization)

Br c 1 c c c 2 ...c(c1)c1cc3c4cccc4c4cccc4c3cc1n2-c1ccc2c(c1)c1cccc1n2-c1cccc1.CCO.
Cc1cccc1.OB(O)c1ccc2ccc3ccnc3c2n1.c1ccc([PH])(c2cccc2)(c2cccc2)[Pd]([PH])(c2cccc2)
(c2cccc2)c2cccc2)([PH])(c2cccc2)(c2cccc2)c2cccc2)[PH](c2cccc2)(c2cccc2)c2cccc2)cc1



SMILES



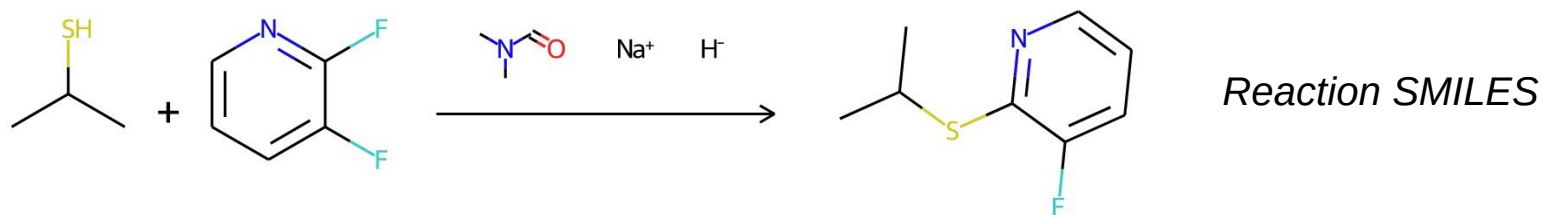
SMILES for ciprofloxacin

SMILES – a text notation of organic molecules designed for chemical information systems.

Reaction SMILES are to depict reactions.

The idea of SMILES is to build a spanning tree in the molecular graph.

CC(C)S.Fc1ccncc1F>CN(C)C=O.[Na+].[H-]>CC(C)Sc1ncccc1F

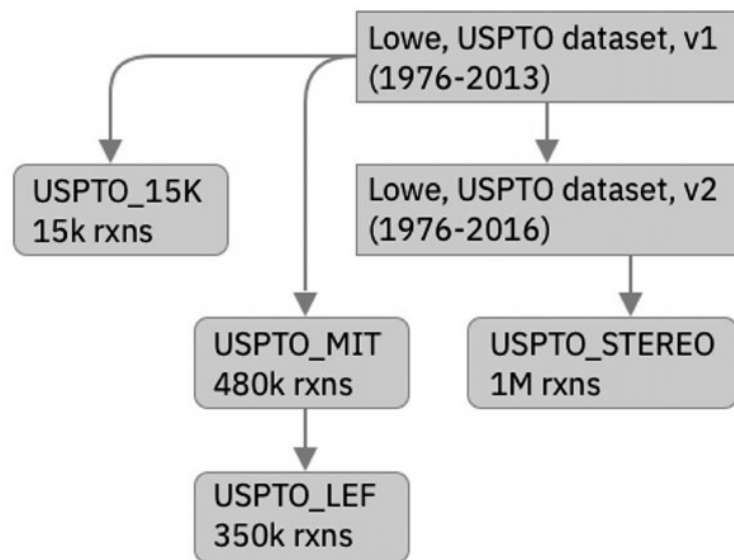


Any part of a reaction SMILES can be predicted by masked language modeling

Chemical reaction data

Chemical reactions from US patents (USPTO dataset, 2012) – the only open chemical reaction dataset.

Consists of 1-2M reactions obtained by text mining, pretty noisy.



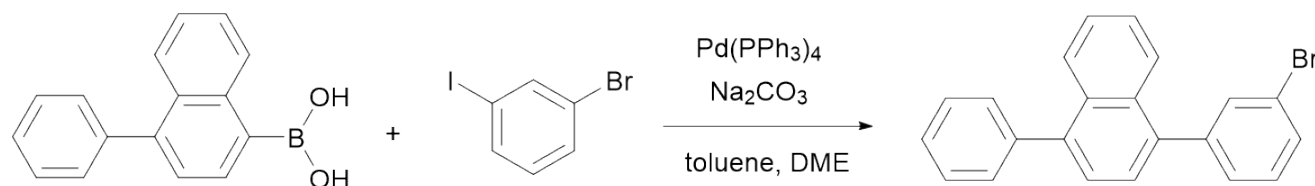
Reaxys – a proprietary expert-curated database from Elsevier, 56M reactions.

