# Fueling the Digital Chemistry Revolution with Language and Multimodal Foundation Models
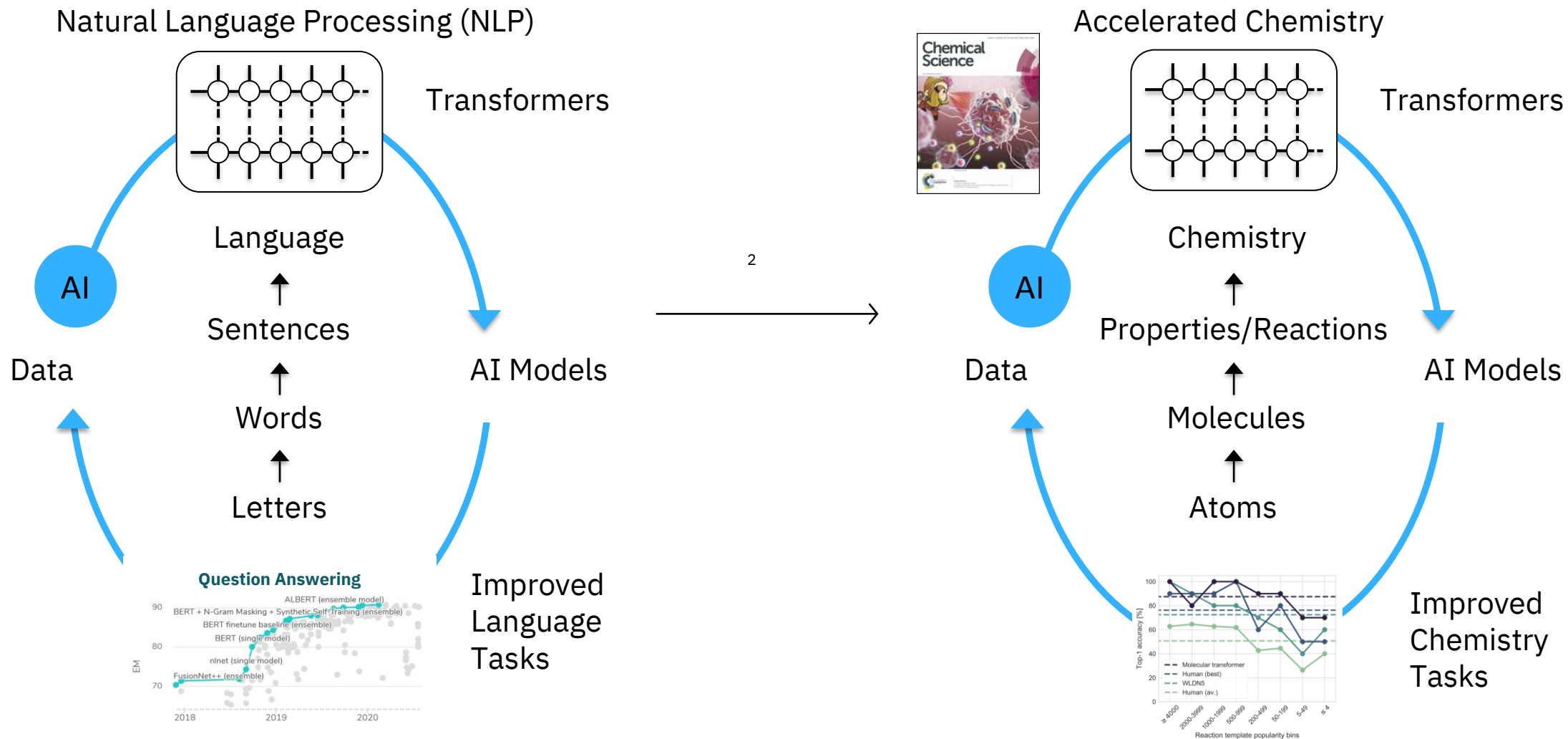
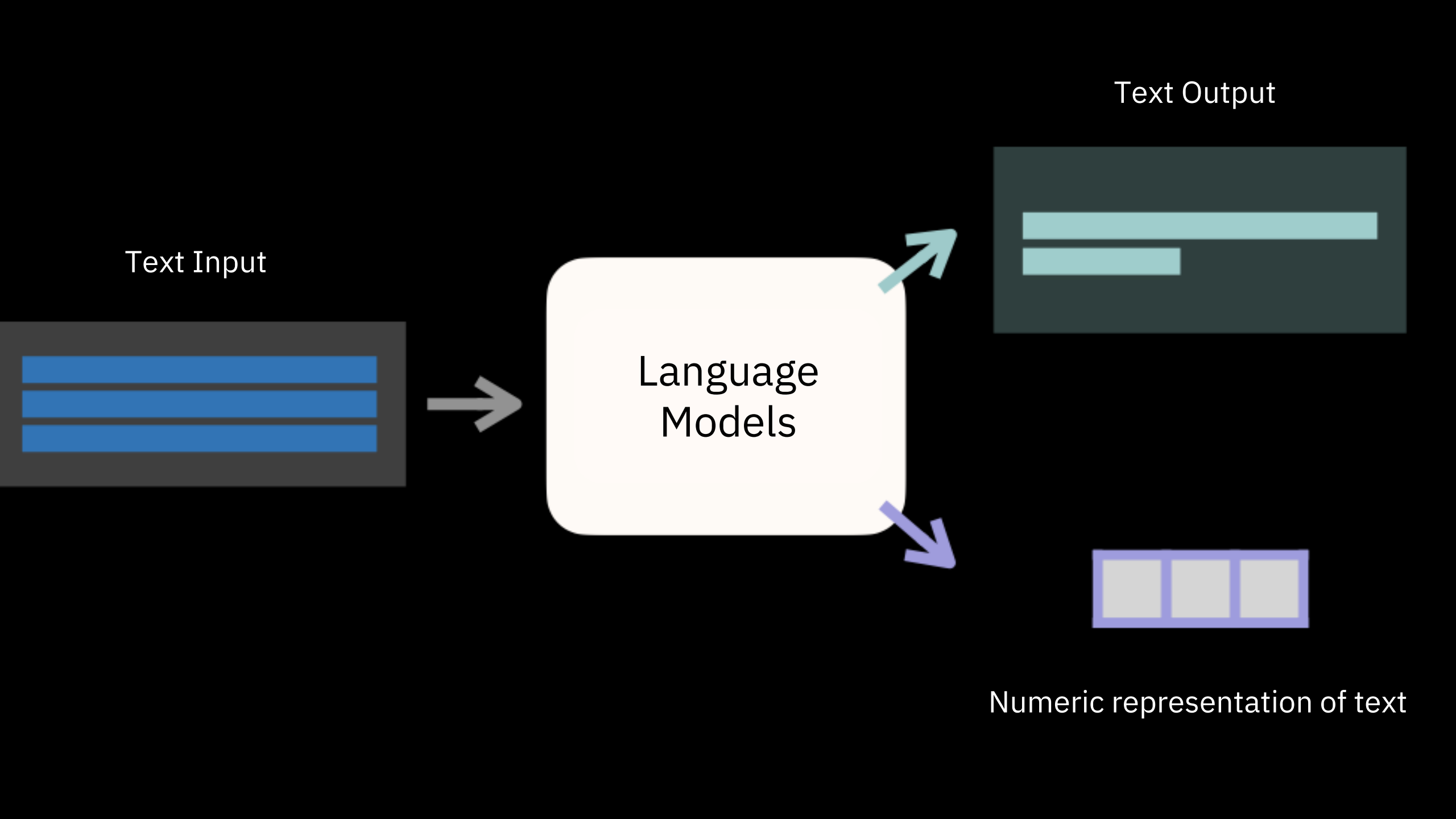Teodoro Laino
IBM Research Europe – Zurich

