

# Fueling the Digital Chemistry Revolution with Language Models

Teodoro Laino  
IBM Research Europe – Zurich  
NCCR Catalysis



@teodorolaino

# Credits

## Current Members



## Previous Members

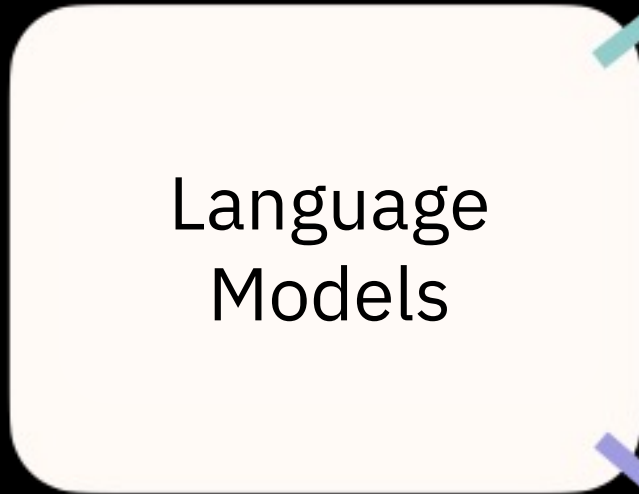


What are language  
models?

Text Input



Language Models



Text Output



Numeric representation of text

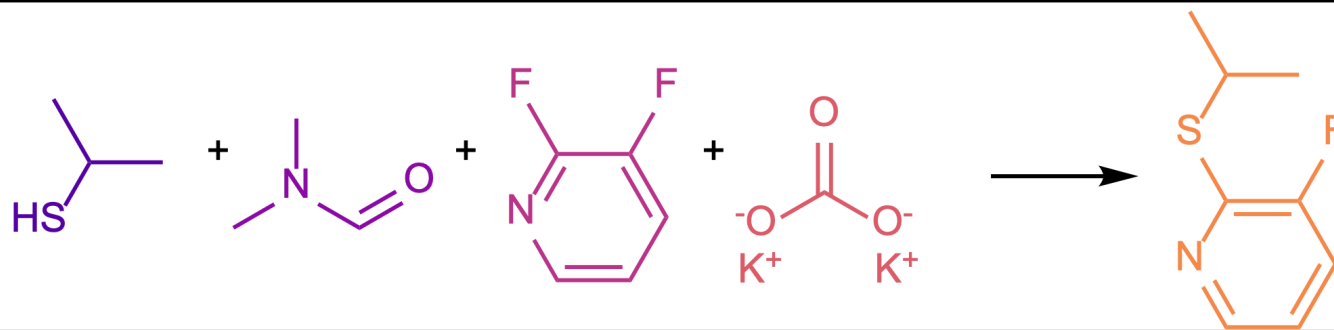
What is a language?

# What is a language?

“A language is a purely human and non-intrinsic method of **communicating** ideas, emotions, and desires by means of **voluntarily produced symbols**” (Sapir, 1921)

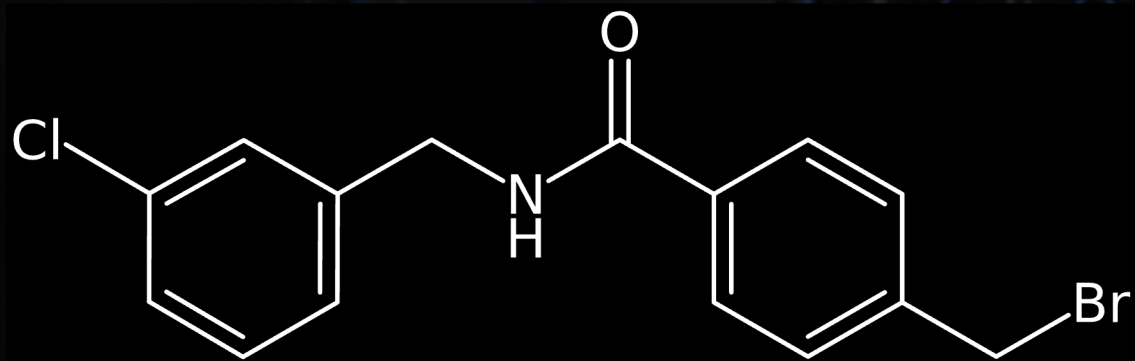
# Language and Chemical Tasks

2.7 g (12.3 mmol) 4,4-Dimethyl-1,2,3,4-tetrahydro-2-oxo-7-quinolinecarboxylic acid were added to a solution of 3.8 g (18.5 mmol) N,N'-dicyclohexylcarbodiimide and 1.1 ml (12.3 mmol) aniline in 80 ml dichloromethane. The reaction mixture was stirred for 4 hours at ambient temperature and the precipitate was filtered off with suction and recrystallised from ethanol. There was obtained 1.2 g of the title compound; m.p. 249-251° C.

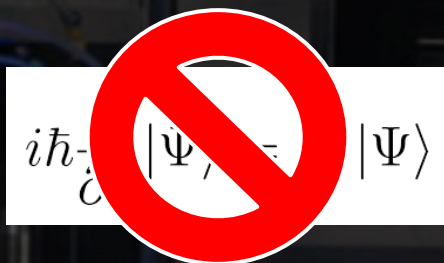


Chemical Synthesis

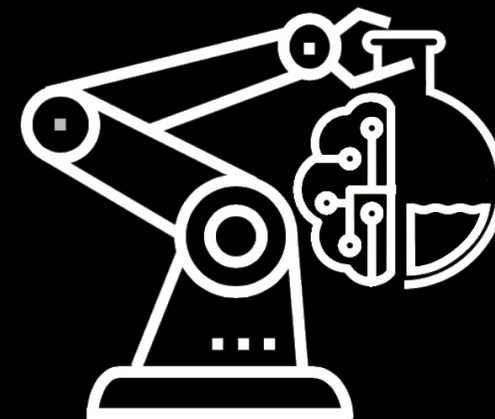
# Data and chemical reactions



Target molecule



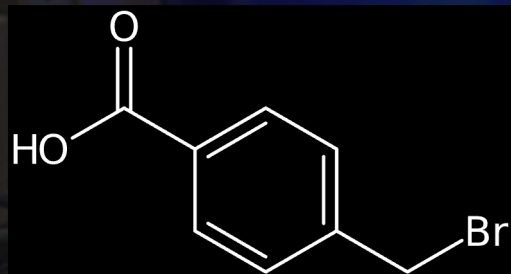
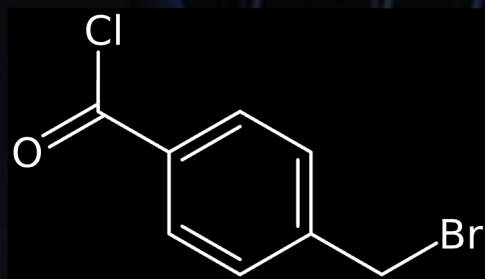
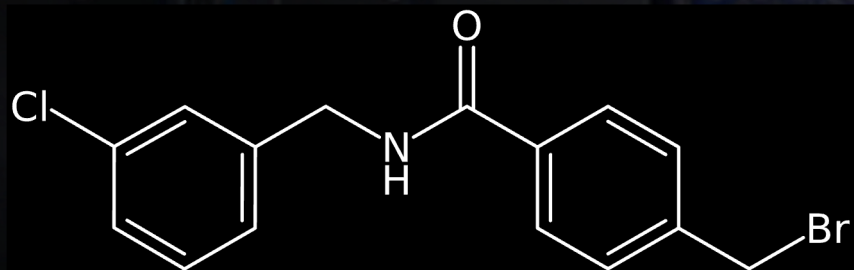
2.7 g (12.3 mmol) 4,4-Dimethyl-1,2,3,4-tetrahydro-2-oxo-7-quinolinecarboxylic acid were added to a solution of 3.8 g (18.5 mmol) N,N'-dicyclohexylcarbodiimide and 1.1 ml (12.3 mmol) aniline in 80 ml dichloromethane. The reaction mixture was stirred for 4 hours at ambient temperature and the precipitate was filtered off with suction and recrystallised from ethanol. There was obtained 1.2 g of the title compound; m.p. 249-251° C.



