

Combining high-quality, humanly curated data  
with language models:

The dawn of on-demand machine learning  
models for digital chemistry

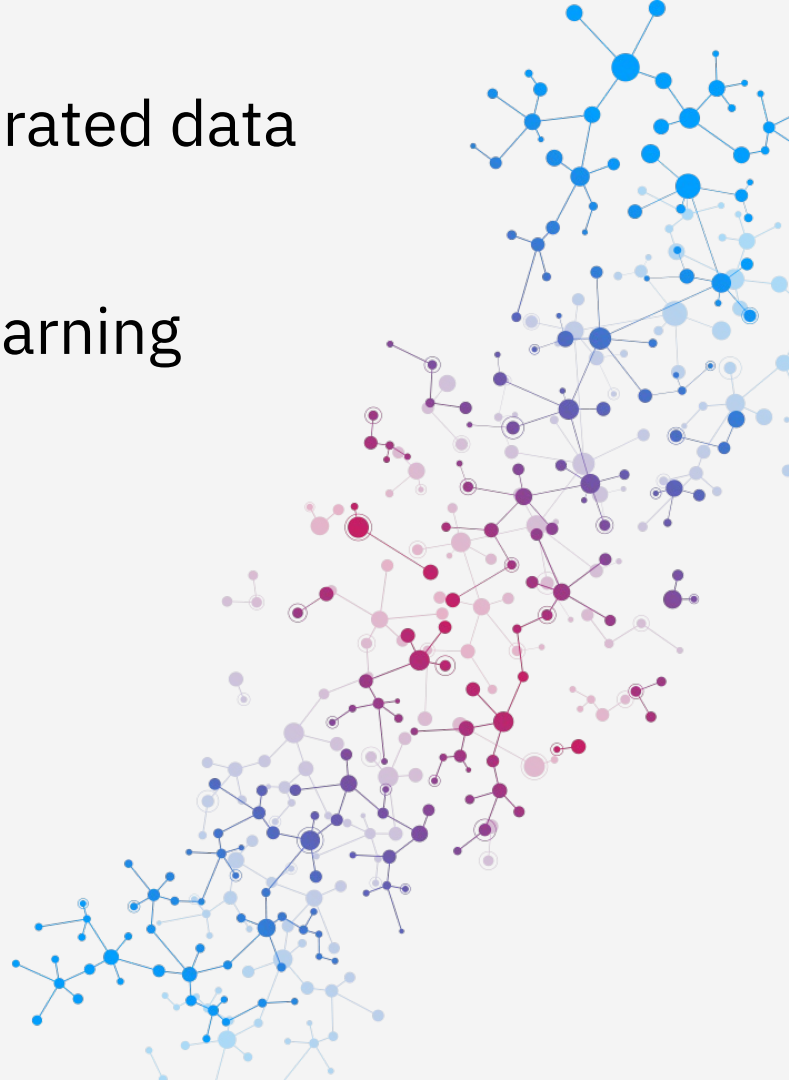
**Alain Vaucher**, *IBM Research*

Fiona Shortt de Hernandez, *Thieme Group*

Sascha Hausberg, *Thieme Group*

Teodoro Laino, *IBM Research*

ACS Fall 2022  
24 August 2022



## OUTLINE

1. Introduction:  
AI for chemical reactivity & data
2. Study (IBM Research & Thieme):  
combining curated data & data-driven models
  - Data & model training
  - Examples, challenges, expert feedback
3. Where do we go from here?

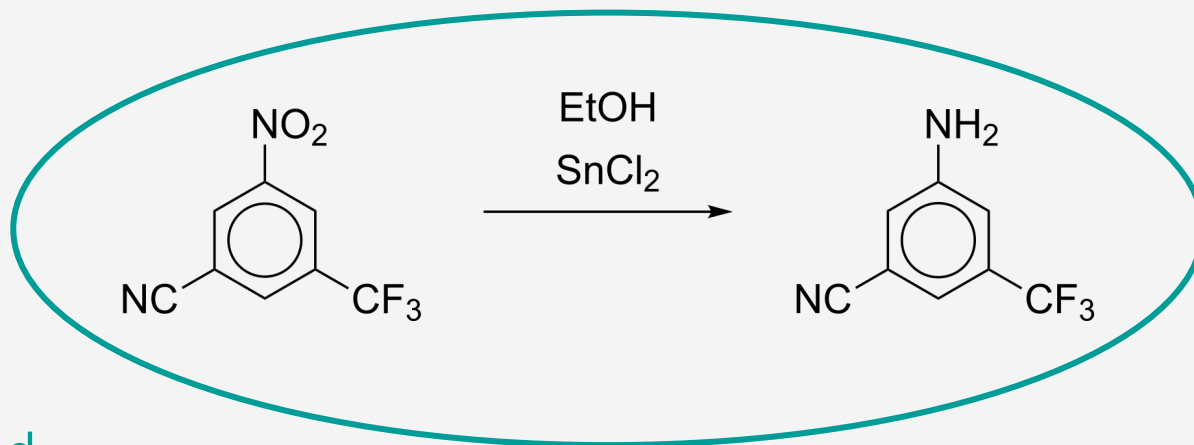
## OUTLINE

1. Introduction:  
AI for chemical reactivity & data
2. Study (IBM Research & Thieme):  
combining curated data & data-driven models
  - Data & model training
  - Examples, challenges, expert feedback
3. Where do we go from here?

# AI and chemical reactivity

Experimental conditions?

Product?



Related reactions?

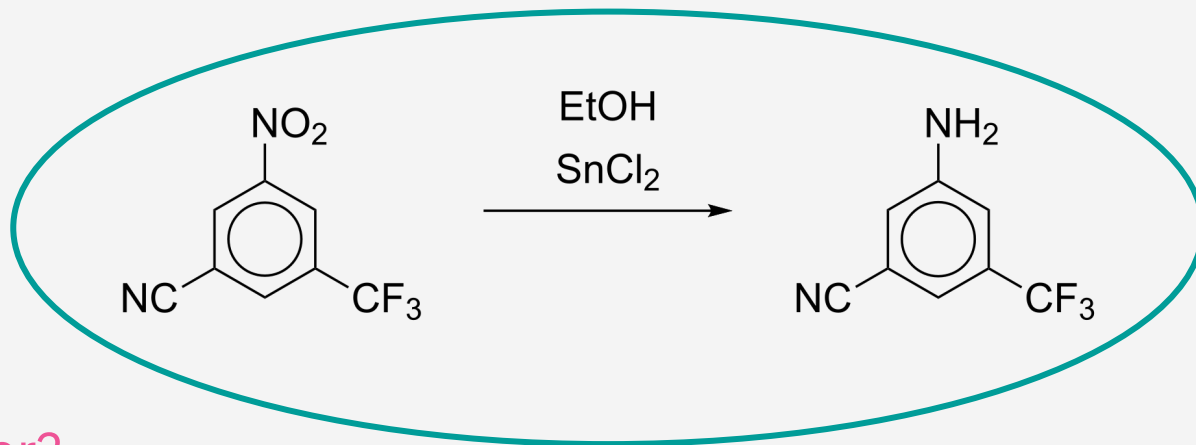
Retrosynthesis?

Yield?