@teodorolaino
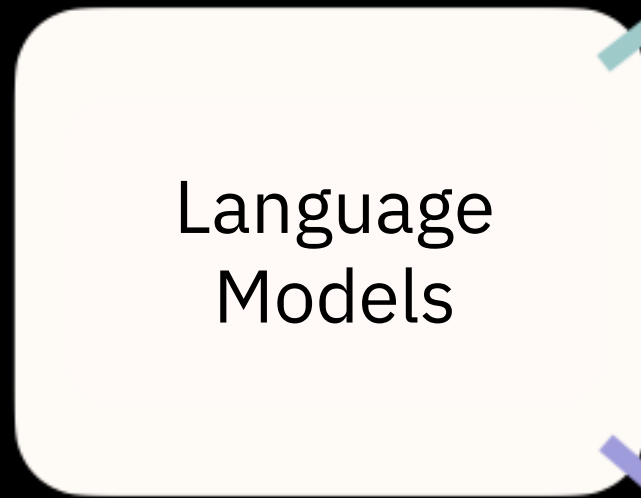
# AI breakthroughs for language are changing scientific discovery

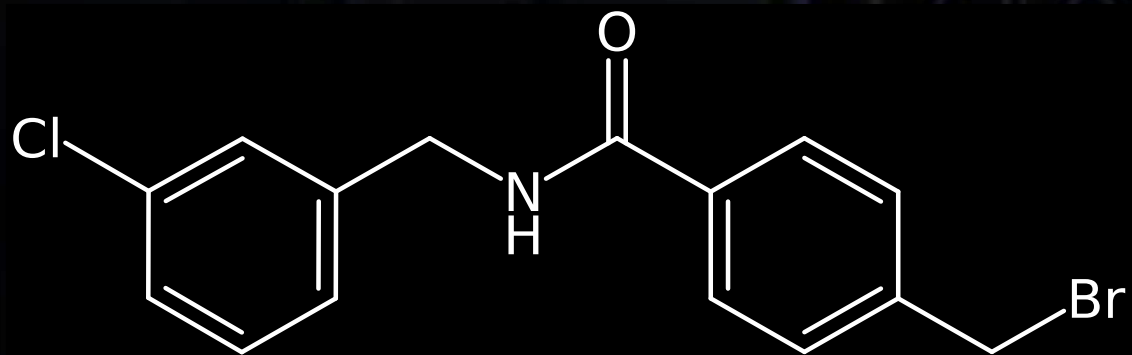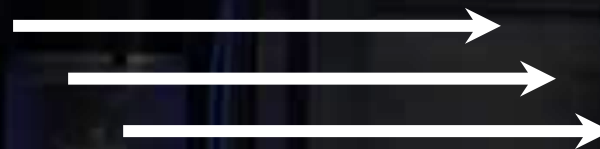Generative modeling and transformers are achieving new breakthroughs in chemistry