The screenshot shows the Reaxys interface with search results for a chemical reaction. The reaction is: CC1=CC=C(C=C1)C2=CC=CC=C2O1C=CC=CC=C1 + CC1=CC=C(C=C1)C2=CC=CC=C2O1C=CC=CC=C1 → CC1=CC=C(C=C1)C2=CC=CC=C2O1C=CC=CC=C1. The reaction ID is 25769146. The interface includes filters on the left, a list of results in the center, and a table of conditions and references on the right.

Conditions	Yield	Reference
With triethylamine in ethyl acetate at 75 - 80°C; for 18h; Reagent/catalyst; Solvent; Experimental Procedure	95.1%	Zhejiang Huahai Pharmaceutical Co., Ltd.; Du, Xiaoqi; Zhou, Lianchao; Liu, Jiegen CN106187857, 2016, A Location in patent: Paragraph 0014; 0026; 0027; 0028; 0029; 0030-0033 Full Text Details Abstract
With sodium acetate; acetic acid at 80°C; for 5h; Reagent/catalyst; Temperature; Experimental Procedure	91.3%	Guangzhou Aige Biological Technology Co., Ltd.; Tan Bin; Zhang Xiantao; He Shengjiang CN107188842, 2017, A Location in patent: Paragraph 0037-0044 Full Text Details Abstract
Stage #1: (S)-1-(3-ethoxy-4-methoxyphenyl)-2-(methylsulfonyl)ethanamine-(S)-2-acetamido-4-methylpentanoate With sodium hydroxide in dichloromethane at 0 - 5°C; for 2h; Stage #2: 3-acetylaminothalic anhydride With perchloric acid; acetic acid in dichloromethane at 45°C; for 3.33333h; Reflux; Experimental Procedure	89.2%	Xinfa Pharmaceutical Industry Limited Company; Qi, Yuxin; Chen, Jun; Zhou, Lishan; Fan, Yansen; Ju, Lizhu; Li, Xinfu CN105348172, 2016, A Location in patent: Paragraph 0070; 0071; 0072 Full Text Details Abstract