# AI and chemical reactivity

Accuracy

State-of-the-art

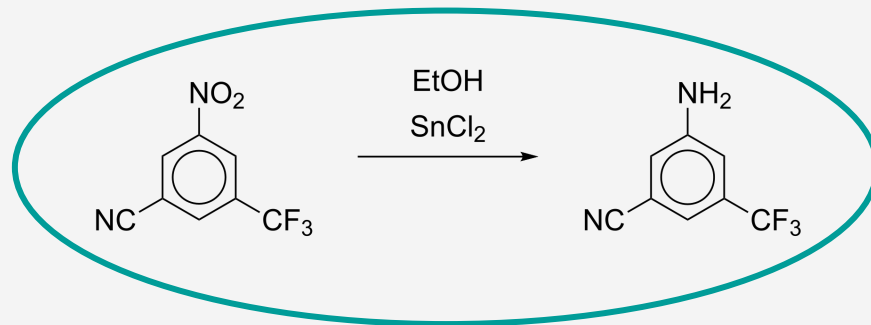


End user?

Applicability?

Usefulness?

# AI and chemical reactivity



## Applicability?

Training data must cover adequate chemistry

## End user?    Usefulness?

What models are needed? How are they used?

# AI and chemical reactivity – data sources

- Patents (USPTO [1], Pistachio [2], SciWalker [3])
- Scientific publications
- Proprietary reactions (industry)
- Publishers (Reaxys, CASFinder, Thieme)
- Others (Open Reaction Database [4])
- Etc.

[1] <https://dx.doi.org/10.6084/m9.figshare.5104873.v1>

[2] <https://www.nextmovesoftware.com/pistachio.html>

[3] <https://www.sciwalker.com>

[4] <https://open-reaction-database.org>



# AI and chemical reactivity – data sources

- **Patents** (USPTO [1], Pistachio [2], SciWalker [3])
- Scientific publications
- Proprietary reactions (industry)
- Publishers (Reaxys, CASFinder, Thieme)
- Others (Open Reaction Database [4])
- Etc.

[1] <https://dx.doi.org/10.6084/m9.figshare.5104873.v1>

[2] <https://www.nextmovesoftware.com/pistachio.html>

[3] <https://www.sciwalker.com>

[4] <https://open-reaction-database.org>





# AI and chemical reactivity – going beyond patents

- Patent data:
  - Good coverage of mid-size organic compounds
  - Good coverage of common organic reactions
  - essential for development of data-driven models
- Limitations:
  - Reporting errors
  - Incorrect extraction
  - Reproducibility concerns
  - Limited coverage of reaction space

# AI and chemical reactivity – going beyond patents

- If patent data is not enough:
  - Fine-tuning models for specific application
  - Push for more general models
- Open questions:
  - What data sources?
  - Data safety?
- Study: reactivity models trained on curated data

# AI and chemical reactivity – IBM RXN & Thieme SOS



Combine **IBM RXN models** with **curated data from Thieme**

- How useful are the predictions for chemists?
- How valuable is the curated data for ML applications?

## OUTLINE

1. Introduction:  
AI for chemical reactivity & data
2. Study (IBM Research & Thieme):  
combining curated data & data-driven models
  - Data & model training
  - Examples, challenges, expert feedback
3. Where do we go from here?

## OUTLINE

1. Introduction:  
AI for chemical reactivity & data
2. Study (IBM Research & Thieme):  
combining curated data & data-driven models
  - Data & model training
  - Examples, challenges, expert feedback
3. Where do we go from here?

# Curated data (Thieme)



## **Science of Synthesis**

~450,000 reactions

Collections & volumes covering  
diverse topics in chemical  
synthesis



## **Synfacts**

~16,000 reactions (2017-2018)

Highlights in chemical synthesis  
Focus on total synthesis

# Science of Synthesis Reference Library

## Stereoselective Synthesis Volume 1, section 1.13.1.2

Synfacts 2021; 17(12): 1314  
DOI: 10.1055/s-0041-1737090

1.13.1 Hydroamination of Simple Alkenes

691

**1-Methyl-3-(octan-2-yl)imidazolidin-2-one (2, R<sup>1</sup> = Me); Typical Procedure**<sup>[21]</sup>  
A suspension of 1-methylimidazolidin-2-one (20 mg, 0.20 mmol), oct-1-ene (1.3 g, 12 mmol),  $(\text{AcCl})_2\text{Os-1}$  (8.0 mg, 5.0  $\mu\text{mol}$ ), and AgOTf (2.6 mg, 10  $\mu\text{mol}$ ) in *m*-xylene (0.5 mL) was stirred at 100 °C for 48 h. The crude mixture was filtered through a plug of silica gel, concentrated, and chromatographed (hexanes/EtOAc 5:1 to 1:1) to give a colorless oil; yield: 37 mg (86%); 76% ee.

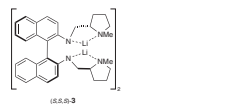
### 1.13.1.2 Cyclization of Aminoalkenes

#### 1.13.1.2.1 Using Chiral Alkali Metal Based Catalysts

The first reports on base-catalyzed additions of amines to alkenes date back 60 years. In particular, multiple catalyst systems utilizing alkali metals have been reported.<sup>[8,14,15]</sup> However, application of chiral alkali metal complexes in the asymmetric hydroamination of nonactivated aminoalkenes has drawn little attention to date.<sup>[16,17]</sup> Attempts to perform asymmetric hydroamination utilizing chiral alkaline earth metal complexes have been thwarted by facile Schlenk equilibria of the metal species in solution.<sup>[14,25]</sup>

The proline-derived dimeric diamidinaphthyl dilithium salt (**1**), which is prepared via deprotonation of the corresponding tetraamine with butyllithium, catalyzes asymmetric intramolecular hydroamination reactions of aminoalkenes **4** at or below ambient temperatures to form pyrrolidines **5** with enantioselectivities of up to 74% ee (Scheme 3).<sup>[22,26]</sup> The enantioselectivities may be improved to up to 85% ee by lowering the reaction temperature to -10 °C.<sup>[26]</sup> The unique reactivity of (**1**) is believed to derive from the close proximity of the two lithium centers chelated by the proline substituents, because more simple lithium amides require significantly higher reaction temperatures and give inferior selectivities.

**Scheme 3** Lithium-Catalyzed Asymmetric Hydroamination/Cyclization of Aminoalkenes<sup>[22,26]</sup>



R <sup>1</sup>	R <sup>2</sup>	Catalyst (mol%)	Temp (°C)	Time (h)	Yield (%)	ee (%)	Ref
Me	Me	2.5	22	45	93	67	[22]
(CH <sub>2</sub> ) <sub>5</sub>	5	5	20	2	82	74	[22]
(CH <sub>2</sub> ) <sub>5</sub>	2	2	-10	22	84	85	[26]

for references see p 727

### Category

Synthesis of Natural Products and Peptides

### Key words

BMS-986158

Stille coupling

Chan-Lam coupling

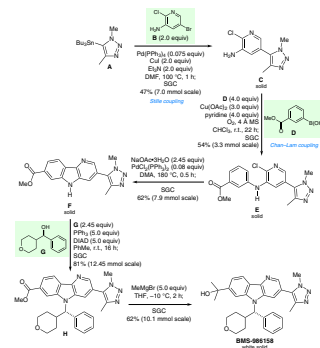
copper catalysis

palladium-catalyzed C-H activation

Mitsunobu reaction

A. V. CAVAI<sup>1</sup> ET AL. (BRISTOL MYERS SQUIBB COMPANY, PRINCETON, USA)  
Discovery and Preclinical Pharmacology of an Oral Bromodomain and Extra-Terminal (BET) Inhibitor Using Scaffold-Hopping and Structure-Guided Drug Design  
*J. Med. Chem.* 2021, 64, 14247-14265, DOI: 10.1021/acs.jmedchem.1c00025.