# 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.
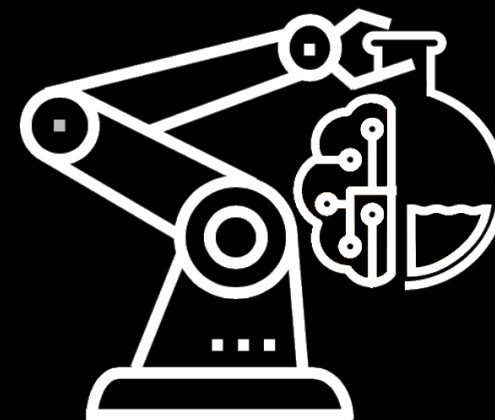
**Synthesis execution**
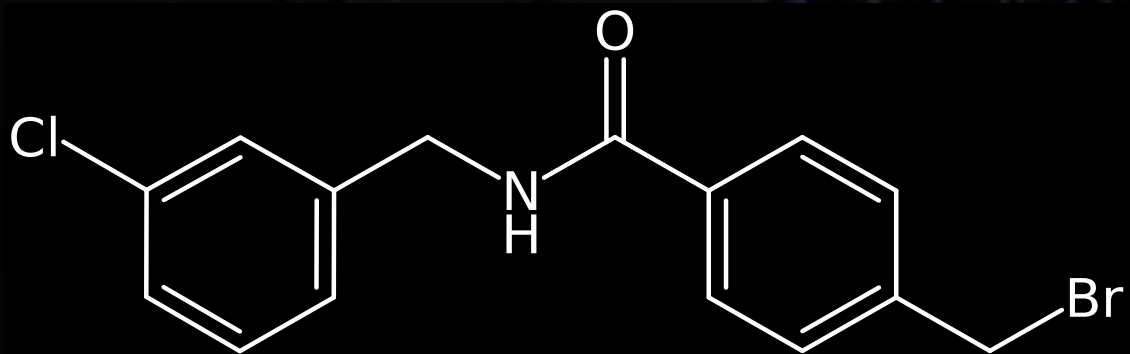
# Data and chemical reactions



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.
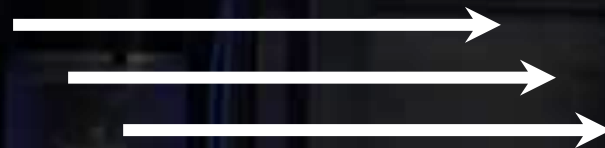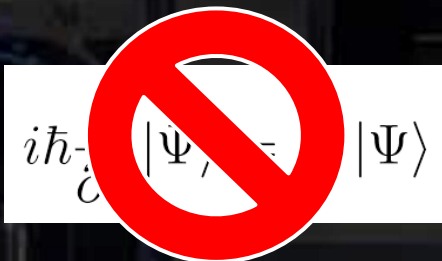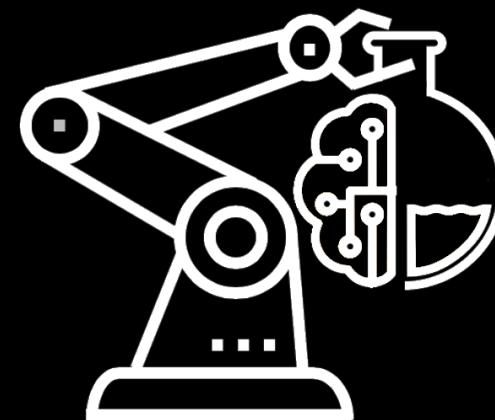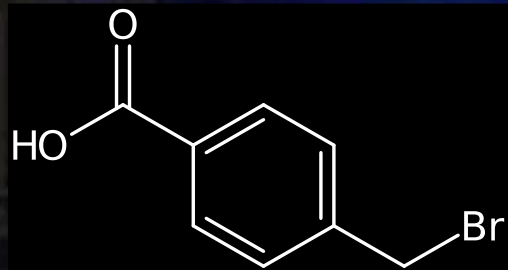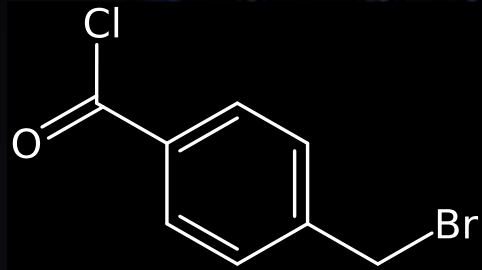
Synthesis execution