Synthesis execution



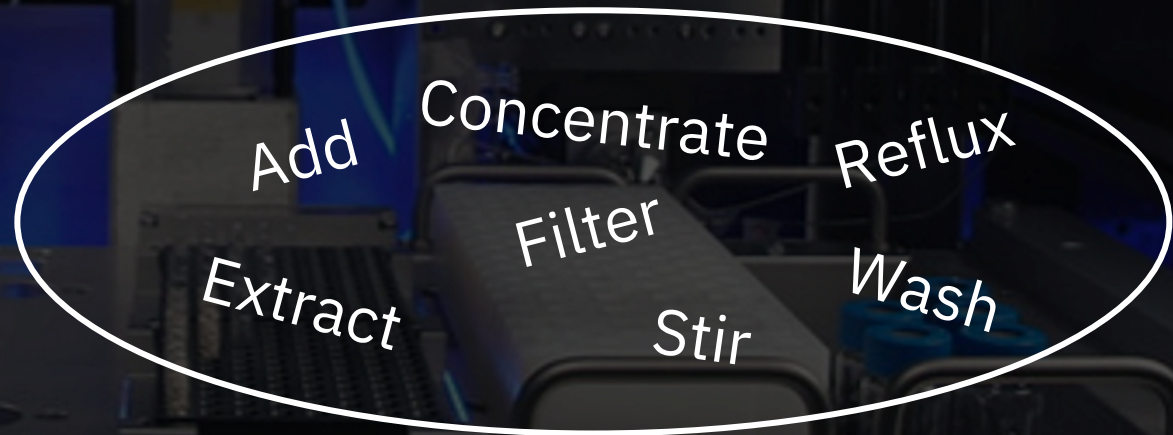
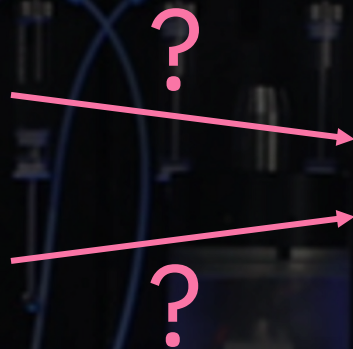
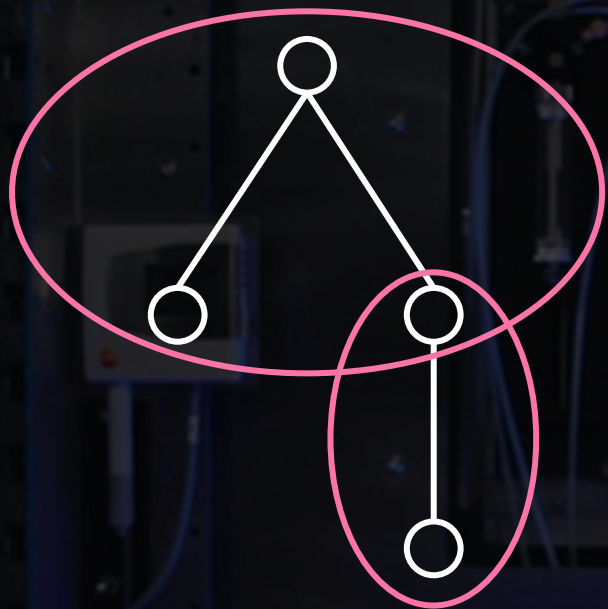
# Synthesis Design



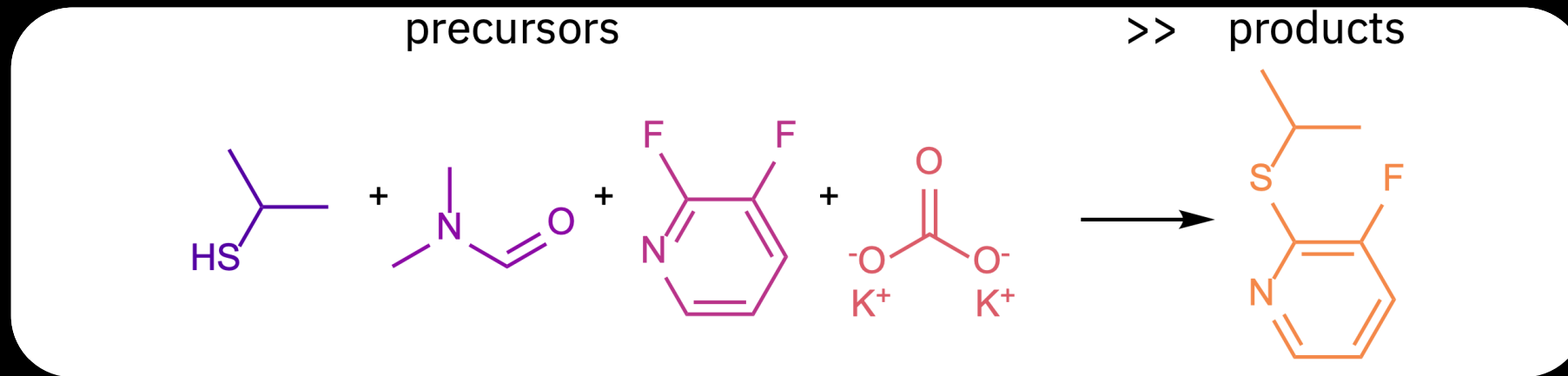
Retrosynthetic tree



# Synthesis Execution



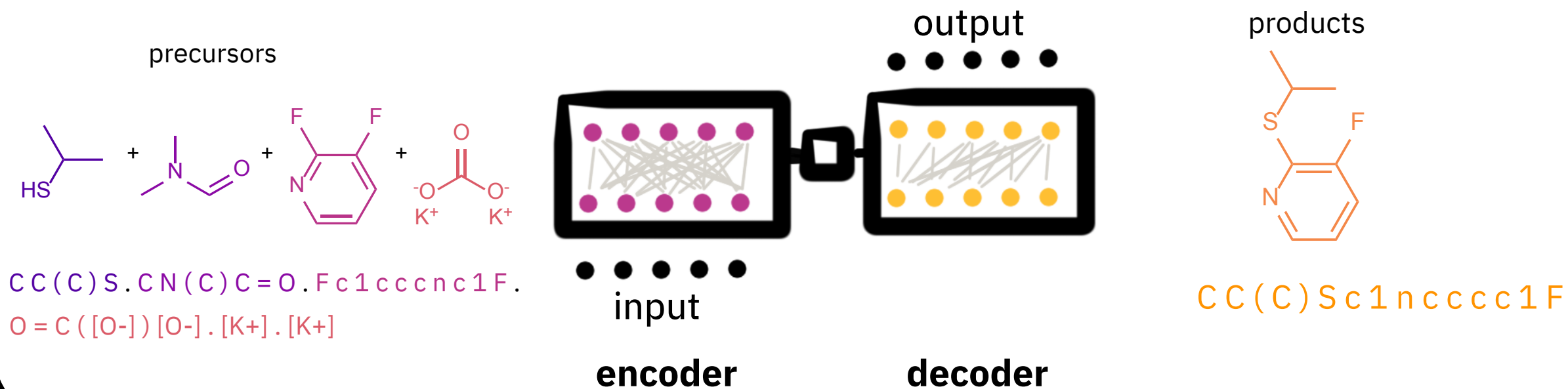
# Atoms as *letters*, molecules as *words*



CC(C)S.CN(C)C=O.Fc1ccnc1F.O=C([O-])[O-].[K+].[K+]>>CC(C)Sc1ncccc1F

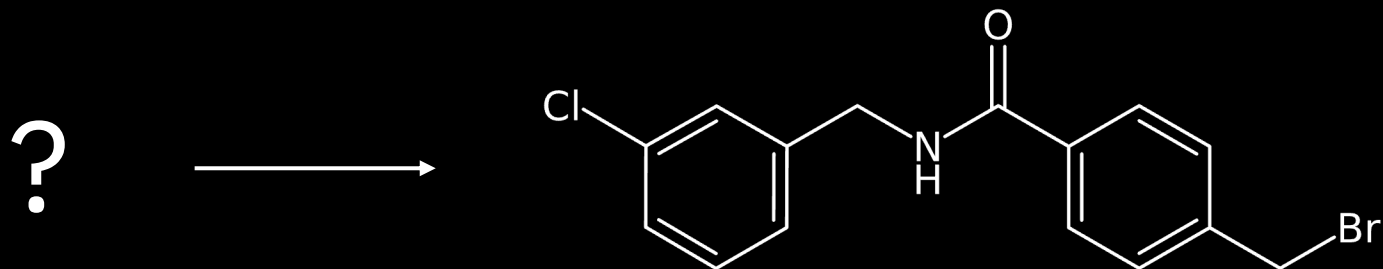
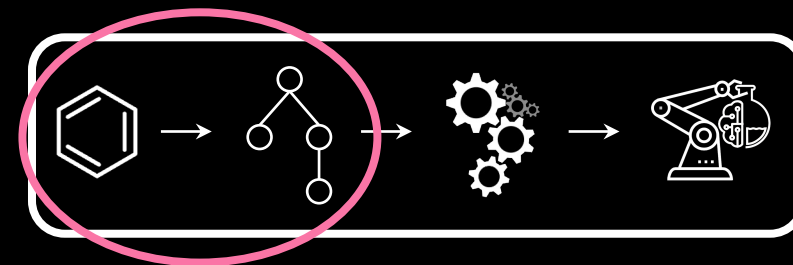
Cast reaction prediction as translation task

# Molecular Transformer

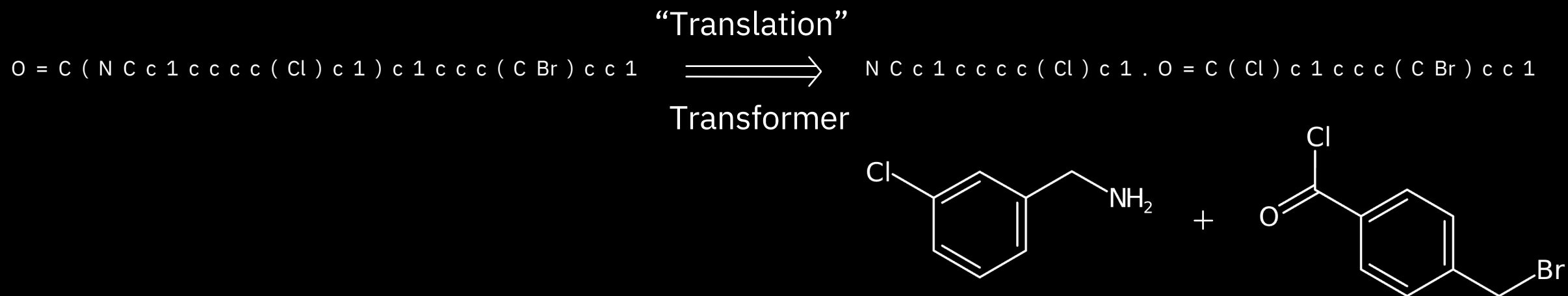


- **No rules** integrated / no chemical knowledge
- **Accurate predictions** on unseen reactions (>90% accuracy on benchmark)
- Better than rule and graph-based approaches