### Synthesis of BMS-986158



**Significance:** BMS-986158 is an inhibitor of the bromodomain and extra-terminal (BET) family of adaptor proteins that are involved in the transcriptional regulation of key oncogenes. It has entered phase 1 [2a] clinical trials in patients with advanced cancers and hematologic indications including myelofibrosis.

**Comment:** Key steps in the small-scale discovery synthesis of the 5H-pyrido[3,2-b]indole core of BMS-986158 are (1) the copper-catalyzed oxidative coupling of the chloropyridine **C** with the boronic acid **D** (Chan-Lam coupling) and (2) the palladium-catalyzed C-H activation reaction **E** → **F**.

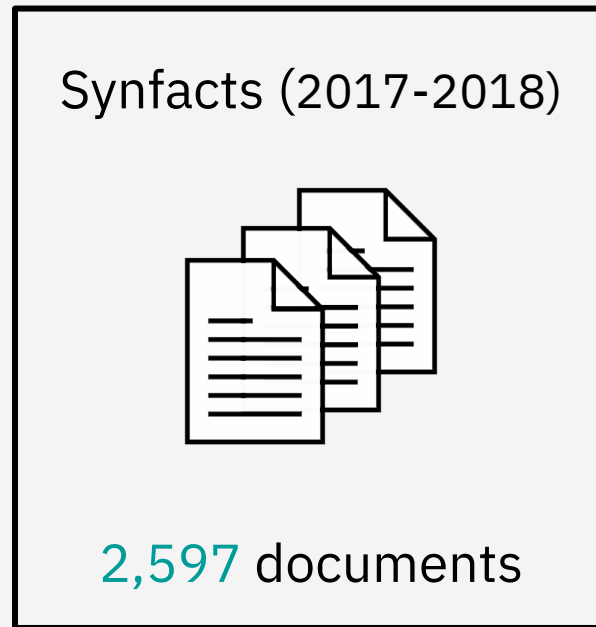
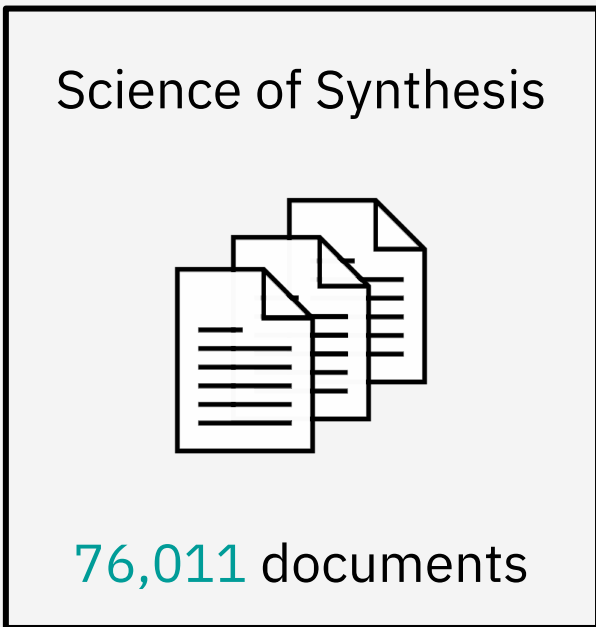
This document was downloaded for personal use only. Unauthorized distribution is strictly prohibited.

SYNFACTS Contributors: Philip Koelsch  
Synfacts 2021, 17(12), 1314. Published online: 17.11.2021  
DOI: 10.1055/s-0041-1737090, Rev. No.: 0073219

© 2021, Thieme. All rights reserved.  
Georg Thieme Verlag KG, Rüdigerstraße 14, 7030 Stuttgart, Germany

1314

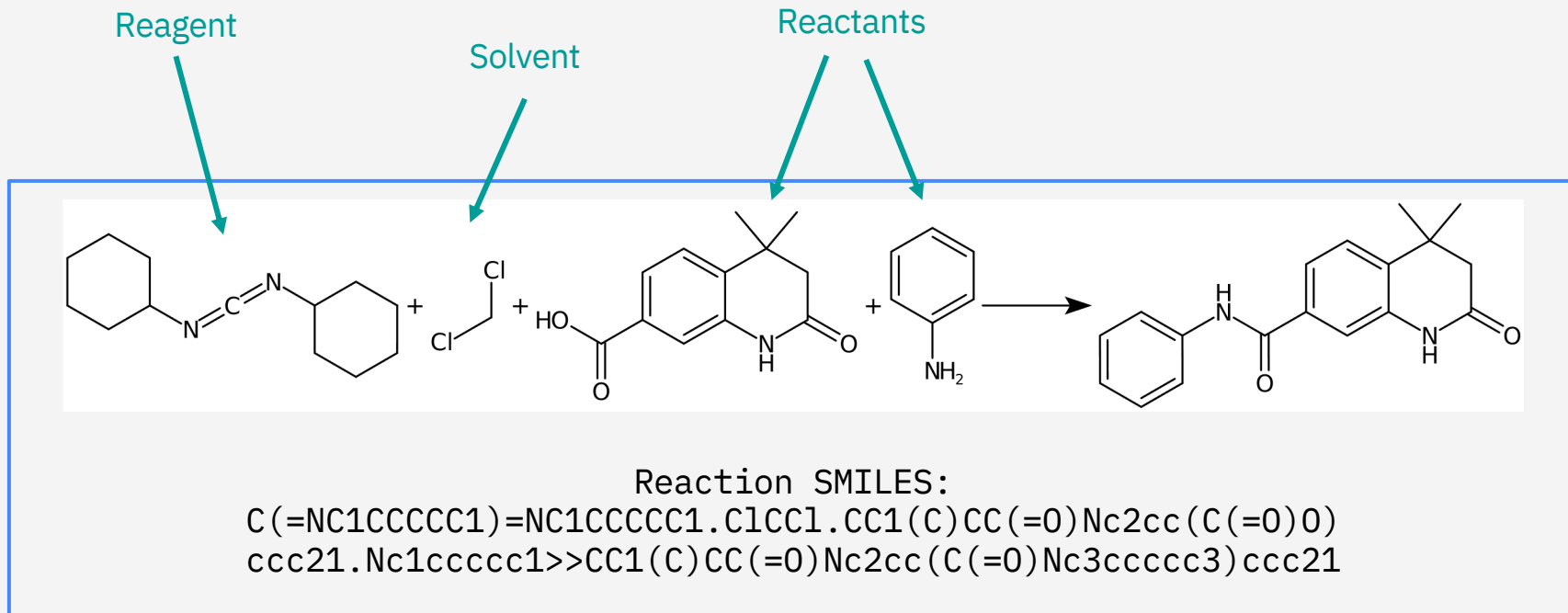
# Data - summary



PDF, XML and RDF files



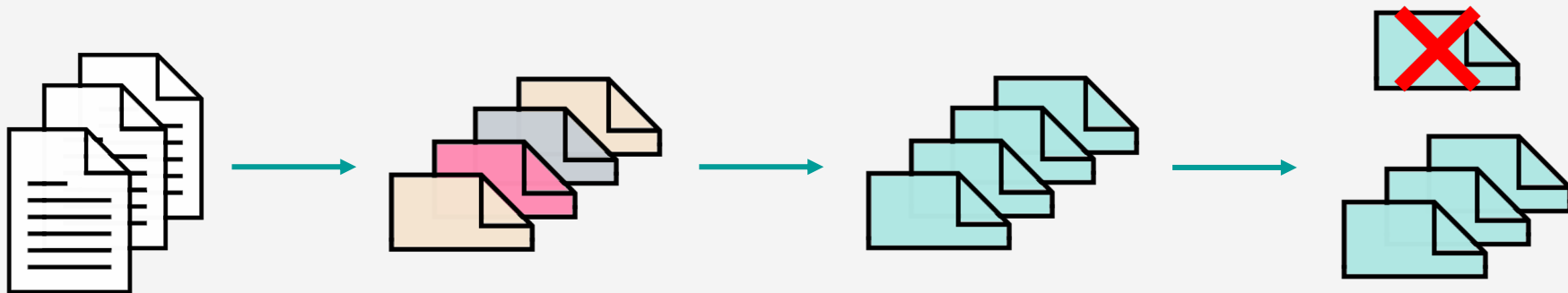
# Target reaction format



# Data Processing

Steps to produce training data:

