

Closing the loop

7 December 2023

Agenda

1. Scientific discovery as a loop
2. Molecular generation
3. Optimization strategies
4. Closing the loop: case studies
5. Hands-on lab



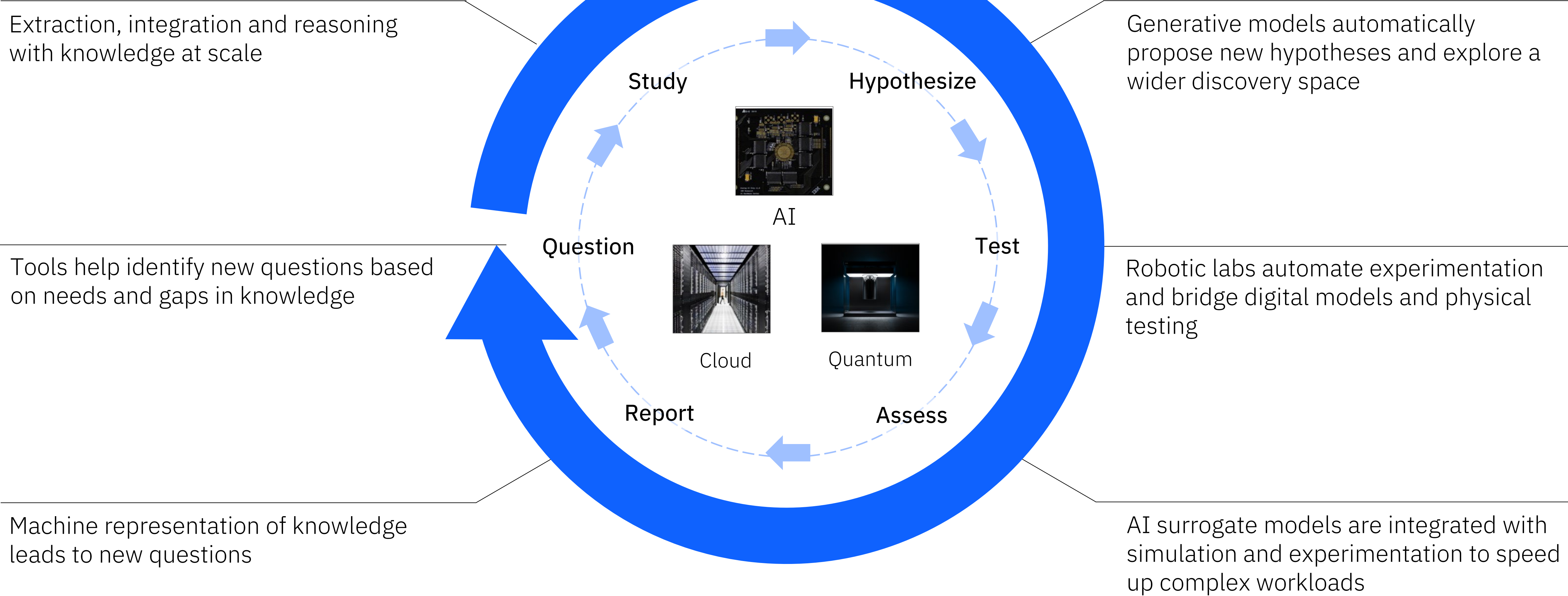
Scientific discovery as a loop

- Evolution of the scientific method
- AI for science
- State of the art

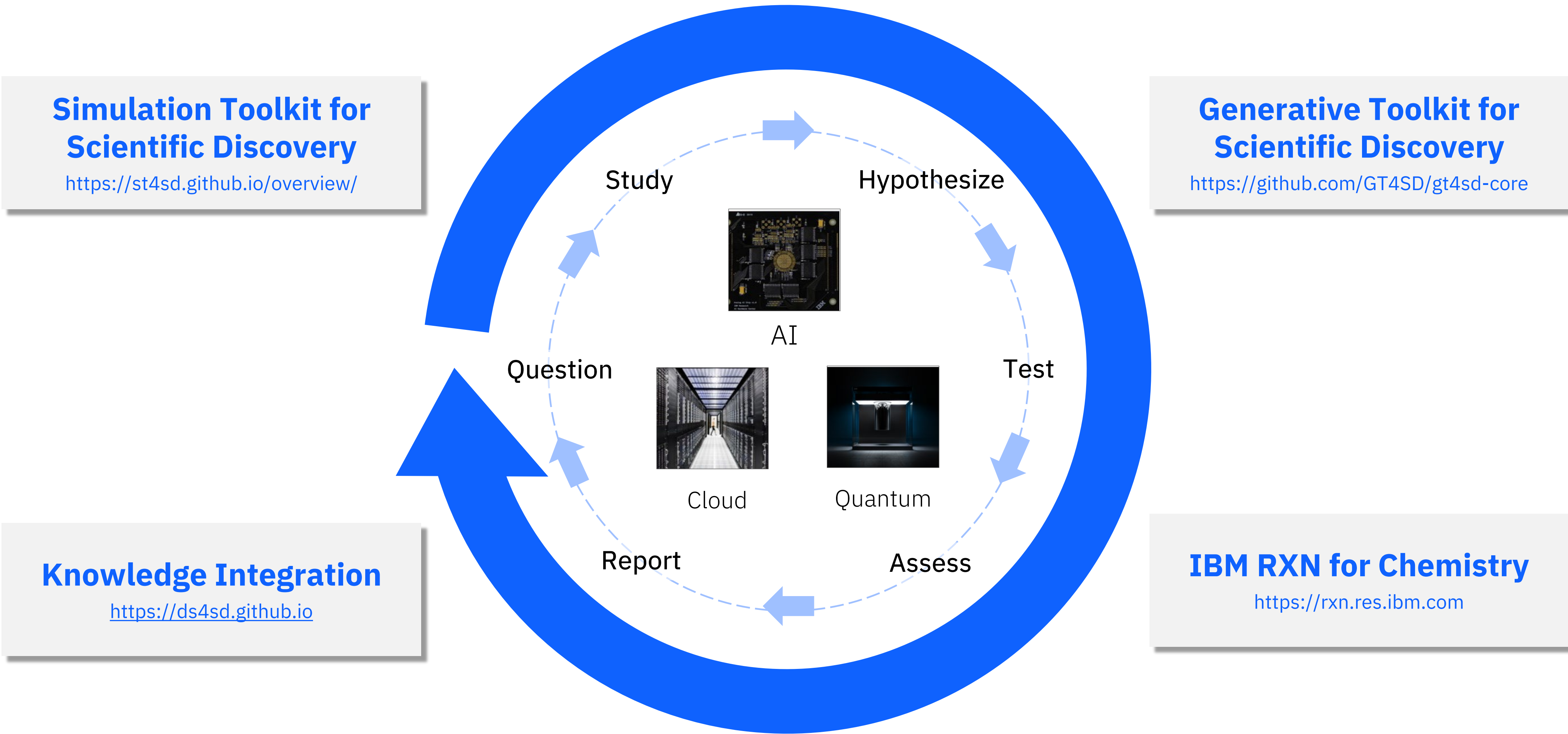
A new era of accelerated discovery

1 st Paradigm	2 nd Paradigm	3 rd Paradigm	4 th Paradigm	5 th Paradigm
Empirical Science	Theoretical Science	Computational Science	Big data-driven Science	Accelerated discovery
Observations Experimentation	Scientific laws Physics Biology Chemistry	Simulations Molecular dynamics Mechanistic models	Big data Machine learning Patterns Anomalies Visualization	Scientific knowledge at scale AI generated hypotheses Autonomous testing
Pre-Renaissance	~1600s	~1950	~2000	~2020