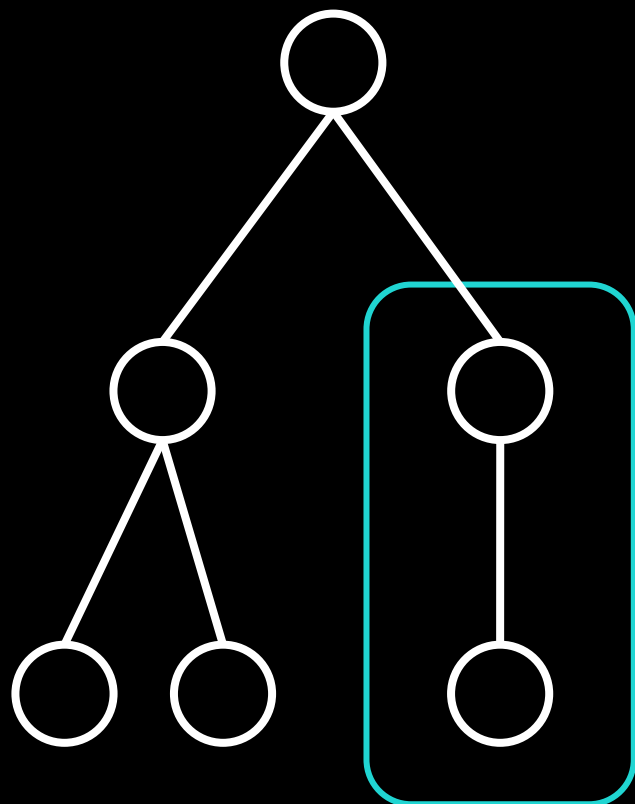
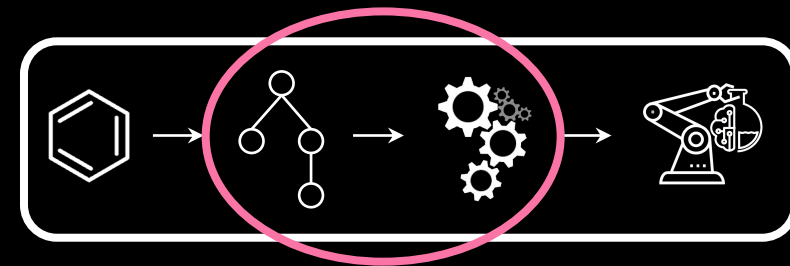
# Synthesis Design



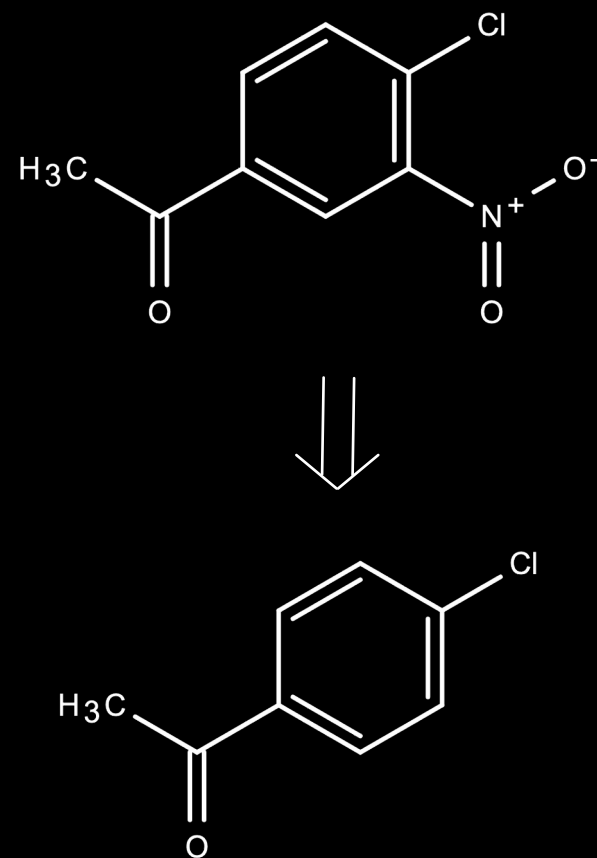
Similar approach, both sides switched



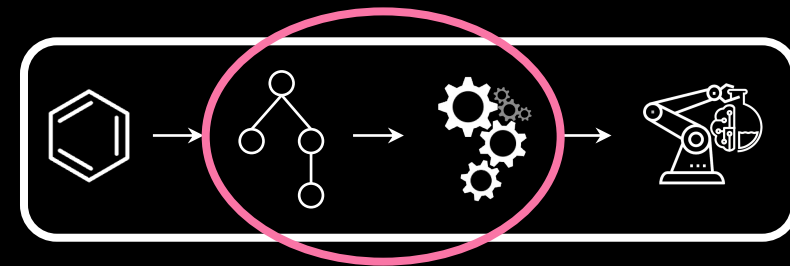
# Synthesis actions



One reaction step



# Building a dataset for ML model



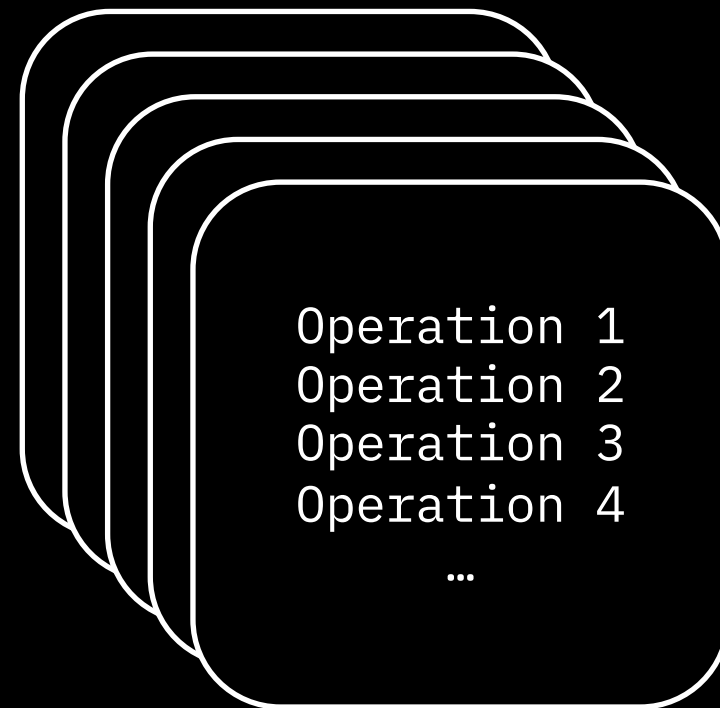
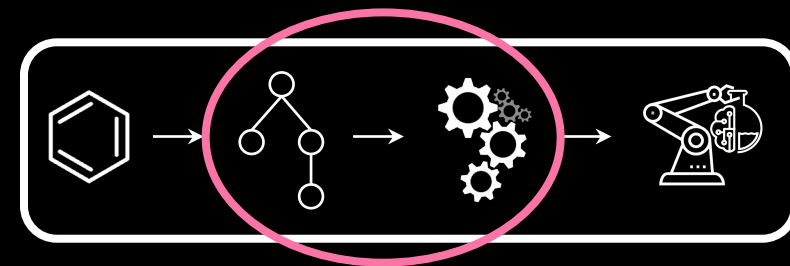
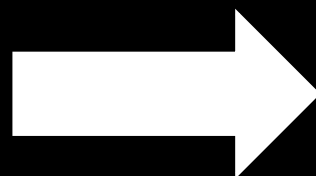
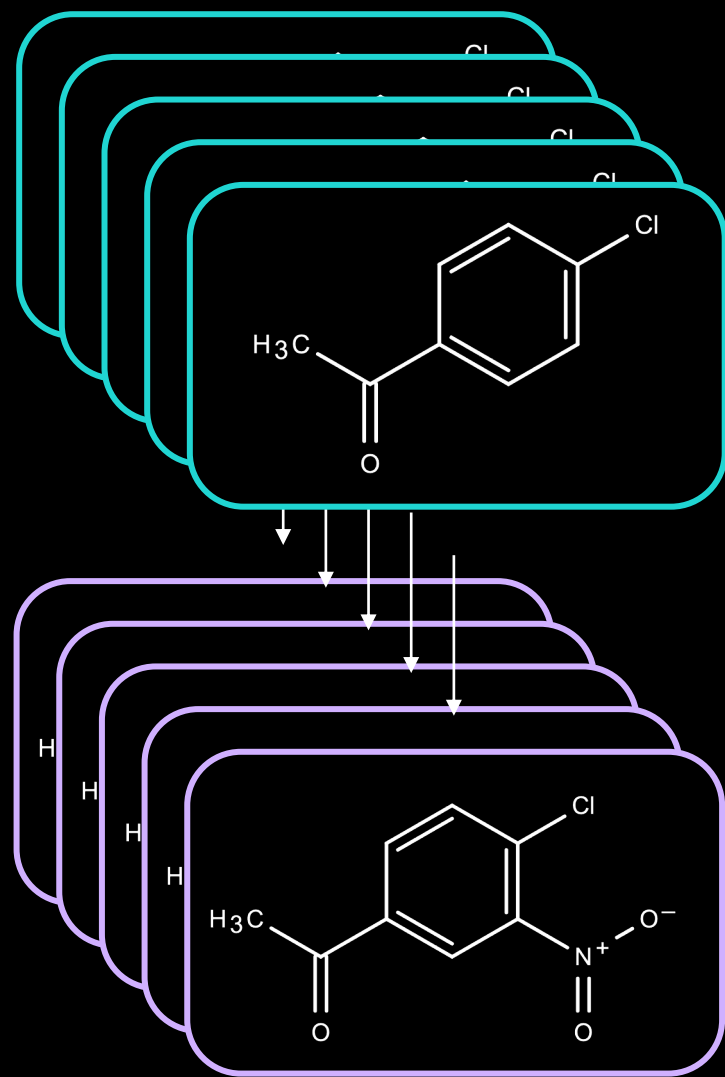
The TFA was removed in vacuo and a saturated solution of NaHCO<sub>3</sub> was added.