1. Conversion to reaction SMILES
2. Standardization
3. Sanity checks and filters



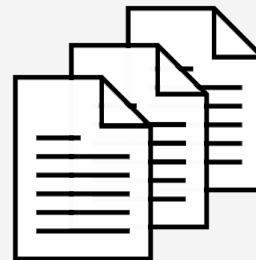
# Processed Data

Science of Synthesis



318,383 reactions

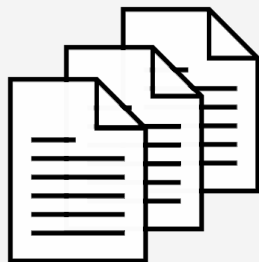
Synfacts (2017-2018)



13,818 reactions

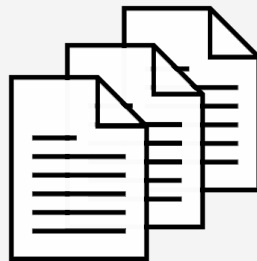
# Processed Data

Science of Synthesis



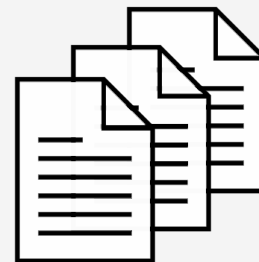
318,383 reactions

Synfacts (2017-2018)



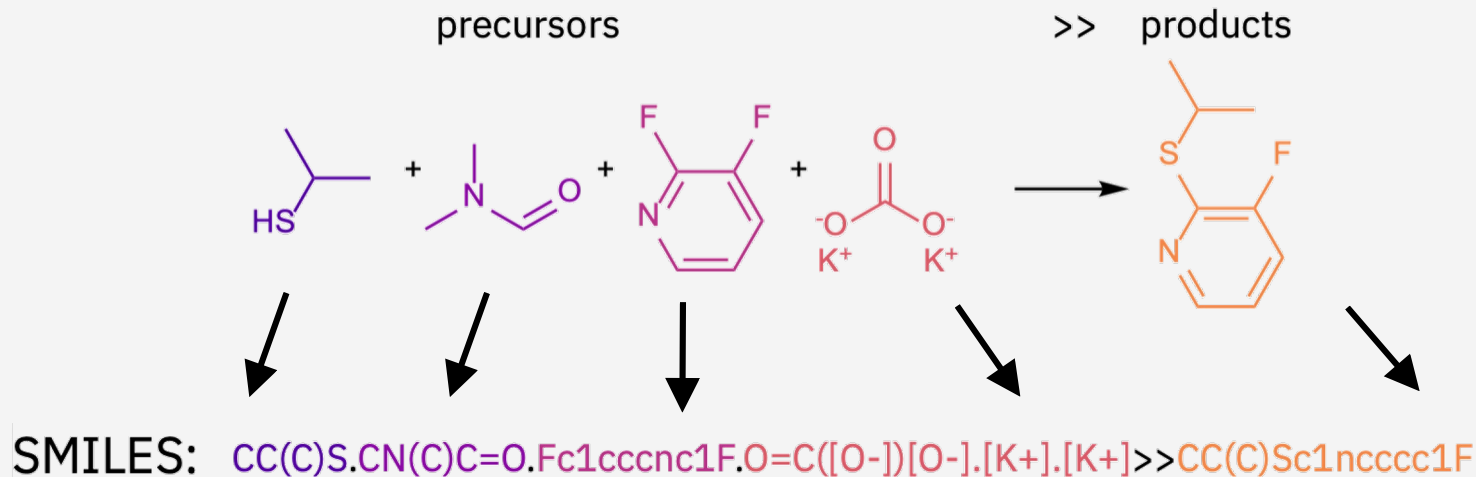
13,818 reactions

Pistachio



~2M reactions

# Atoms as *letters*, molecules as *words*

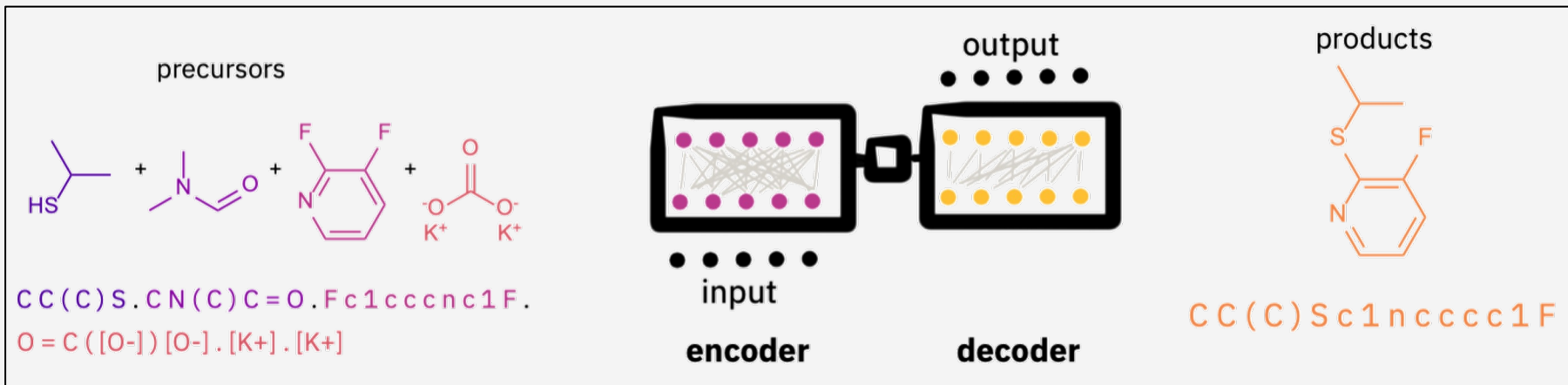


Split -> "sequences of atoms"