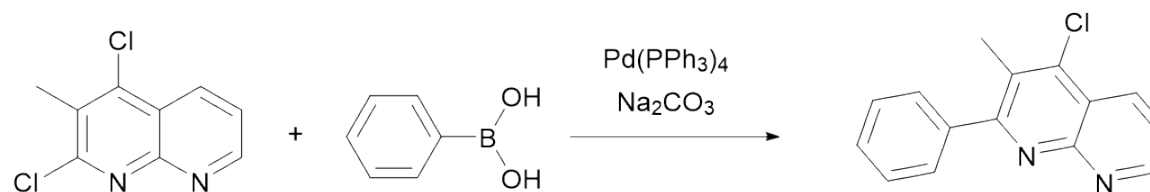
USPTO noise

US07985491B2, 2011



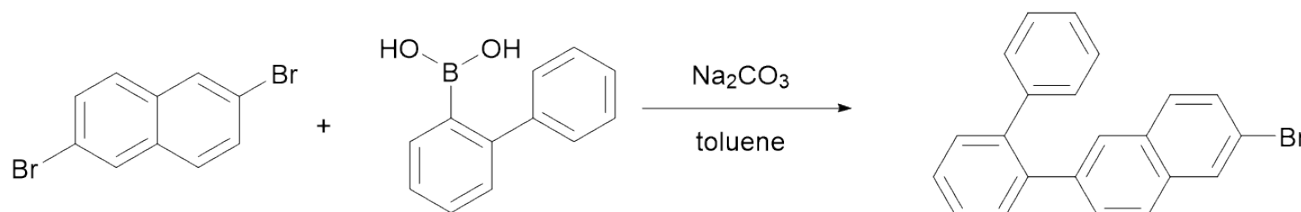
All reagents are specified

US08765940B2, 2014



No solvent is specified

US08853675B2, 2014



The Pd catalyst is missing

US09394290B2, 2016



All reagents are missing

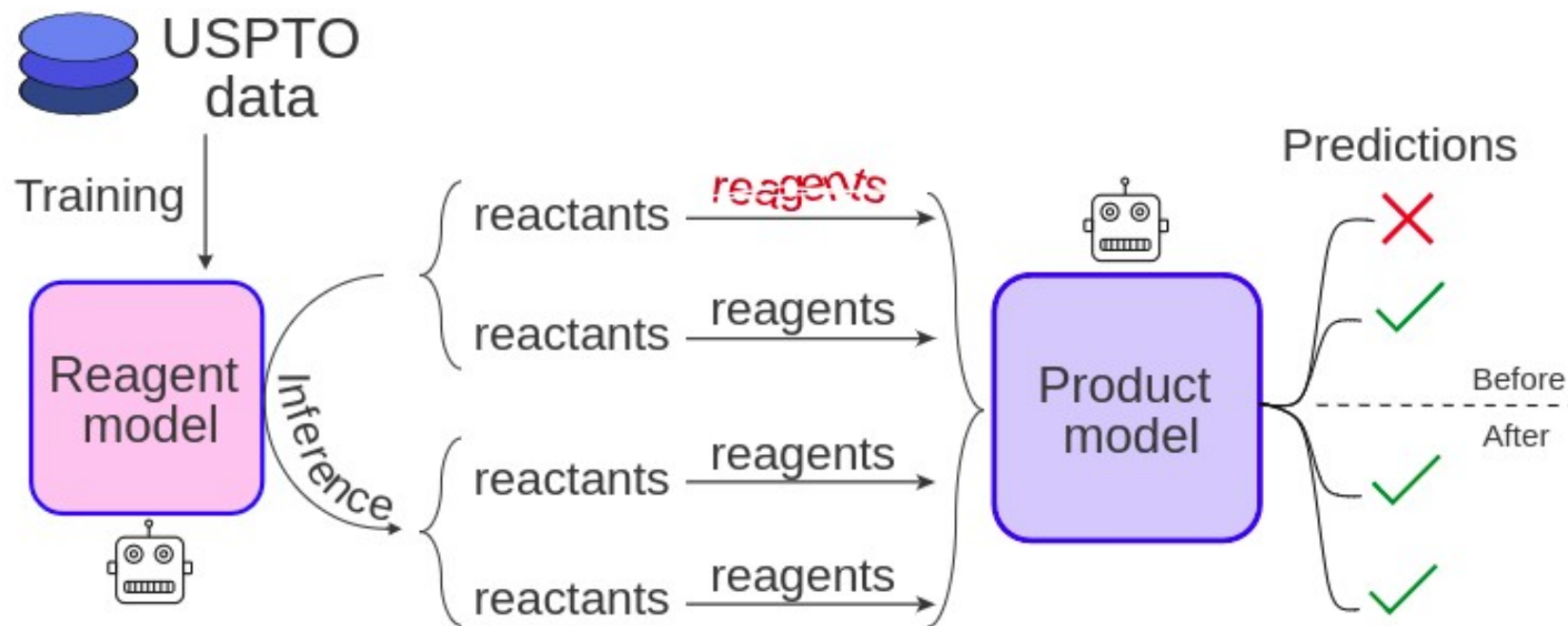
A catalyst, a base and a solvent are necessary.

Paper idea

Reagent and product models are transformers.

We can use a reagent model to improve product prediction models.