# Synthesis Design



Retrosynthetic tree

# Synthesis Execution

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



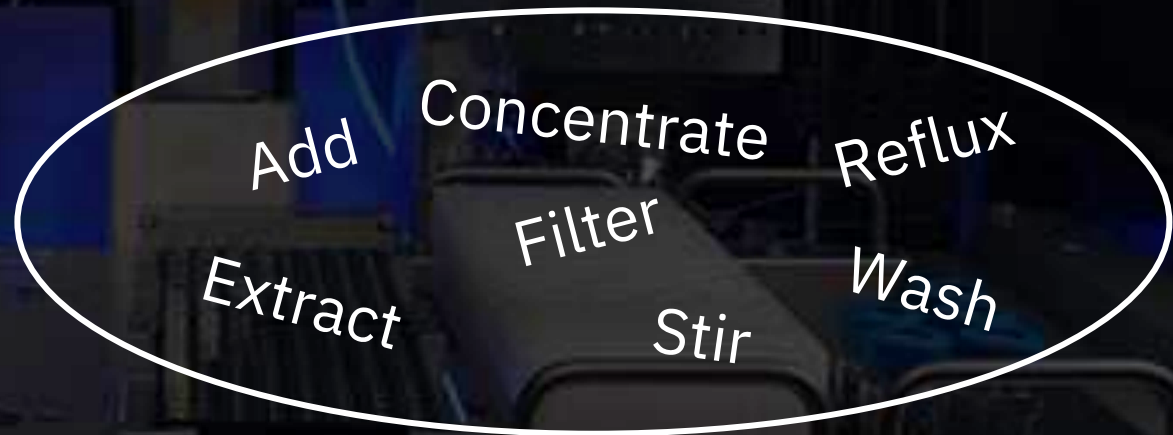Cast reaction prediction as translation task

# *Molecular* Transformer



precursors

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

input

output

encoder  decoder

products

CC(C)Sc1nccccc1F

- **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

"Translation"

O = C ( N C c 1 c c c c ( Cl ) c 1 ) c 1 c c c ( C Br ) c c 1  ⟹  N C c 1 c c c c ( Cl ) c 1 . O = C ( Cl ) c 1 c c c ( C Br ) c c 1

Transformer

+

Chem. Sci., 2020, 11, 3316-3325

# Synthesis actions



One reaction step

# Building a dataset for ML model



The TFA was removed in vacuo and a saturated solution of NaHCO3 was added.