Translation



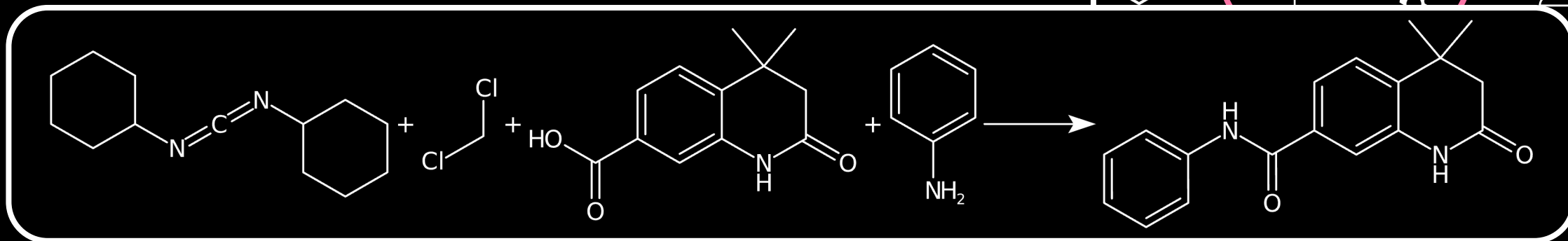
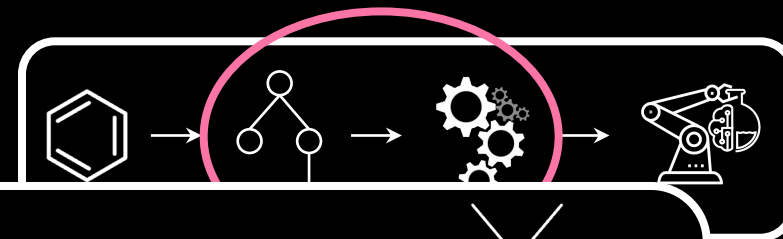
```
Concentrate(),  
Add(name='saturated solution of NaHCO3')
```

# SMILES-to-actions dataset





# SMILES-to-actions



C(=NC1CCCCC1)=NC1CCCCC1 . ClCCl . CC1(C)CC(=O)Nc2cc(C(=O)O)ccc21 . Nc1ccccc1 >> CC1(C)CC(=O)Nc2cc(C(=O)Nc3ccccc3)ccc21

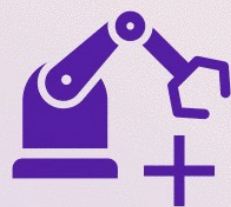
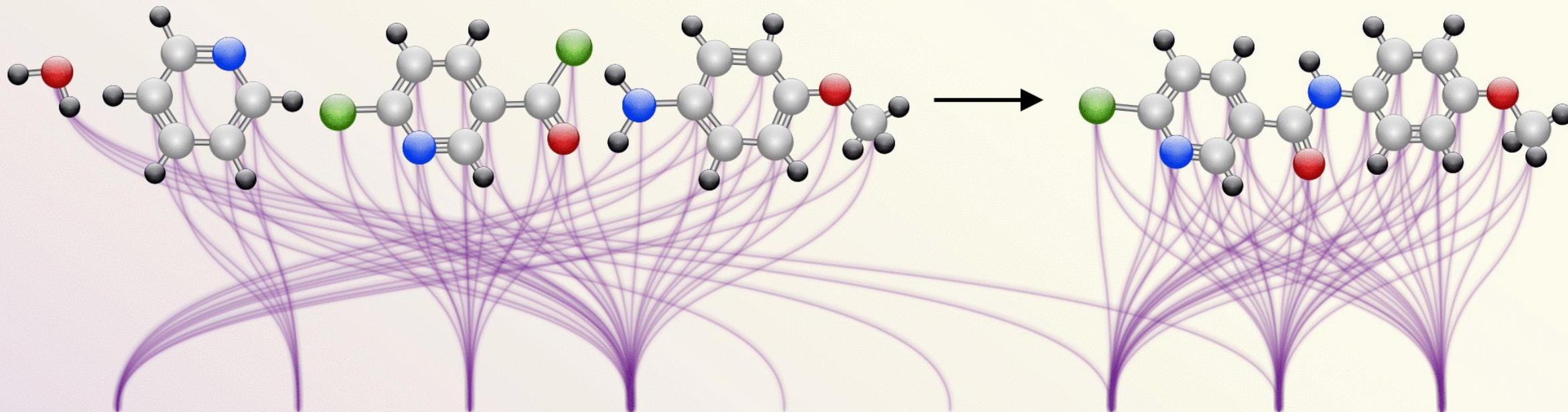
2.7 g (12.3 mmol) 4,4-Dimethyl-1,2,3,4-tetrahydro-2-oxo-7-quinolinecarboxylic acid were added to a solution of 3.8 g (18.5 mmol) N,N'-dicyclohexylcarbodiimide and 1.1 ml (12.3 mmol) aniline in 80 ml dichloromethane. The reaction mixture was stirred for 4 hours at ambient temperature and the precipitate was filtered off with suction and recrystallised from ethanol. There was obtained 1.2 g of the title compound; m.p. 249-251° C.

ML model

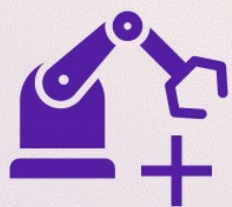
1. MAKESOLUTION with N,N'-dicyclohexylcarbodiimide (3.8 g, 18.5 mmol) and aniline (1.1 ml, 12.3 mmol) and dichloromethane (80 ml)
2. STIR for 4 hours at ambient temperature
3. ADD 4,4-Dimethyl-1,2,3,4-tetrahydro-2-oxo-7-quinolinecarboxylic acid (2.7 g, 12.3 mmol)
4. STIR for 4 hours at ambient temperature
5. FILTER keep precipitate
6. RECRYSTALLIZE from ethanol
7. YIELD title compound (1.2 g)

1. ADD \$1\$
2. ADD \$4\$
3. ADD \$2\$
4. ADD \$3\$
5. STIR for @3@ at #4#
6. FILTER keep precipitate
7. RECRYSTALLIZE from ethanol
8. YIELD \$-1\$

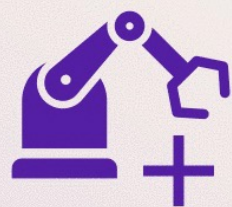
O . c1ccncc1 . O=C(Cl)c1ccc(Cl)nc1 . COc1ccc(N)cc1 >> COc1ccc(NC(=O)c2ccc(Cl)nc2)cc1