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

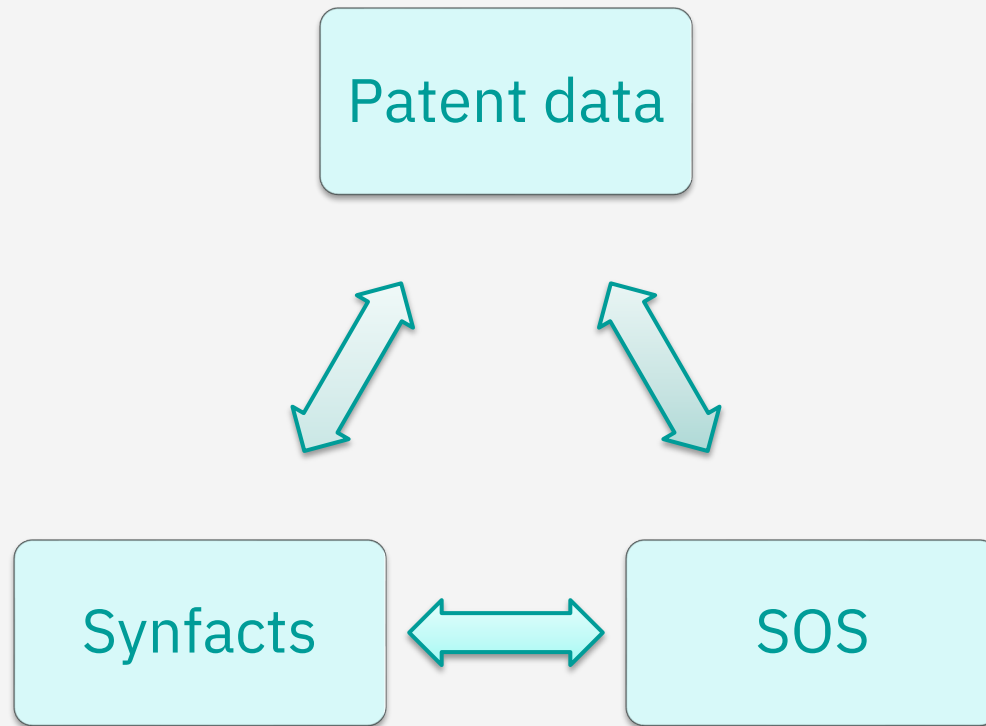
→ Borrow methods developed for human languages

# Training the models



Forward reaction prediction and retrosynthesis

# Training the models



3 sources included for training (weighted differently)

# Model accuracy (top-1)

## Forward reaction prediction

Patent-only: 22.9%



Patent + Thieme: **69.5%**

## Retrosynthesis

Patent-only: 1.1%



Patent + Thieme: **17.8%**

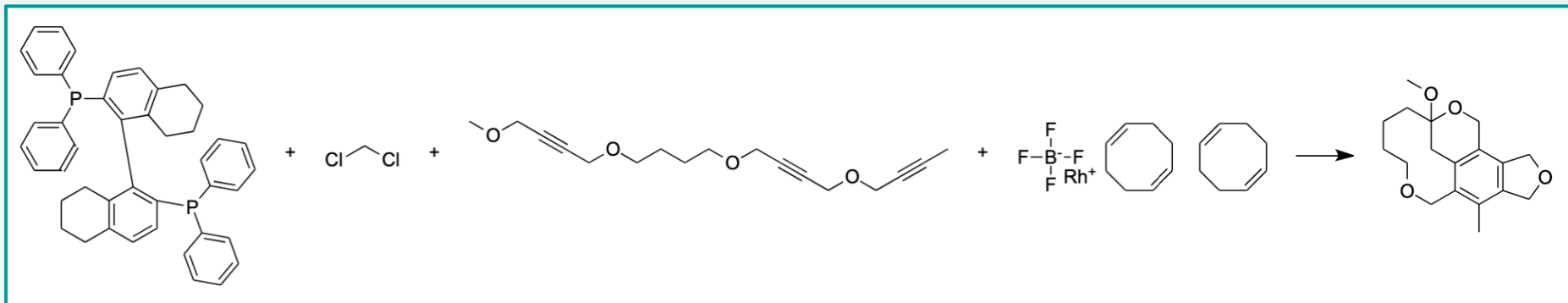


## OUTLINE

1. Introduction:  
AI for chemical reactivity & data
2. Study (IBM Research & Thieme):  
combining curated data & data-driven models
  - Data & model training
  - Examples, challenges, expert feedback
3. Where do we go from here?

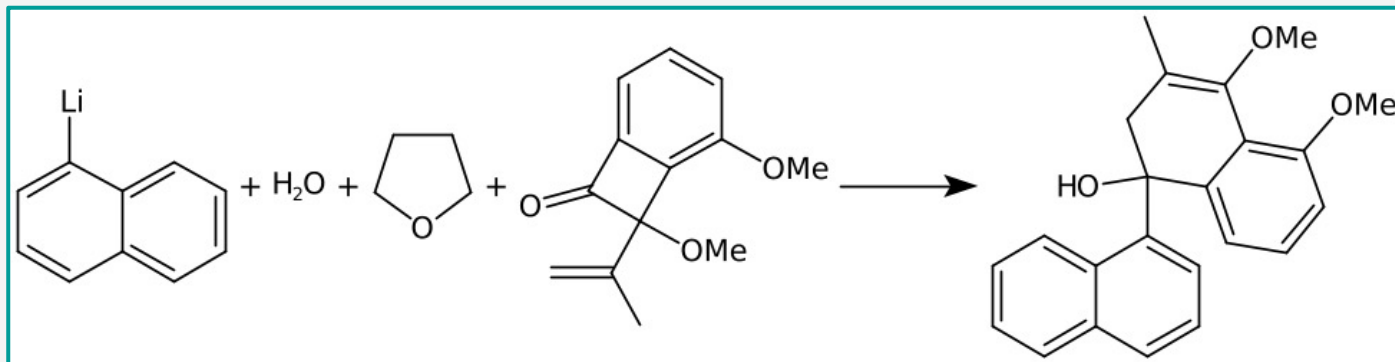
# Forward model: inspection of a few predictions

Correct product predicted



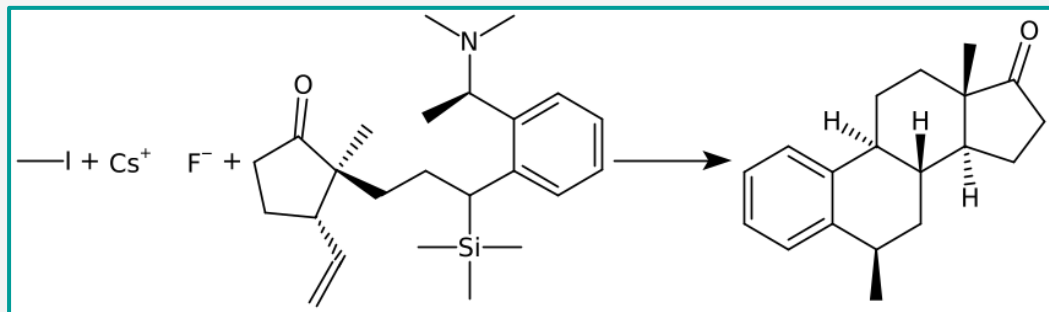
# Forward model: inspection of a few predictions