**Translation** →

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

# SMILES-to-actions dataset



```
Operation 1
Operation 2
Operation 3
Operation 4
    …
```

# 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

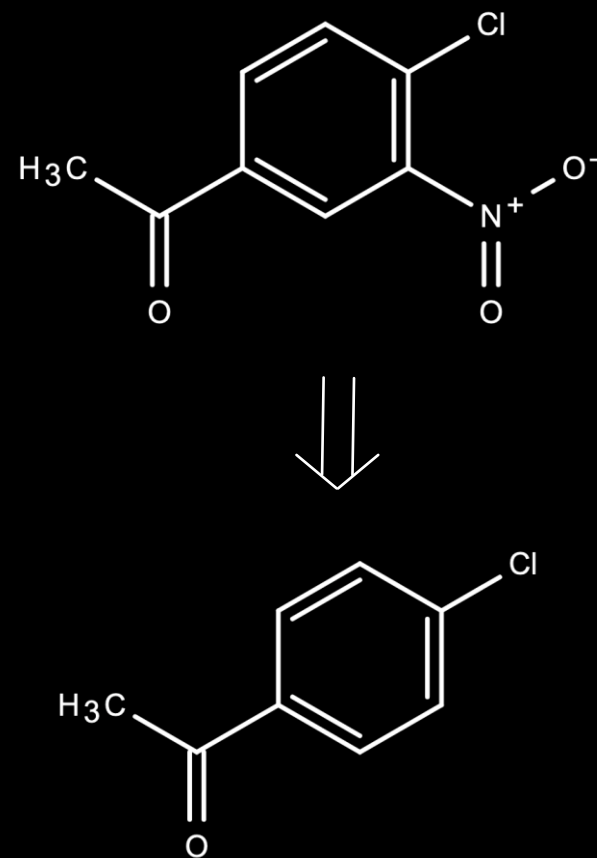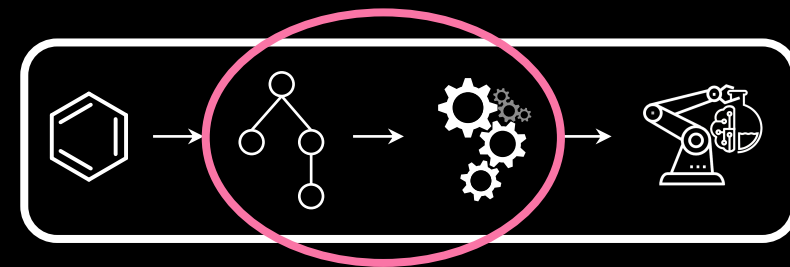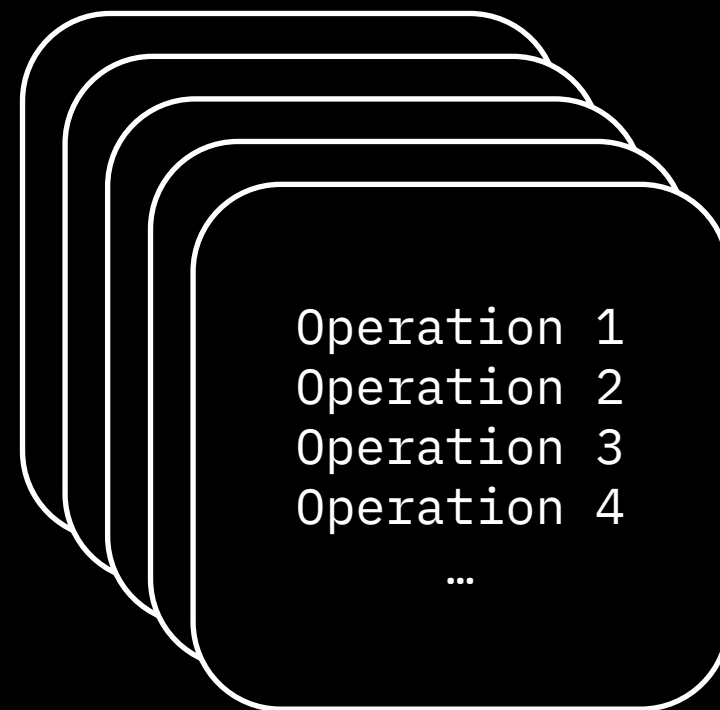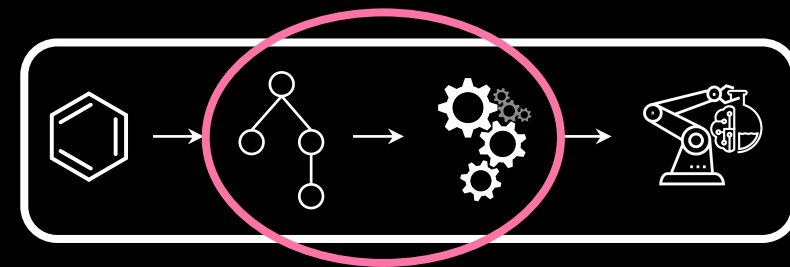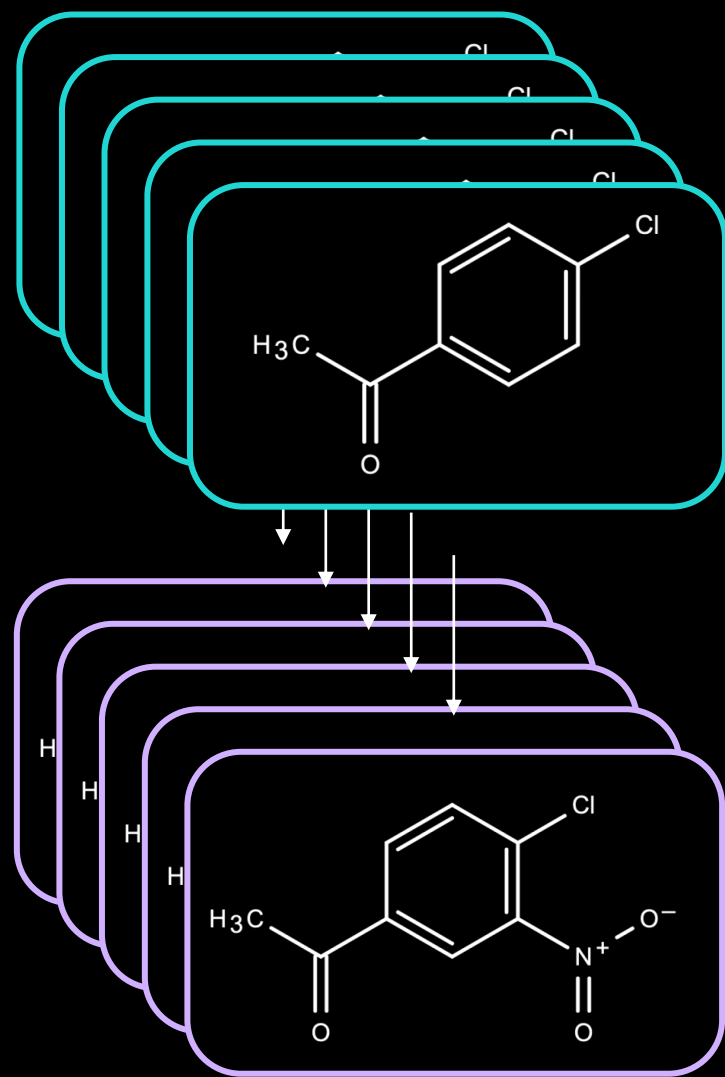2.7 g (12.3 mmol) 4,4-Dimethyl-1,2,3,4-tetrahydro-2-oxo-7-quinolinecarboxylic acid were added to a solution of 3.8 g (18.5 mmol) N,N'-dicyclohexylcarbodiimide and 1.1 ml (12.3 mmol) aniline in 80 ml dichloromethane. The reaction mixture was stirred for 4 hours at ambient temperature and the precipitate 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. ADD
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 $-1$