Model-agnostic in principle

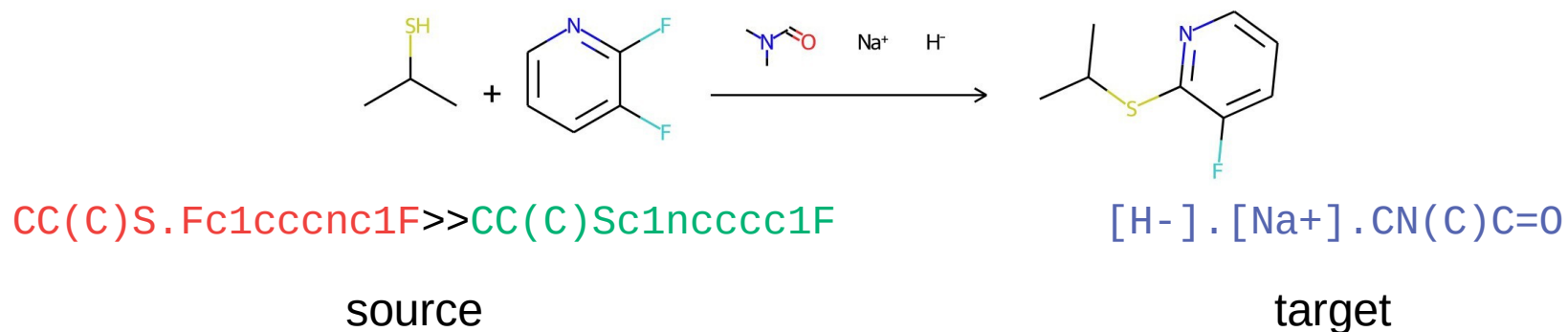


Training set

Training on full USPTO without USPTO MIT test. Final size ~ 1M reactions

Preprocessing:

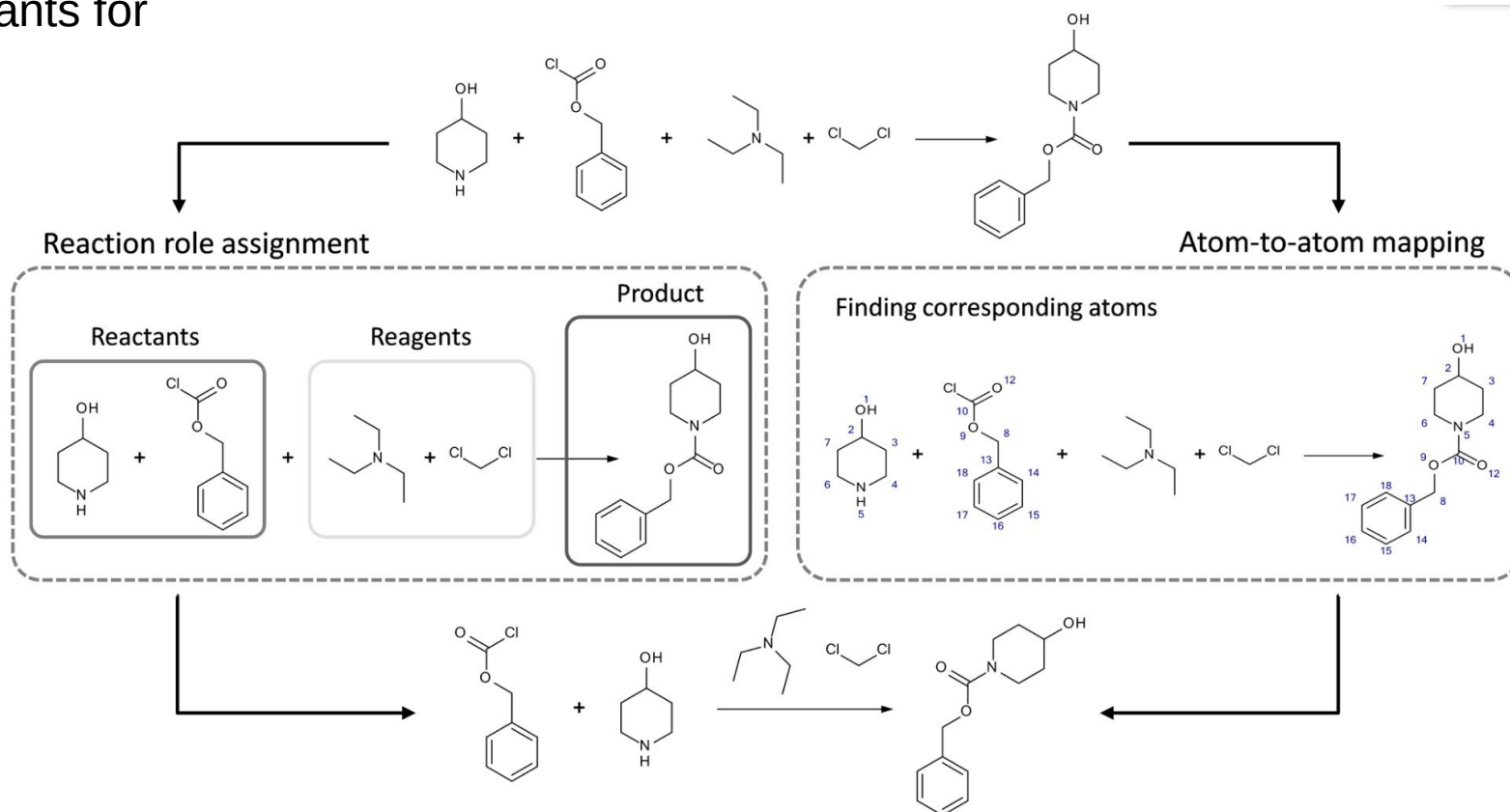
- 1). Delete atom mapping.
- 2). Mix up precursors and extract reagents with RDKit.
- 3). Remove reagents which are too rare.
- 4). Augment data.
- 5). Sort reagents by roles (catalyst, solvent, etc.) using heuristics.