Correct product predicted

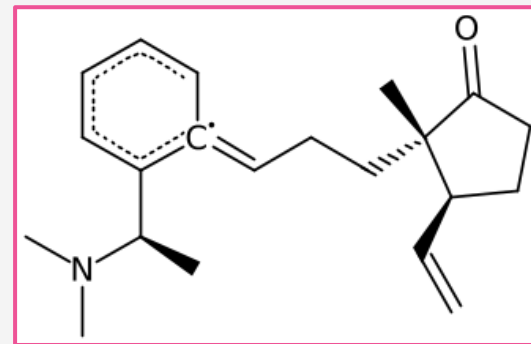


# Forward model: inspection of a few predictions

Thieme reaction

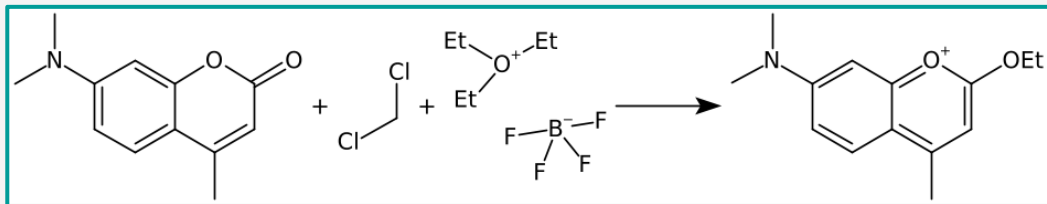


Predicted product

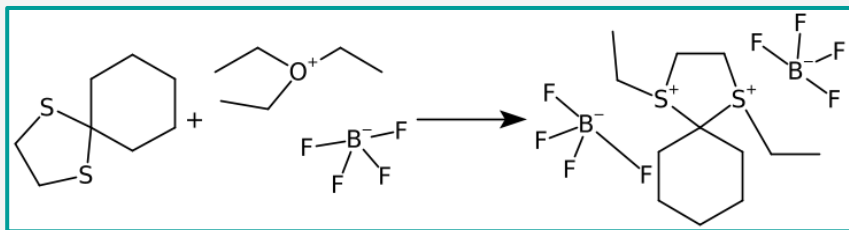
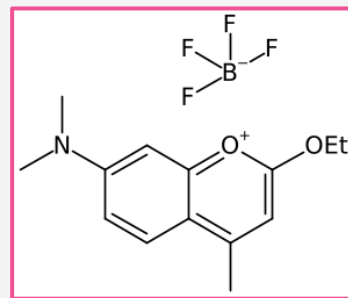


Invalid SMILES: 1.0% of predictions

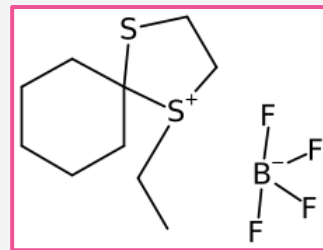
# Forward model: inspection of a few predictions



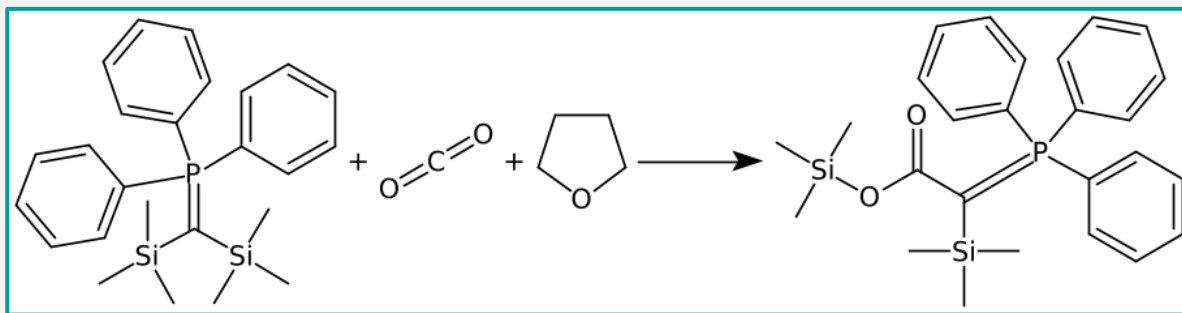
Predicted product



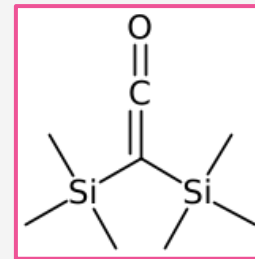
Predicted product



# Forward model: inspection of a few predictions

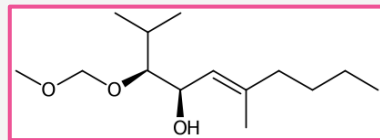
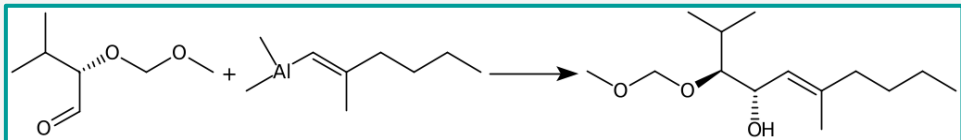


Predicted product

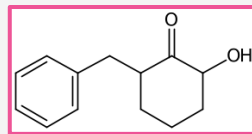
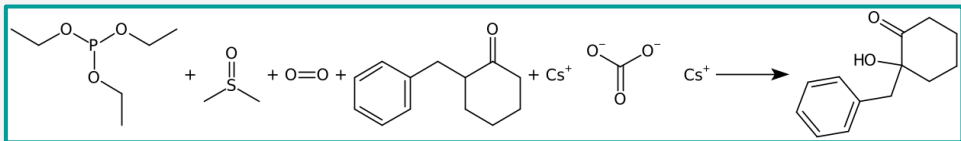


# Forward model: inspection of a few predictions

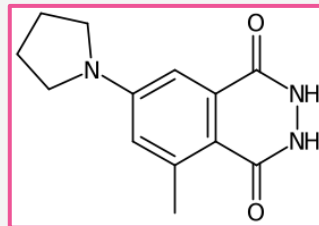
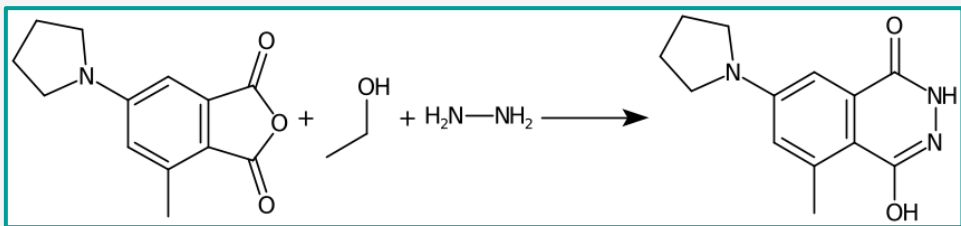
27%: identical molecular formula



stereochemistry



regiochemistry



tautomers

# Expert feedback

- What we asked:
  - Does curated data help?
  - What would make the models even more useful?
- Feedback:
  - Curated data helps! Considerable difference to patent-only.
  - There are still some errors
  - Comments on diversity, usability (see next slides)