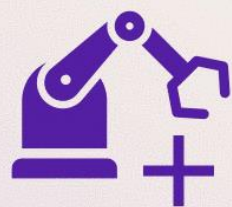O . c1ccncc1 . O=C(Cl)c1ccc(Cl)nc1 . COc1ccc(N)cc1 >> COc1ccc(NC(=O)c2ccc(Cl)nc2)cc1
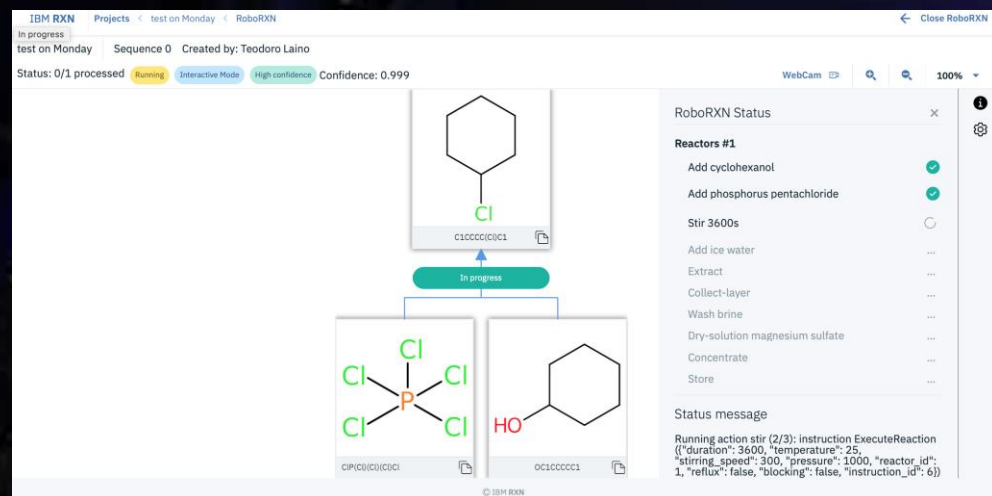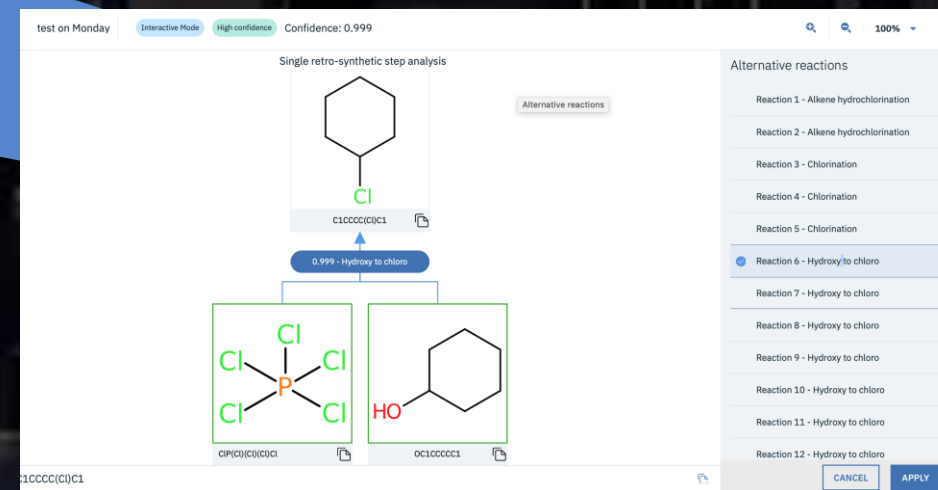
**Add**  **Add**  **Add**  **Stir**  **Add**  **Filter**  **Wash**  **Dry**  **Yield**

# SMILES-2-ACTIONS

# Flex Autoplant



FLEX AUTOPLANT robotic platform

# Analytics

# Synthesizing new molecule

Started: Nov 30 2020, 6:49am PT

Live from IBM RoboRXN

Action 2

## Adding C₂H₃F₃O₃S ⊕          Overview

In this action, the molecule methyl trifluoromethane sulfonate is added to Reactor 2.

Methyl trifluoromethane          2D  3D
C₂H₃F₃O₃S

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 C₂H₃F₃O₃S

00:06:00   ⏸  🔊   ● LIVE

# Enzymatic catalysis



GreenCatRXN

NCCR Catalysis

EC 1.4.3.-

$+$ $H_2O$ $+$ $O_2$

NH$_2$

TEMPO   NaClO   DCM

# The research lab is the central element in scientific discovery

Up to 70% of experimentation is not **reproducible** because of flawed experimental data or metadata[1]

Only one-third to one-half of original **findings** are confirmed in replication studies[2]

**Studies can last years instead of months** due to small differences in protocols[3]

More than half of researchers report being **unable to reproduce** their own experiments[4]

[1] Goncalves, R.S., Musen M.A., *Scientific Data* **6**, 190021 (2019)
[2] Aarts, A. A. *et al., Science* **349**, 943 (2015)
[3] Hines, W. C. *et al., Cell Rep.* **6**, 779–781 (2014)
[4] Baker, M., *Nature* **533**, 452-454 (2016)