Add



Add



Add



Stir



Add



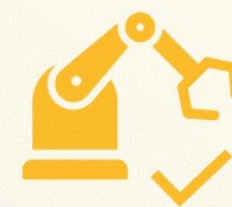
Filter



Wash

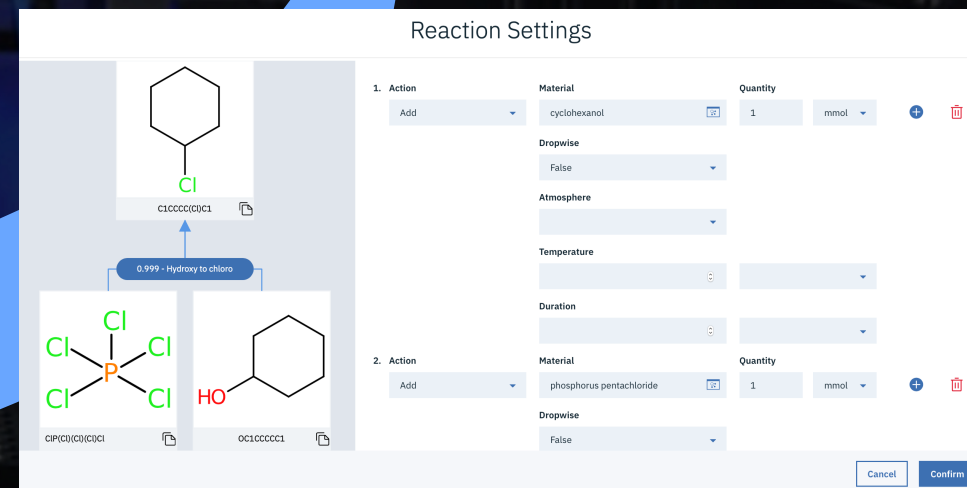
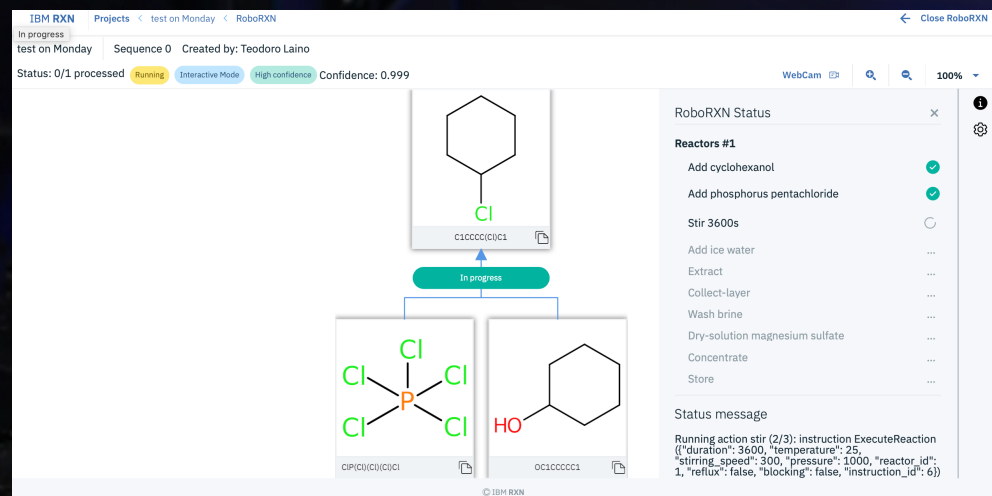
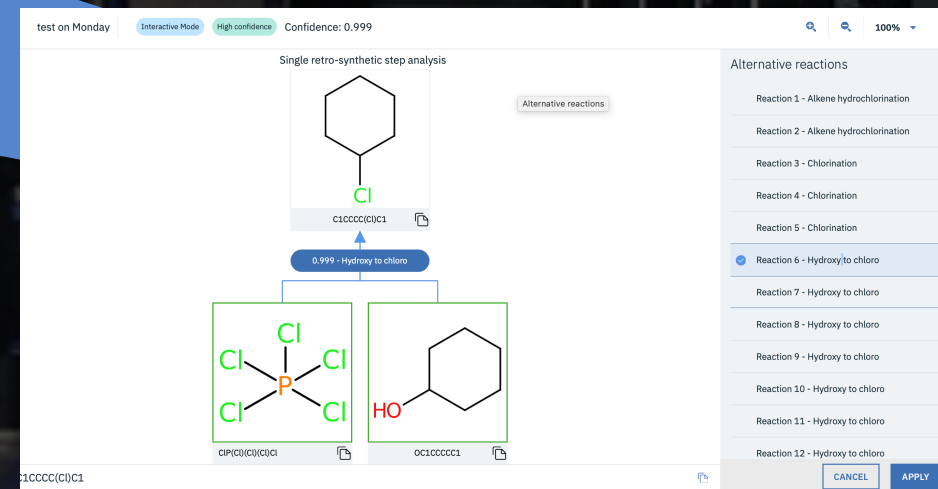
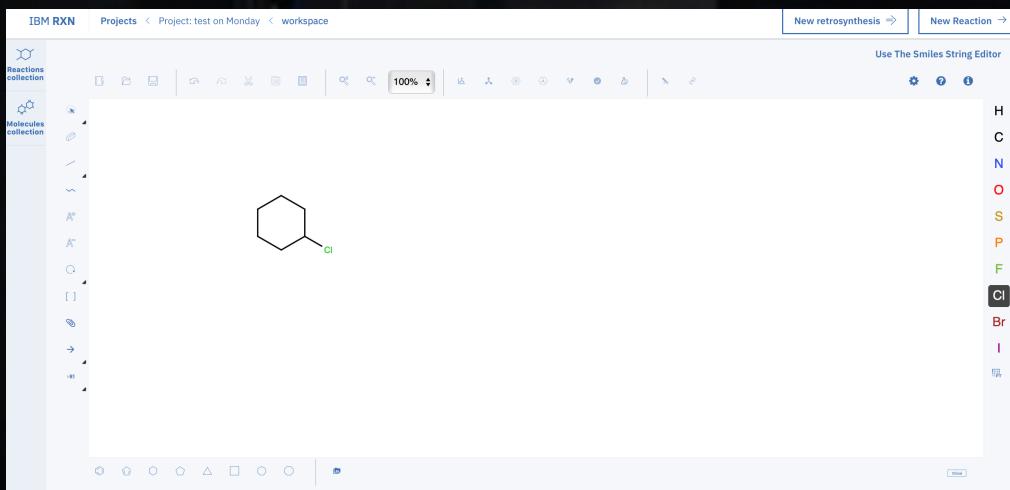


Dry

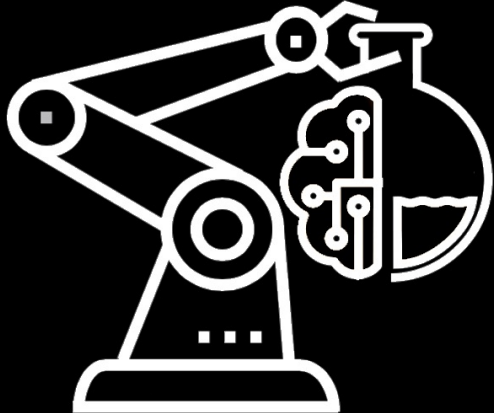
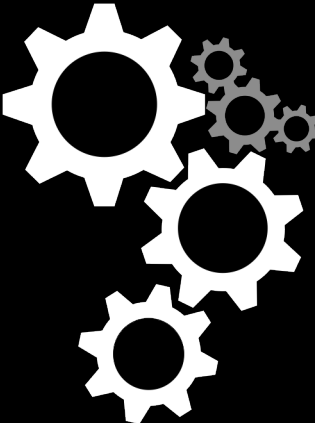
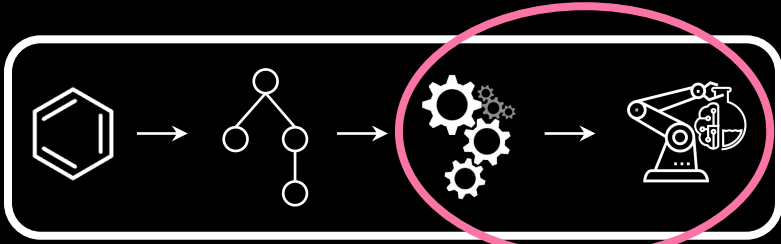


Yield

**SMILES-2-ACTIONS**



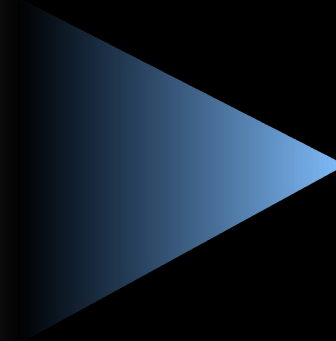
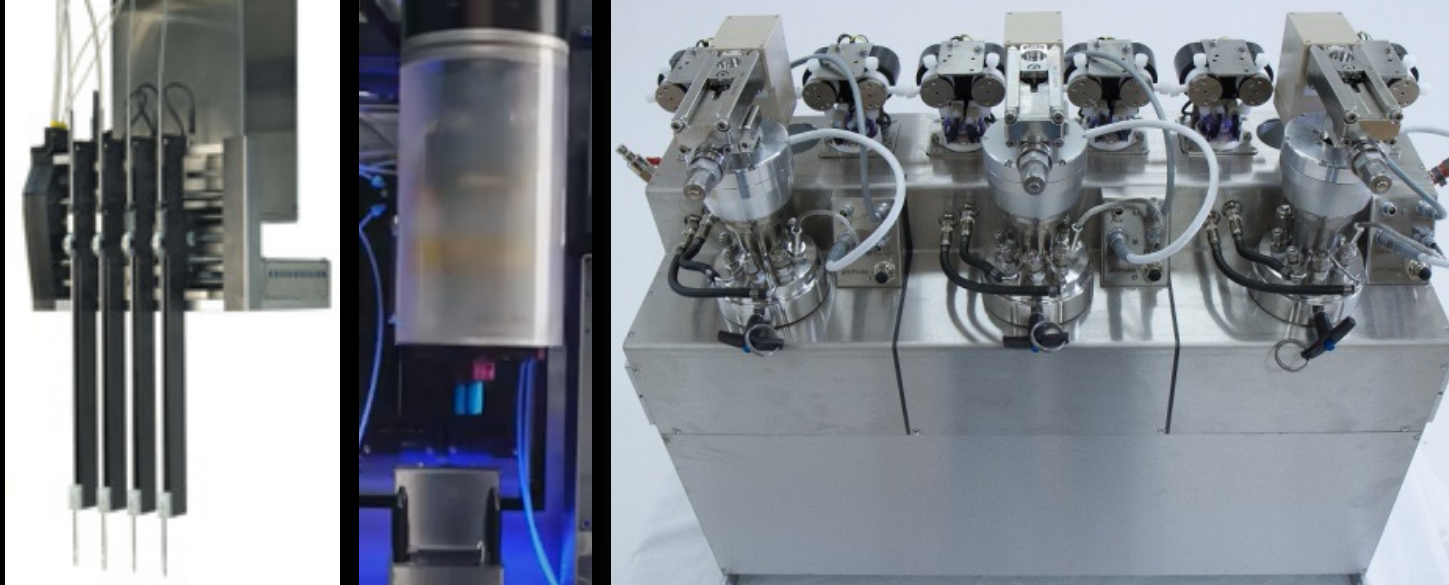
# Execution on chemical robot