Reaction role assignment

An RDKit procedure to separate reactants for reagents (Schneider et al. 2016).

Atom mapping not needed.



Test set

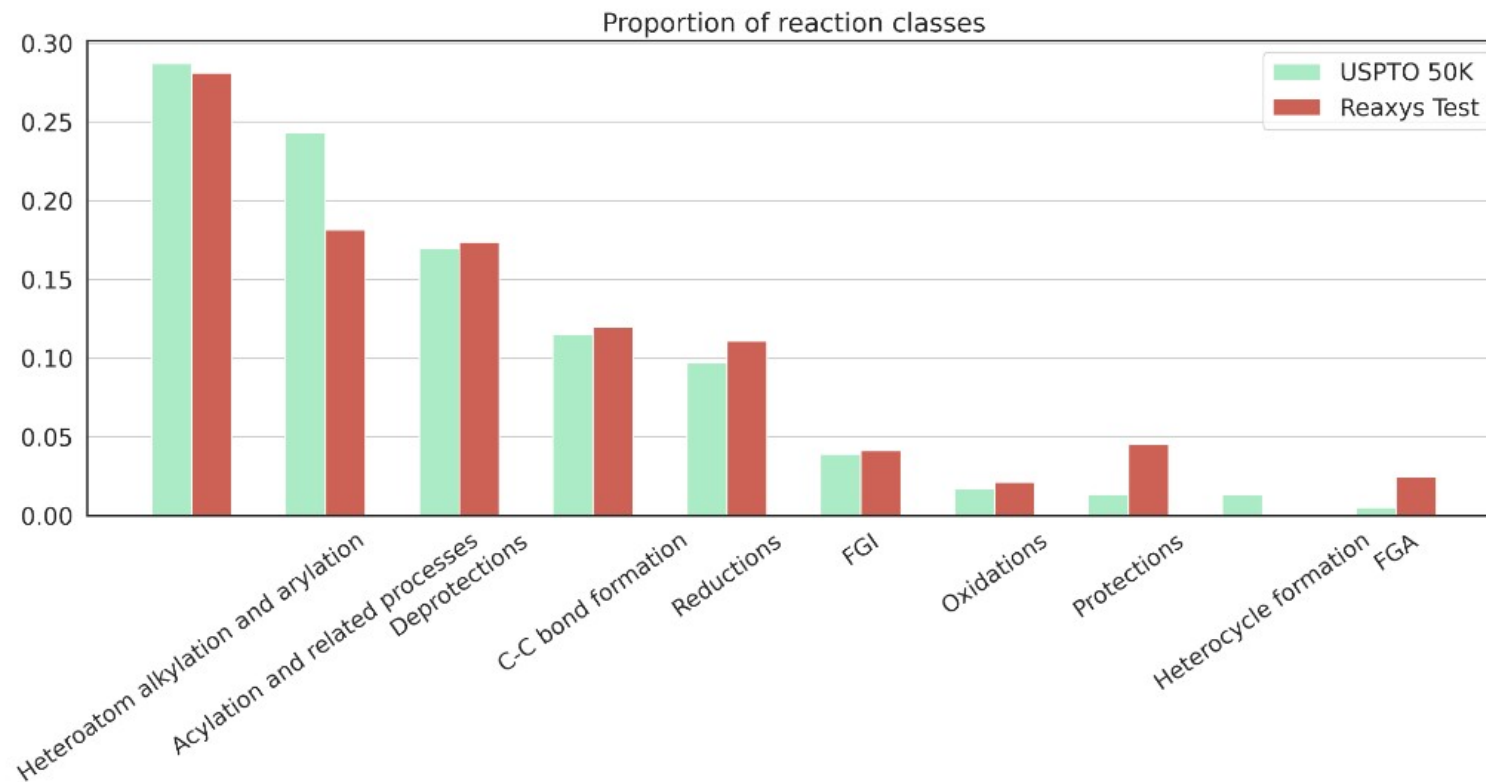
We used a subset of Reaxys for testing purposes.