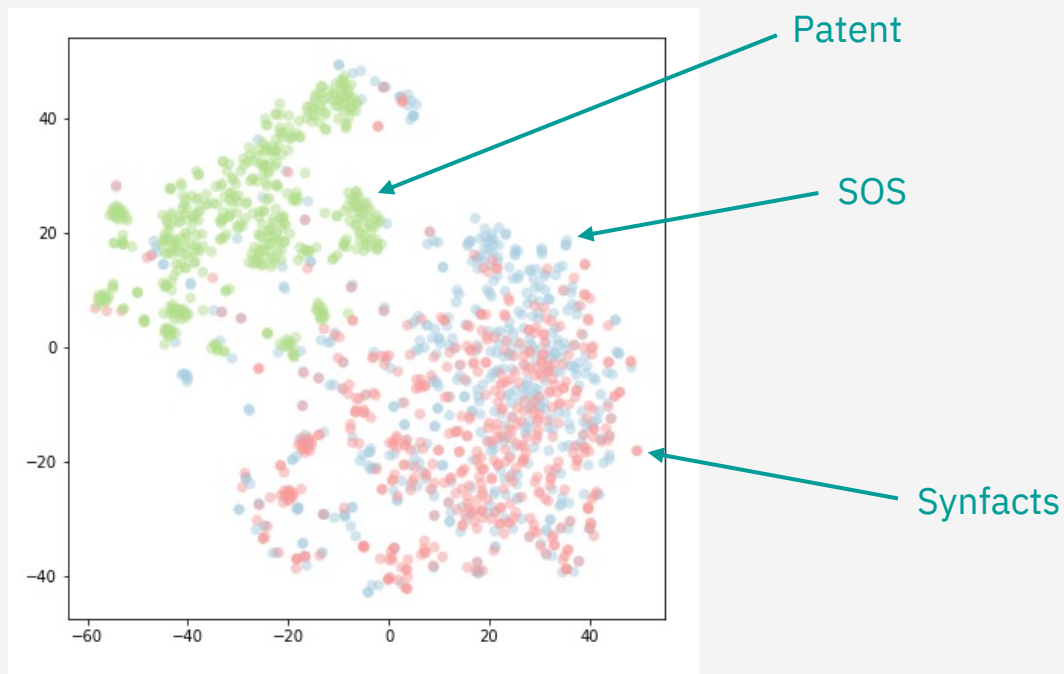
# Related reactions

“It would be nice to see literature references for similar/related reactions.”

The screenshot displays a chemical reaction analysis interface. At the top left, it shows 'Confidence: 0.848' and 'High confidence'. The main area features a central target molecule, (S)-1-phenylethylamine hydrochloride, with the SMILES string NC(=O)[C@H]1C=CC=C1. Below this, a blue box labeled '0.848 - Ester Schotten-Baumann' indicates the reaction type. Three related molecules are shown below: dichloromethane (ClCCl), thionyl chloride (NC(=O)Cl), and (S)-1-phenylethanol (O[C@H]1C=CC=C1). To the right, a 'Similar reactions list' is provided. It contains two entries: 'Reaction 1' with a score of 0.978 and 'Reaction class' N/A, and 'Reaction 2' with a score of 0.969 and 'Reaction class' N/A. Each reaction entry includes a chemical reaction scheme, a document ID (SD-018-00635 and SD-011-00565), a URL, and the dataset name 'thieme'. Navigation arrows are visible at the bottom of each reaction entry.

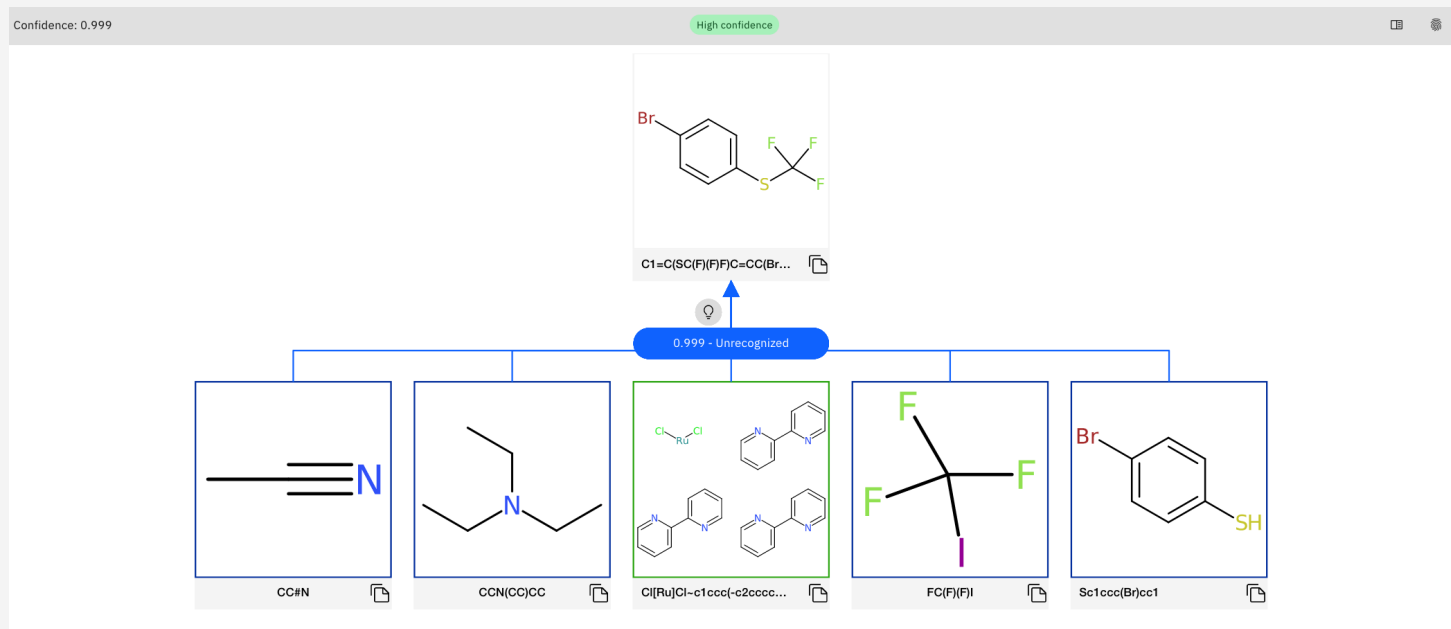
# Diversity of predictions

“The retrosynthetic sequences lack variety.”

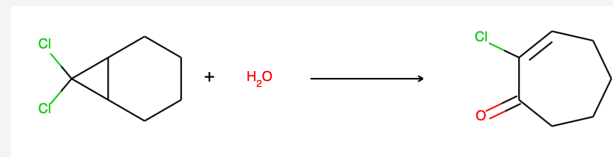
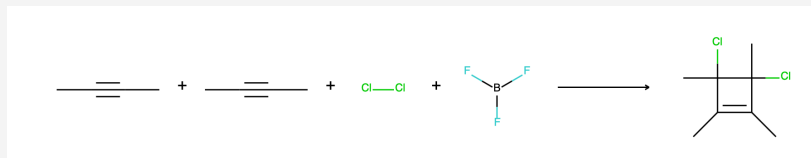
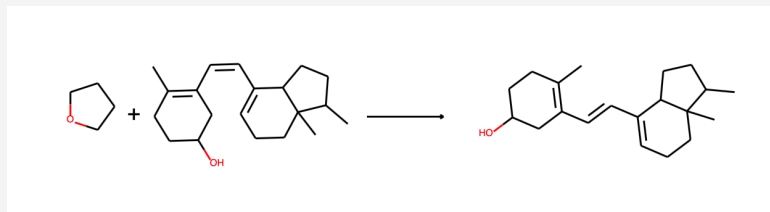
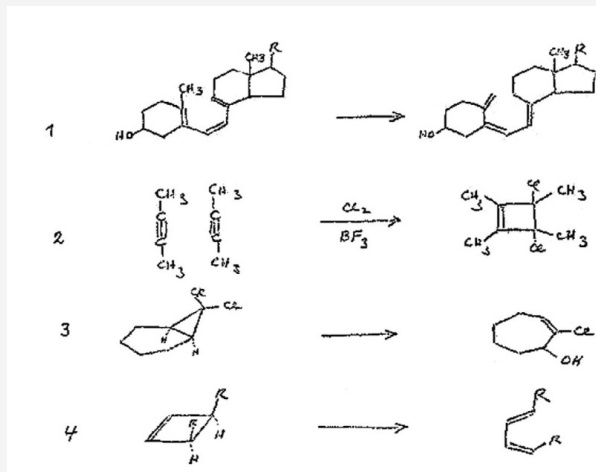


# Photochemical and thermal reactions

“Photochemistry and unimolecular reactions do not work.”