Add Concentrate Reflux  
Extract Filter Stir Wash

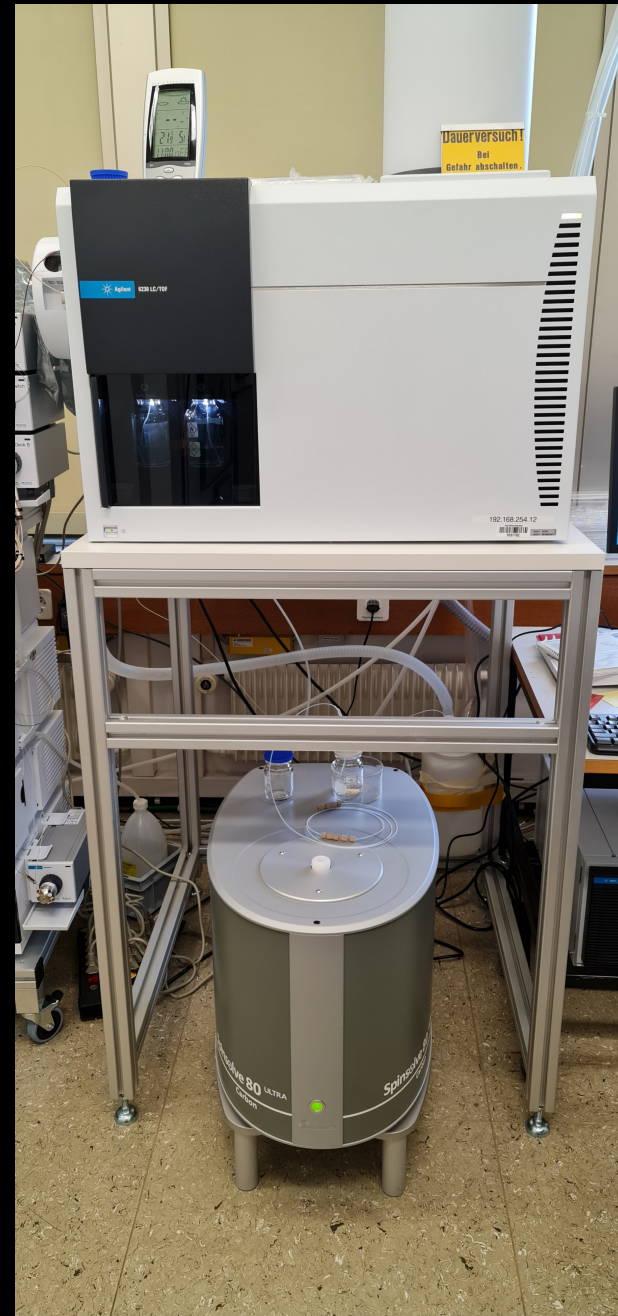
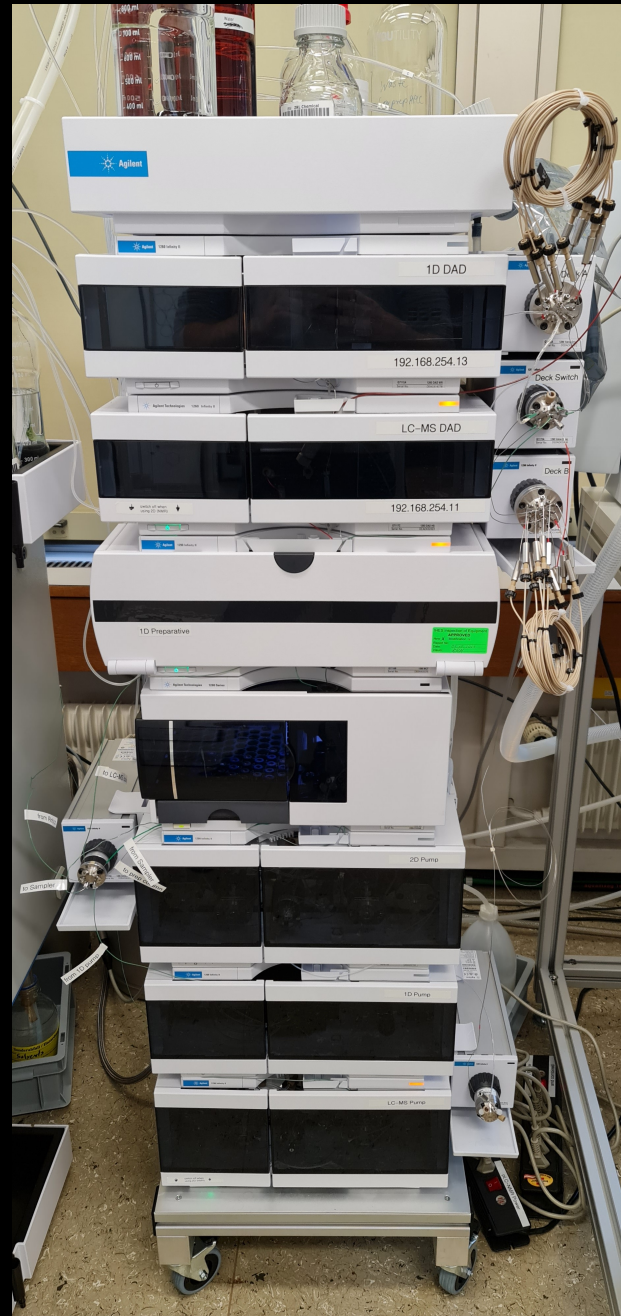
Start heating reactor #2  
Transfer 3 ml of solvent #5 to reactor #2  
Set pressure

# Flex Autoplant

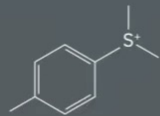


FLEX AUTOPLANT robotic platform

# Analytics



## Synthesizing new molecule



Started: Nov 30 2020, 6:49am PT

Live from IBM RoboRXN

Action 2

Overview

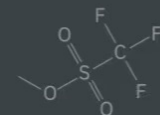
### Adding $C_2H_3F_3O_3S$ +

In this action, the molecule methyl trifluoromethane sulfonate is added to [Reactor 2](#).

Methyl trifluoromethane

$C_2H_3F_3O_3S$

2D 3D



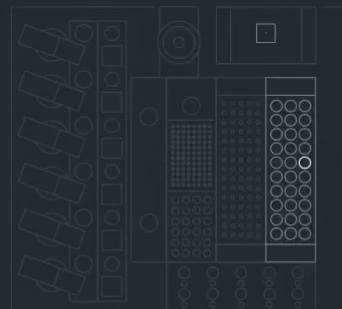
Methyl trifluoromethane sulfonate is a brown liquid. Insoluble in water. This material is a very reactive methylating agent, also known as methyl triflate.

**NOW**

10 ml of [reagent](#) containing methyl trifluoromethane sulfonate is being moved from [Vial 61](#) and added to [Reactor 2](#).