# Multi-modal foundation models

## Capturing end-to-end business workflows for mining, optimization, generation, and automation



Multi-modal data

Vision

Audio

Text
*Description of actions in natural language*

Bytes
010011101101
110101011100
101110110101
100011110011

Foundation model generates workflows

Action 1
Metadata

Action 2
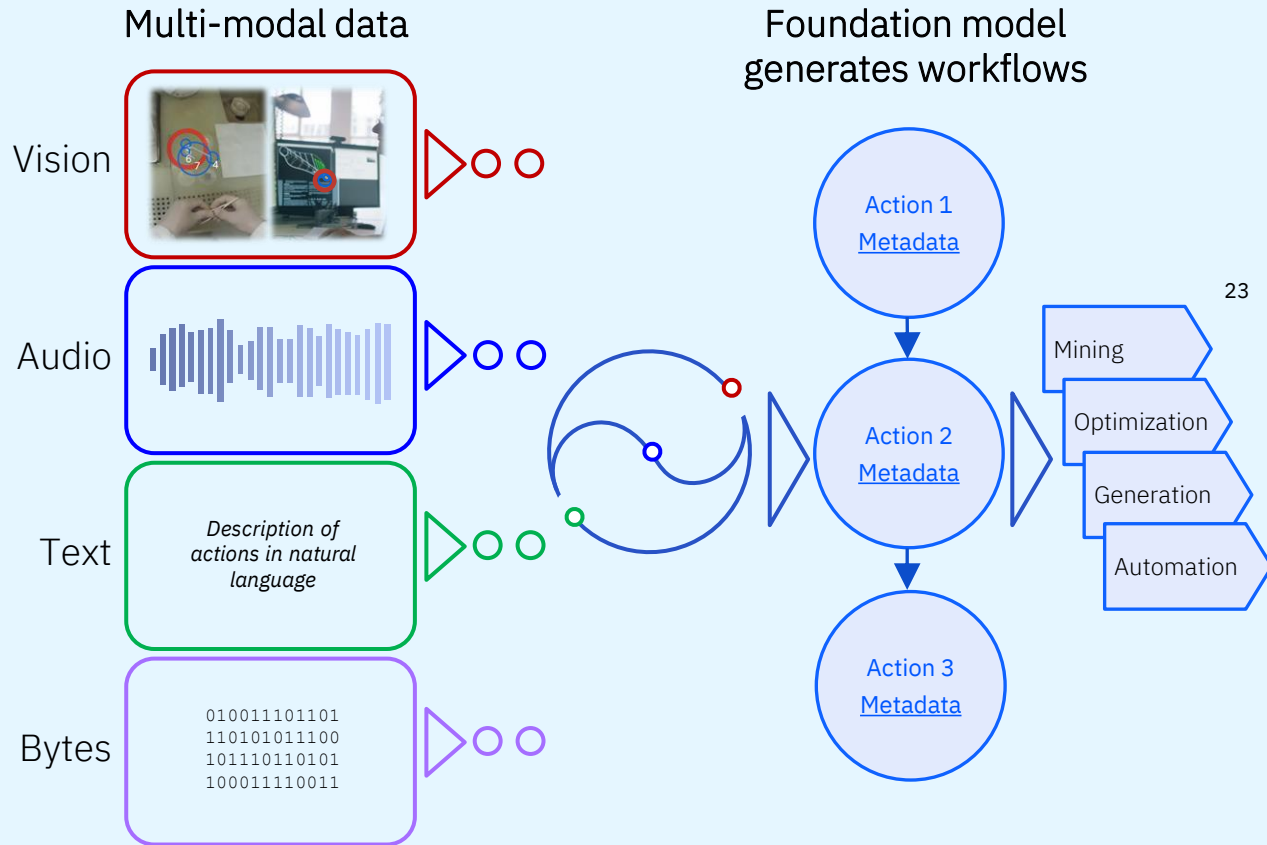Metadata

Action 3
Metadata

Mining

Optimization

Generation

Automation

23

## Key innovations

- AI foundation models for automatic documentation of manual procedures and validation of outcomes
- Hybrid and multi-cloud computing to automatically integrate all data and metadata of any experiment