Size: 96972 reactions.

Reagent SMILES determined by *PubChemPy*.

Reaction types determined by *NameRXN*.

Design goal: similarity to USPTO 50K in terms of types distribution



Discussion: overall performance

Exact match accuracy

Predicted only the molecules in the ground truth and all of them. A.C.B. ~ A.B.C

Partial match accuracy

Some of the molecules are predicted correctly. A.B ~ A.C.D.

Recall

$\#(\text{correctly predicted}) / \#(\text{molecules in target})$

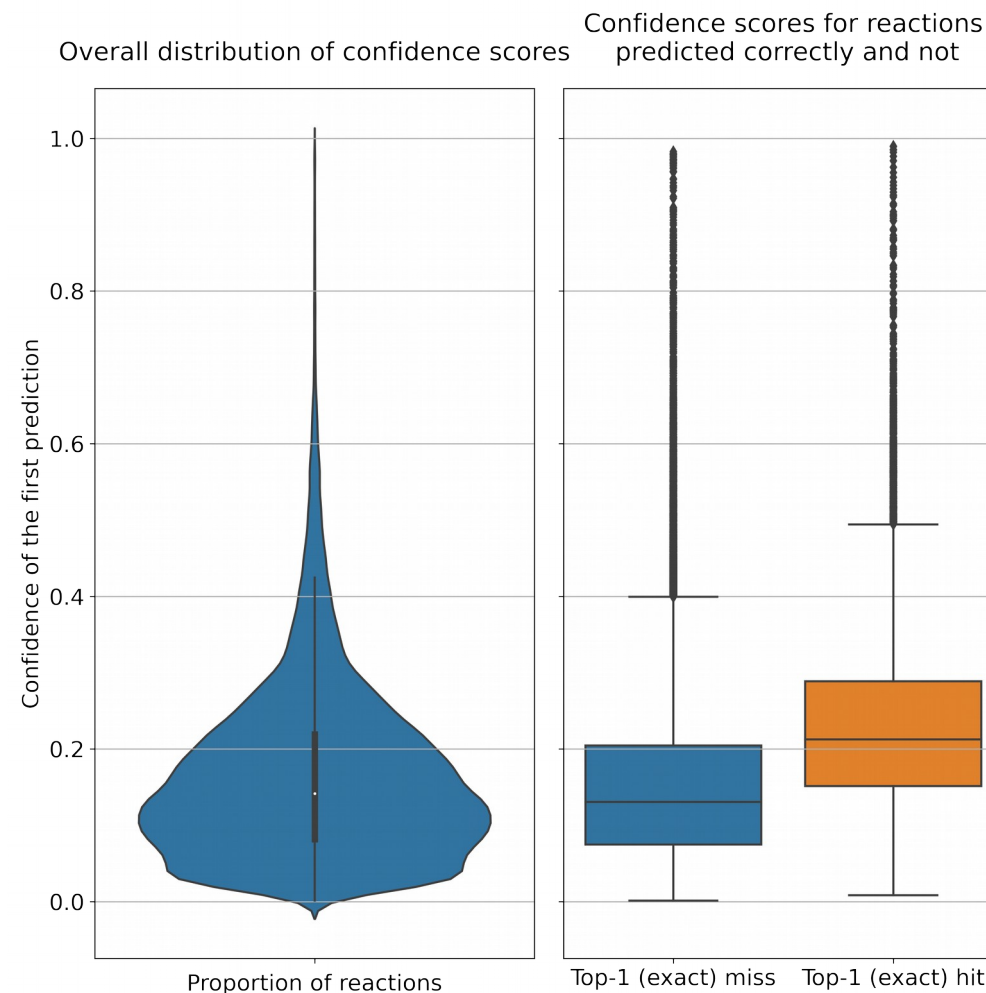
Metric	Top-1	Top-2	Top-3	Top-4	Top-5
Exact match accuracy	17.0	24.7	29.2	31.8	33.5
Full recall	19.2	28.4	31.5	39.3	42.8
Partial match accuracy	70.9	80.5	89.4	87.3	88.9