Position of the robot arm

Moving to Vial 61



Live view module

2 Adding  $C_2H_3F_3O_3S$

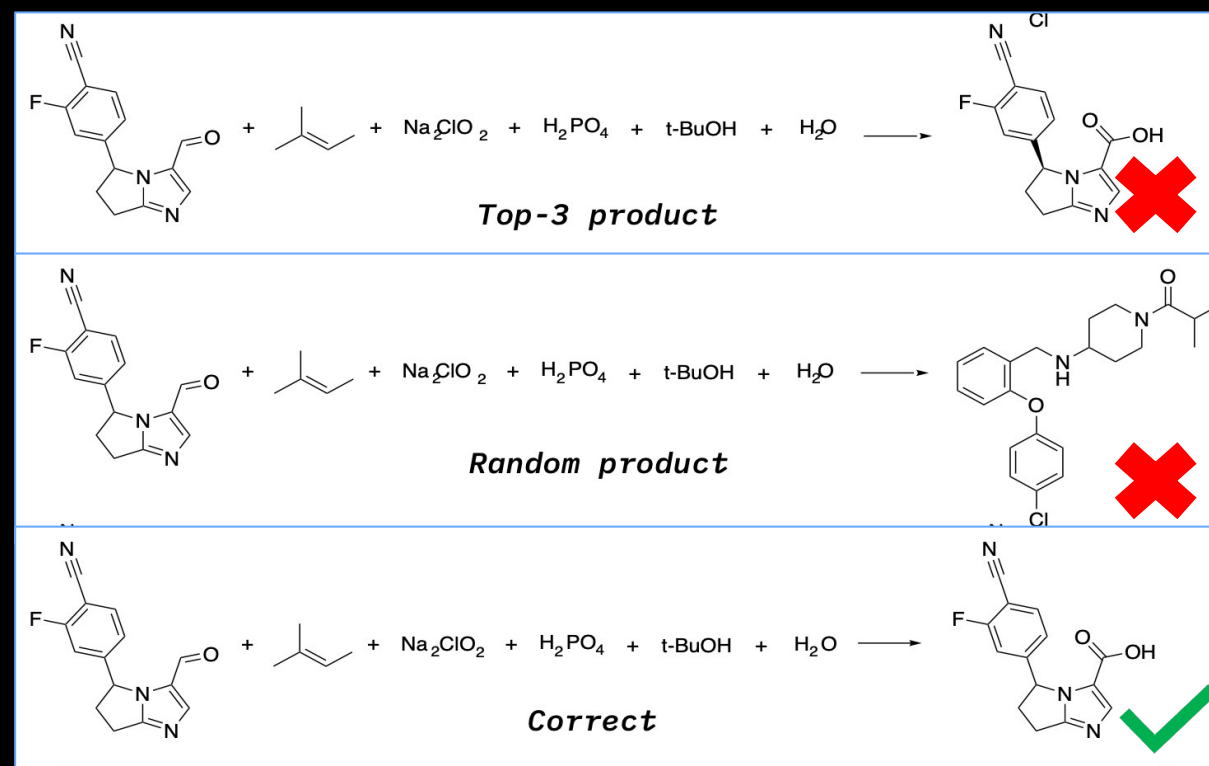
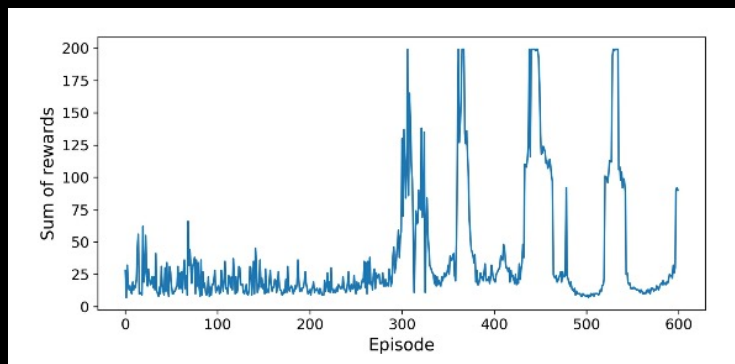
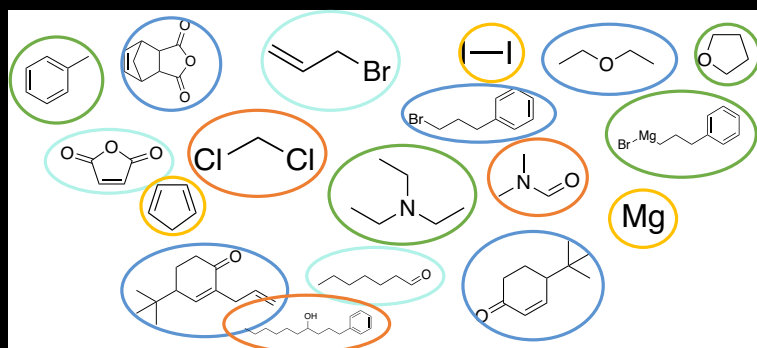
00:06:00



LIVE

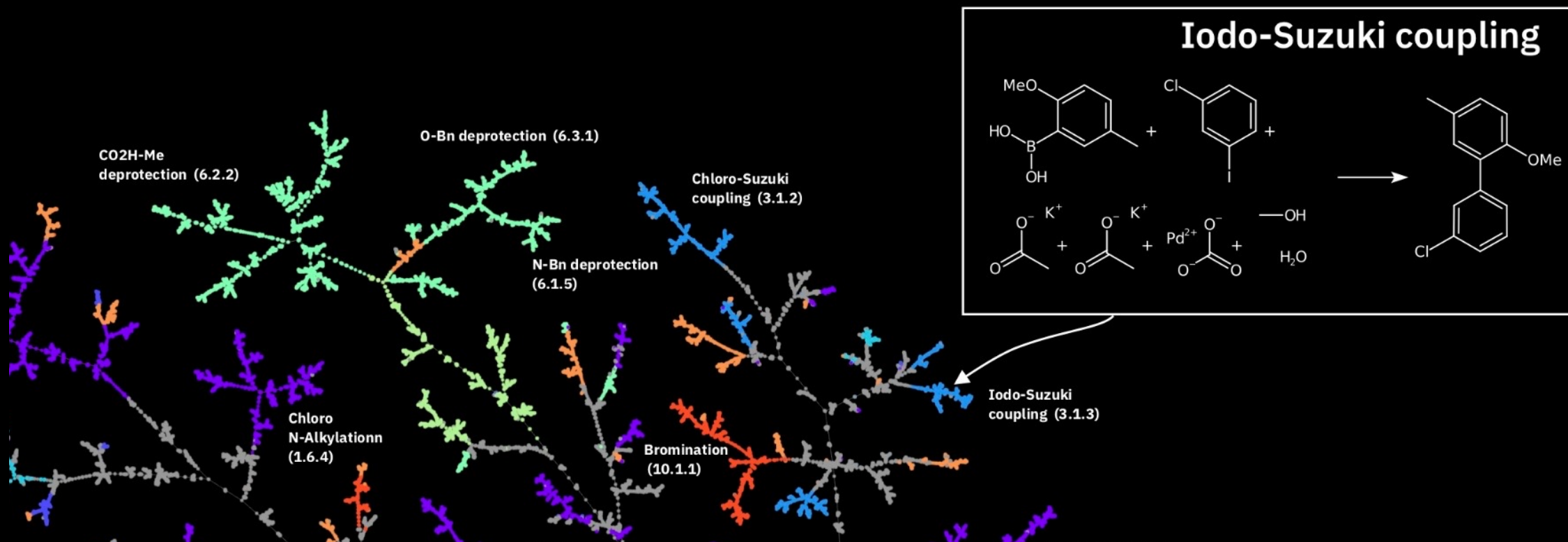
# Unassisted *data sets curation*

The most **difficult examples to learn** during training of reactions prediction models = likely **examples of wrong chemistry**.