## Benefits for the lab

- Capture all details needed to fully capture and describe an experiment
- Minimize the time wasted using different tools to organize data
- Reproduce any version of an experiment at any point in time
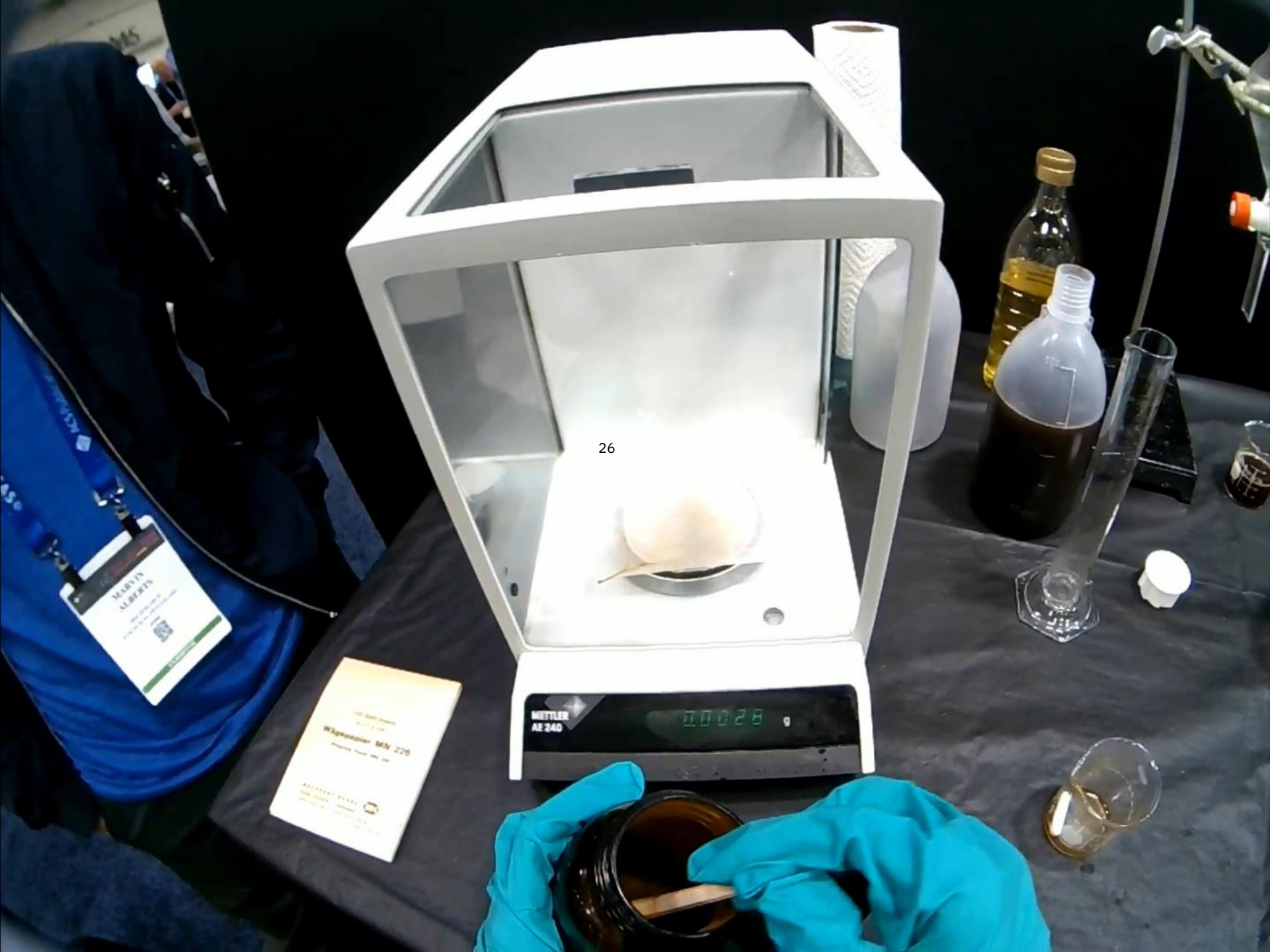- Discover patterns by continuously learning over all experimental data

25

26

# Workflow captured in the Lab that Learns



IBM **Lab that Learns**     Workflows

Workflow: Chih   Lab: ACS Booth Lab                              Documents    Start recording

Step

MeasureSolid

MeasureLiquid

Stir

AnalyticalMeasurement

Zoom: ———○——— 100%

☑ Settings    ☑ **Data**    Notes

**Device data**

pH measurement

5.257

Close                Save

IBM **Research** | © 2024 IBM Corporation

# Multimodal foundation models for the Lab that Learns



Pretraining Foundation Model for Egocentric Videos

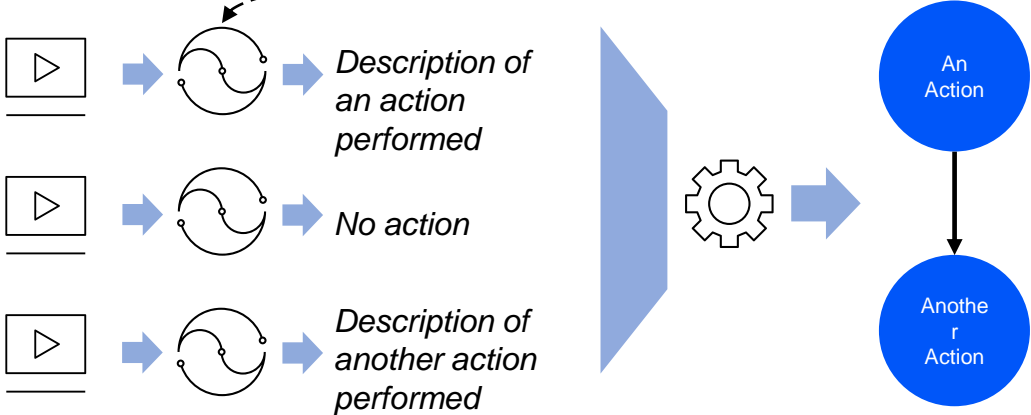ViT encoder → ViT decoder

Pretrained Foundation Model for Language

LM

Fine-tuning Vision-Language Foundation Model for Laboratory Procedures

ViT encoder → LM → *Description of an action performed*

Workflow Consolidation

*Description of an action performed*

*No action*

*Description of another action performed*

An Action → Another Action

# 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

Collaborators:

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