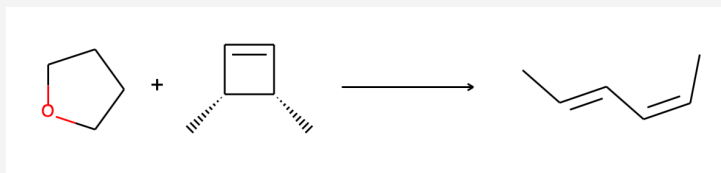


# Photochemical and thermal reactions

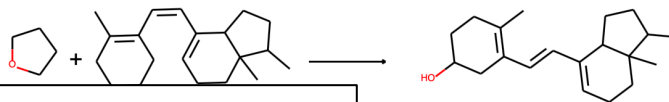


## 3. Woodward's Four Mysterious Reactions

Tantillo, Seeman, Chem. Eur. J.  
2021, 27, 7000 – 7016



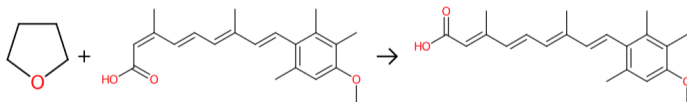
# Photochemical and thermal reactions



## Similar reactions list

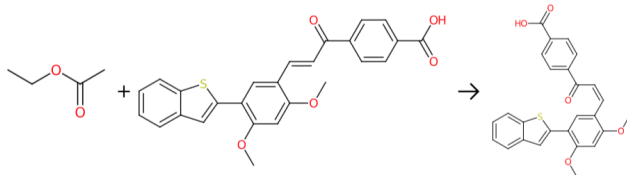
### Reaction 1

Score 0.989 Reaction class N/A



### Reaction 2

Score 0.977 Reaction class N/A

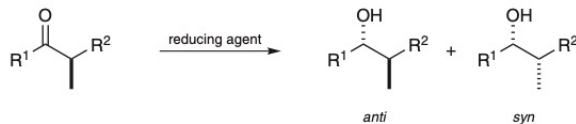


Closest reactions come from patents and are questionable



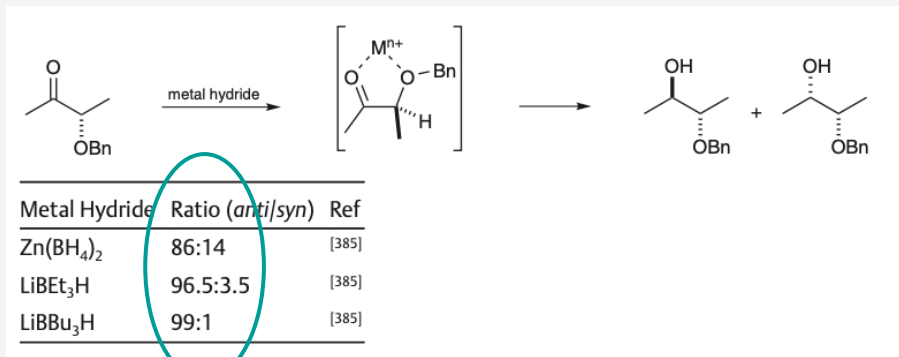
# Challenges: Multiple products

**Scheme 59** Diastereoselective Reduction of  $\alpha$ -Substituted Acyclic Ketones with Borohydrides and Disiamylborane<sup>[364,365]</sup>



R <sup>1</sup>	R <sup>2</sup>	Reducing Agent	Ratio ( <i>anti</i> / <i>syn</i> )	Ref
Ph	Et	NaBH <sub>4</sub> /MeOH	57:43	[364]
Ph	CH=CH <sub>2</sub>	NaBH <sub>4</sub> /MeOH	70:30	[364]
Ph	C≡CH	NaBH <sub>4</sub> /MeOH	89:11	[364]
Me	CHMePh	LiB <sub>s</sub> -Bu <sub>3</sub> H/THF	96:4	[365]
Me	CHMePh	Sia <sub>2</sub> BH/THF	20:80	[365]

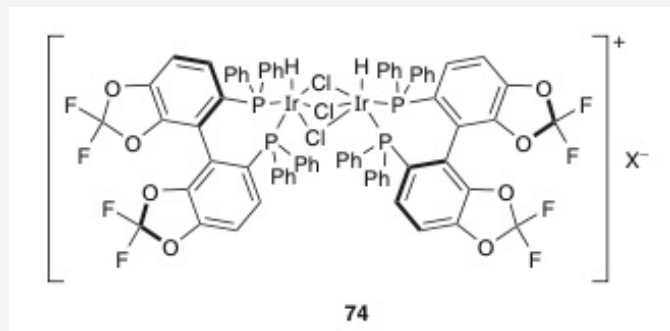
<sup>a</sup> Sia = siamyl (3-methylbutan-2-yl).



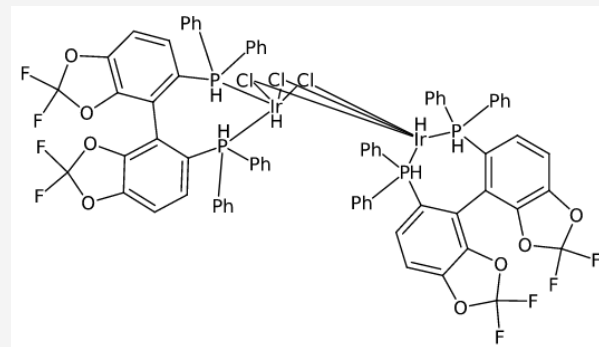
Metal Hydride	Ratio ( <i>anti</i> / <i>syn</i> )	Ref
Zn(BH <sub>4</sub> ) <sub>2</sub>	86:14	[385]
LiEt <sub>3</sub> H	96.5:3.5	[385]
LiBBu <sub>3</sub> H	99:1	[385]

# Challenges: Compound representation

PDF



Extracted structure



- Hydrogen atoms added to P
- Standardization failing (hypervalent chlorine)



# Challenges

- Multiple products
- Compound representation
- Standardization, normalization
- Reagents with no SMILES
- Multiple steps drawn as one reaction
- Etc.

Note: Most of these challenges apply to all datasets!

## OUTLINE

1. Introduction:  
AI for chemical reactivity & data
2. Study (IBM Research & Thieme):  
combining curated data & data-driven models
  - Data & model training
  - Examples, challenges, expert feedback
3. Where do we go from here?

# Where do we go from here?

- Curated data does have value for ML
- Handling of proprietary data for ML:
  - Data safety must be preserved
  - Train models on trusted servers? Federated learning?
  - Address concern of data leaking
- Possibility to select which data to train new models on
- Keep listening to and learning from chemists!

# Thank you for your attention!

**If you have any questions:**

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

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

## Acknowledgments:

### **IBM Research**

Teodoro Laino  
Philippe Schwaller

### **Thieme Group**

Fiona Shortt de Hernandez  
Sascha Hausberg  
Klaus Köberlein

### **Experts (evaluation of models)**

Prof. Dame Margaret Brimble (University of Auckland, New Zealand)  
Prof. Alois Fürstner (MPI Mülheim, Germany)  
Prof. Karl Gademann (University of Zurich, Switzerland)  
Prof. Ang Li (Shanghai Institute of Organic Chemistry, China)  
Prof. Cristina Nevado (University of Zurich, Switzerland)  
Prof. Richmond Sarpong (University of California, Berkeley, USA)  
Prof. Dirk Trauner (University of Pennsylvania, USA)  
... and their research groups