Model confidence

Confidence: product of the probabilities of all tokens in the generated sequence.

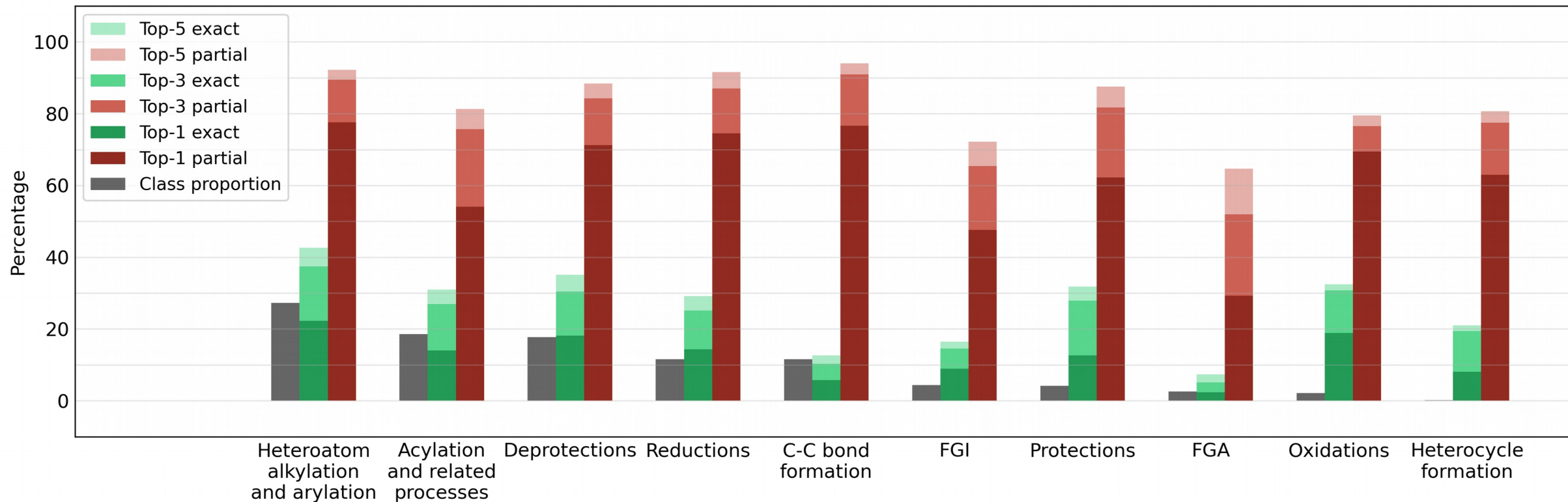
The reagent model is much less confident than a product model.

For the latter, it is close to 1 almost all of the time.