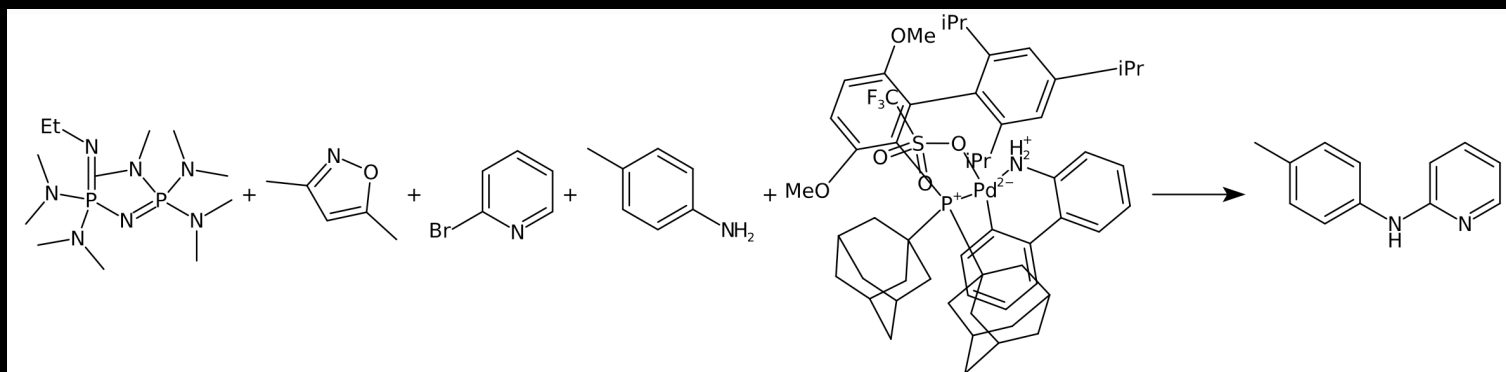




# Reaction Classification, atlases and fingerprints



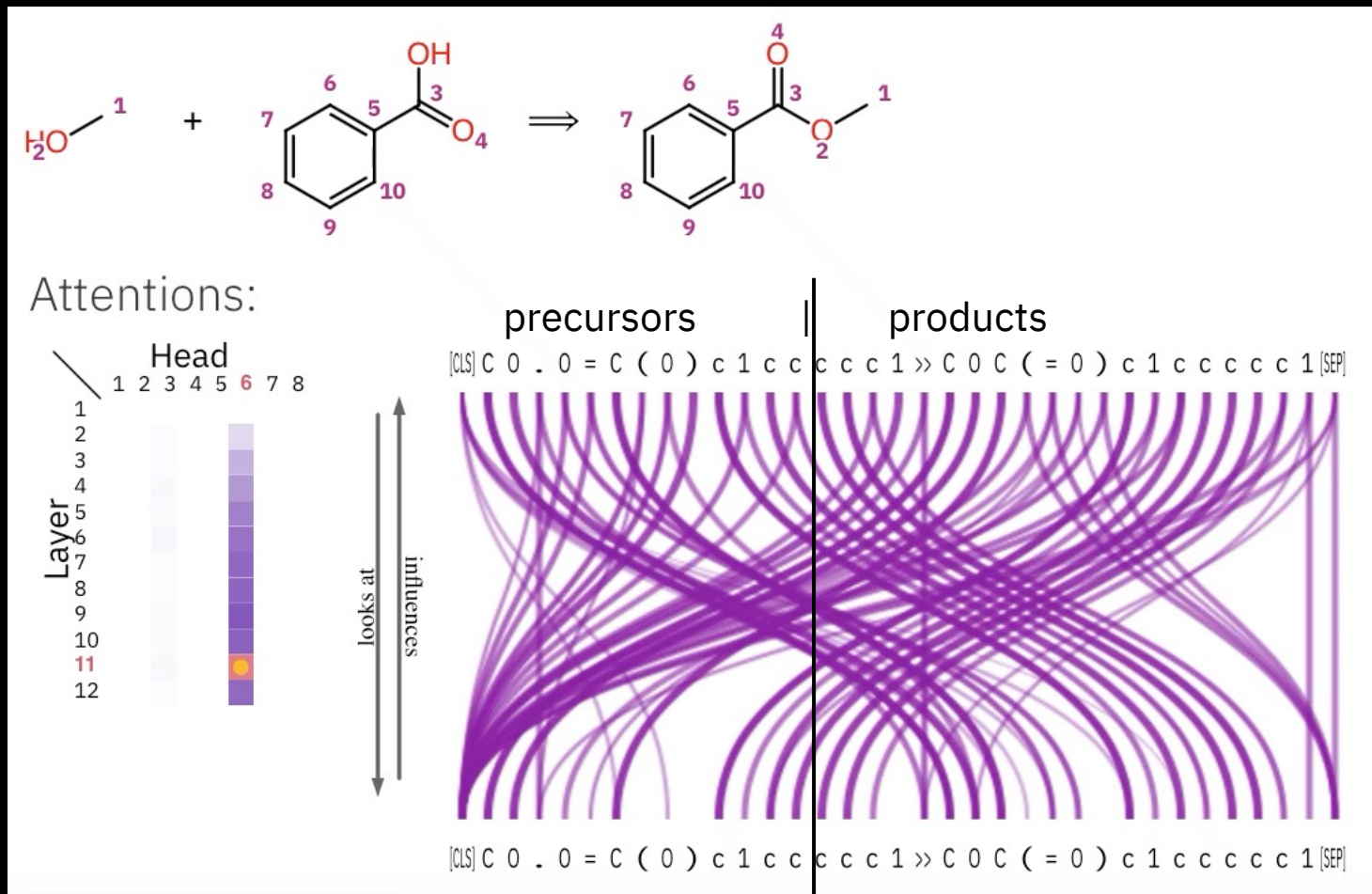
# Prediction of *reaction yields*



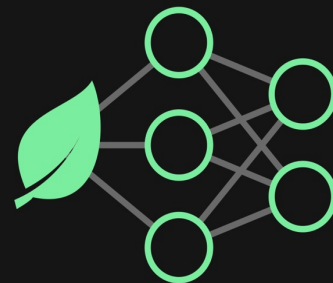
77.6 %

Brc1cccn1.CCN=P(N=P(N(C)C)(N(C)C)N(C)C)(N(C)C)N(C)C.COc1ccc(OC)c([P+](C23CC4CC(CC(C4)C2)C3)(C23CC4C(CC(C4)C2)C3)[Pd-2]2(OS(=O)(=O)C(F)(F)F)[NH2+]c3ccccc3-c3ccccc32)c1-c1c(C(C)C)cc(C(C)C)cc1C(C)C.Cc1cc(C)on1.Cc1ccc(N)cc1>>Cc1ccc(Nc2cccn2)cc1

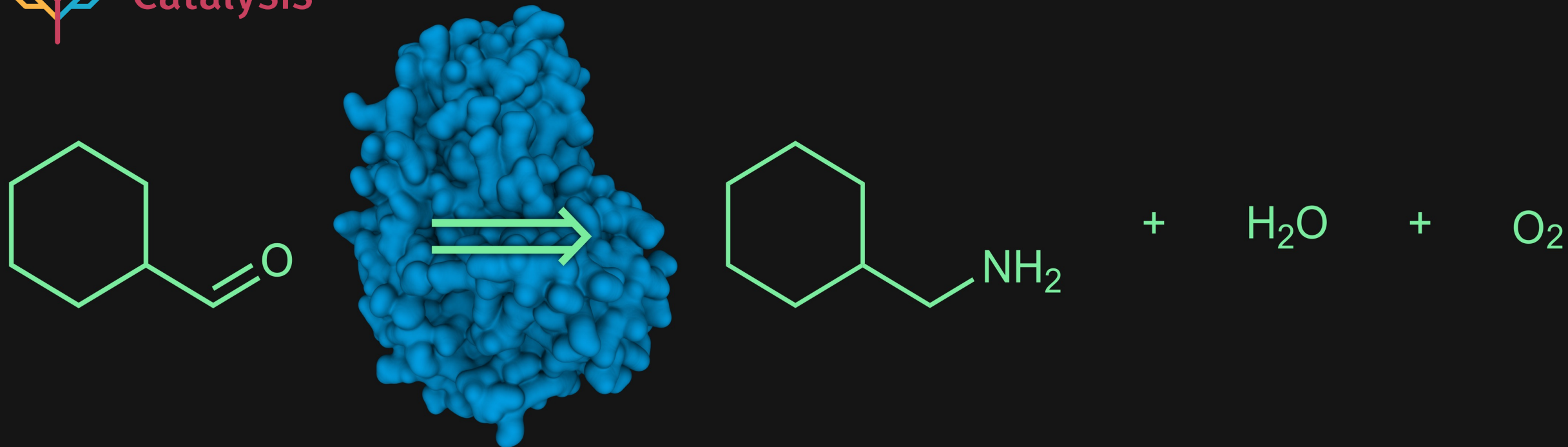
# Atom-Mapping and the learning of chemical reaction grammar (RXNMapper)



# Enzymatic catalysis



## GreenCatRXN



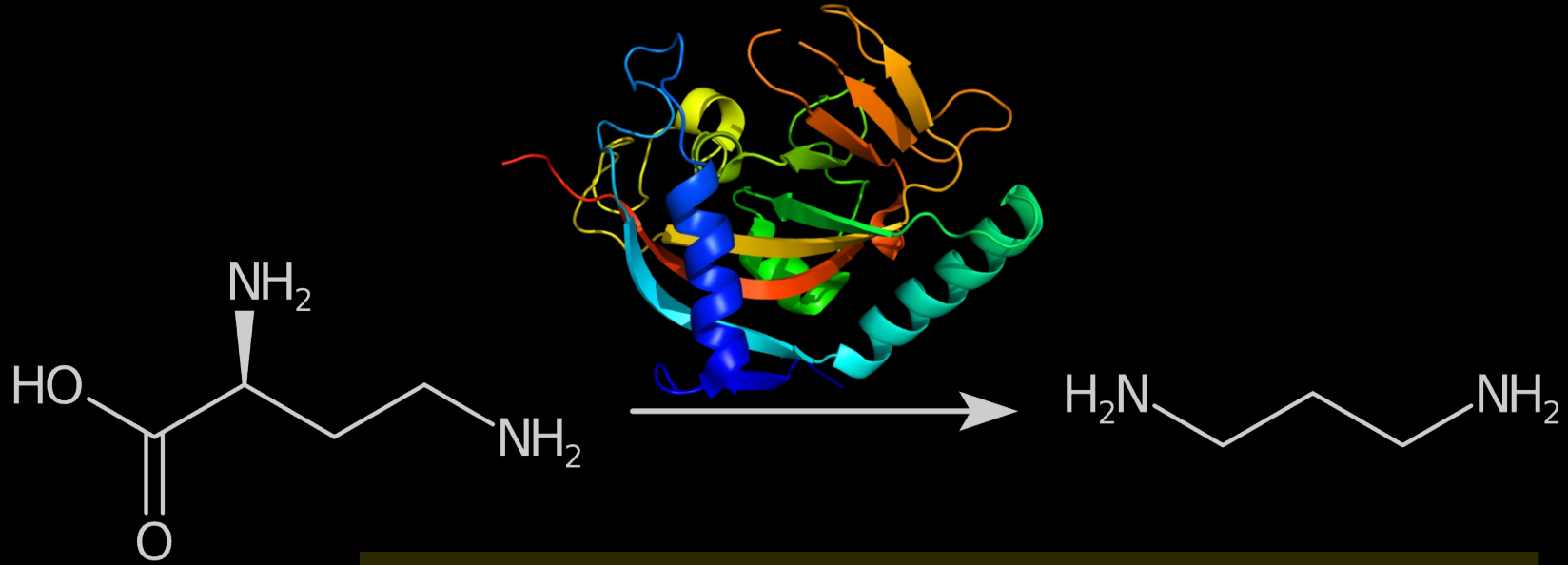
EC 1.4.3.-

~~TEMPO~~

~~NaClO~~

~~DCM~~

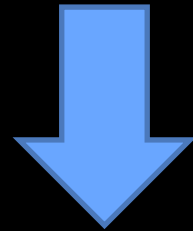
# Recovering active sites (3D) from sequential data (1D)



NCC[C@H](N)C(=O)O

KTYF...SQT SQIHKKNHIRGQARFCP...YVLK

NCCCN



KTYFVL...LLELNENDPGIFVTQSVHKQQAGFSQT SQIHKKNHIRGQARFCPHKR...GYVLK

# References

Chem. Sci., 2018, 9, 6091-6098  
ACS Cent. Sci. 2019, 5, 9, 1572-1583  
Chem. Sci., 2020, 11, 3316-3325  
Nat. Commun., 2020, 11, 3601  
Nat. Commun., 2020, 11, 4874  
Nat. Mach. Intell., 2021, 3, 144–152  
Nat. Mach. Intell. 2021, 3, 485–494  
Adv. Science, 2021, 7, 15, eabe4166  
Mach. Learn.: Sci. Technol., 2021, 2, 015016  
Nat. Commun., 2021, 12, 2573  
Nat. Commun. 2022, 13, 964

Watch the story of RoboRXN (short): <https://youtu.be/ewE1wh7sTUE>

Watch the story of RoboRXN (long): <https://youtu.be/i2-LgHjgDTs>

More information and access/test: <https://rxn.res.ibm.com>

Collaborators:



NCCR  
**Catalysis**

IBM