Performance across reaction types

Reagent prediction scores across reaction classes in the Reaxys test set



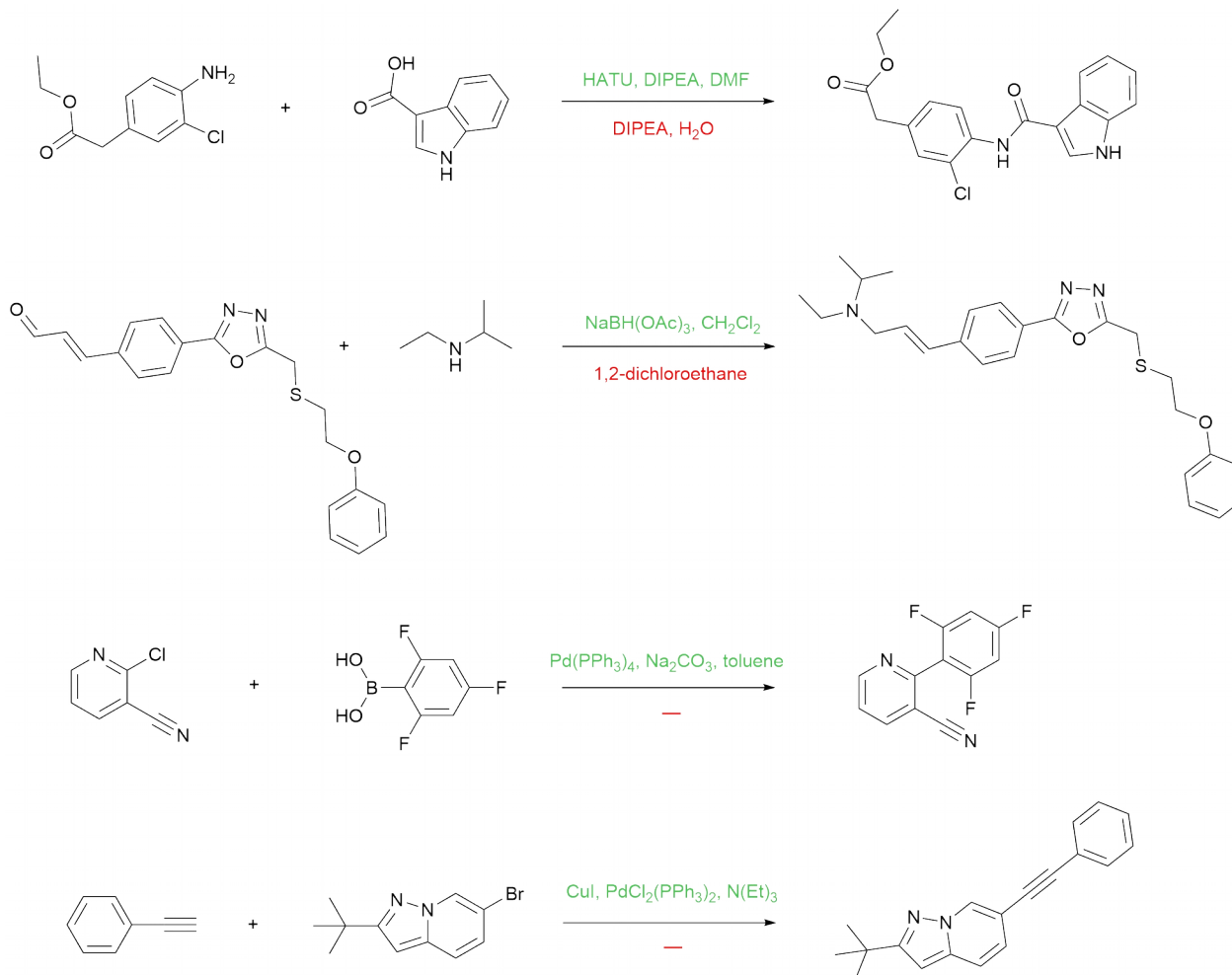
Reagent improvement

Strategy

Replace if more molecules were predicted than there was reported.

Reagents changed in ~25% of reactions

Restored catalysts, reducing agents, etc.



Product prediction

The new model performs better than the old model on both Reaxys and USPTO in both separated and mixed settings.

	Reaxys	USPTO MIT
MT, no reagents	77.3	84.0
MT base, mixed	82.0	87.7
MT new, mixed	83.0	88.3
MT base, separated	84.3	89.2
MT new, separated	84.6	89.6

MT base: trained on basic USPTO.
MT new: trained on USPTO with reconstructed reagents.

Statistical significance: McNemar's test

	F_1 incorrect	F_1 correct
F_2 incorrect	A	B
F_2 correct	C	D

$$x = \frac{(|B - C| - 1)^2}{B + C}$$

Chi-squared distribution (1 degree of freedom).

Null hypothesis: difference is accidental

Conclusion

- > Transformer can be successfully used to suggest reagents for organic reactions.
- > We used the strategy to train a model on USPTO and test it on Reaxys.
- > We used a reagent model to improve a product model in a *self-supervised* and *model-agnostic* fashion.
- > We beat the score of the Molecular Transformer on USPTO MIT.



Thank you for your attention!