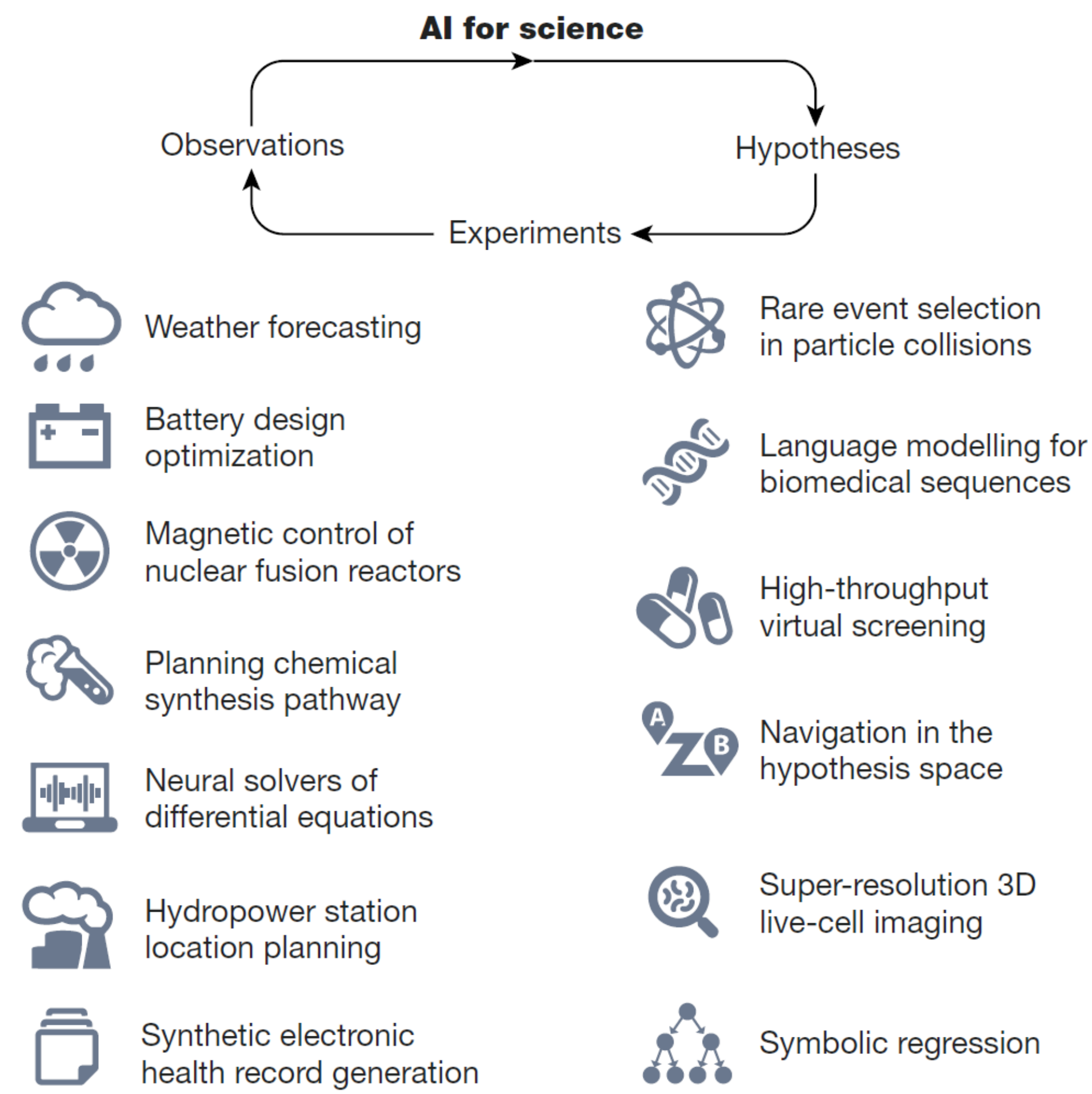
The scientific method is
humanity’s best model for
discovery



Scientific discovery is
accelerated by AI tools



The field of AI for scientific discovery

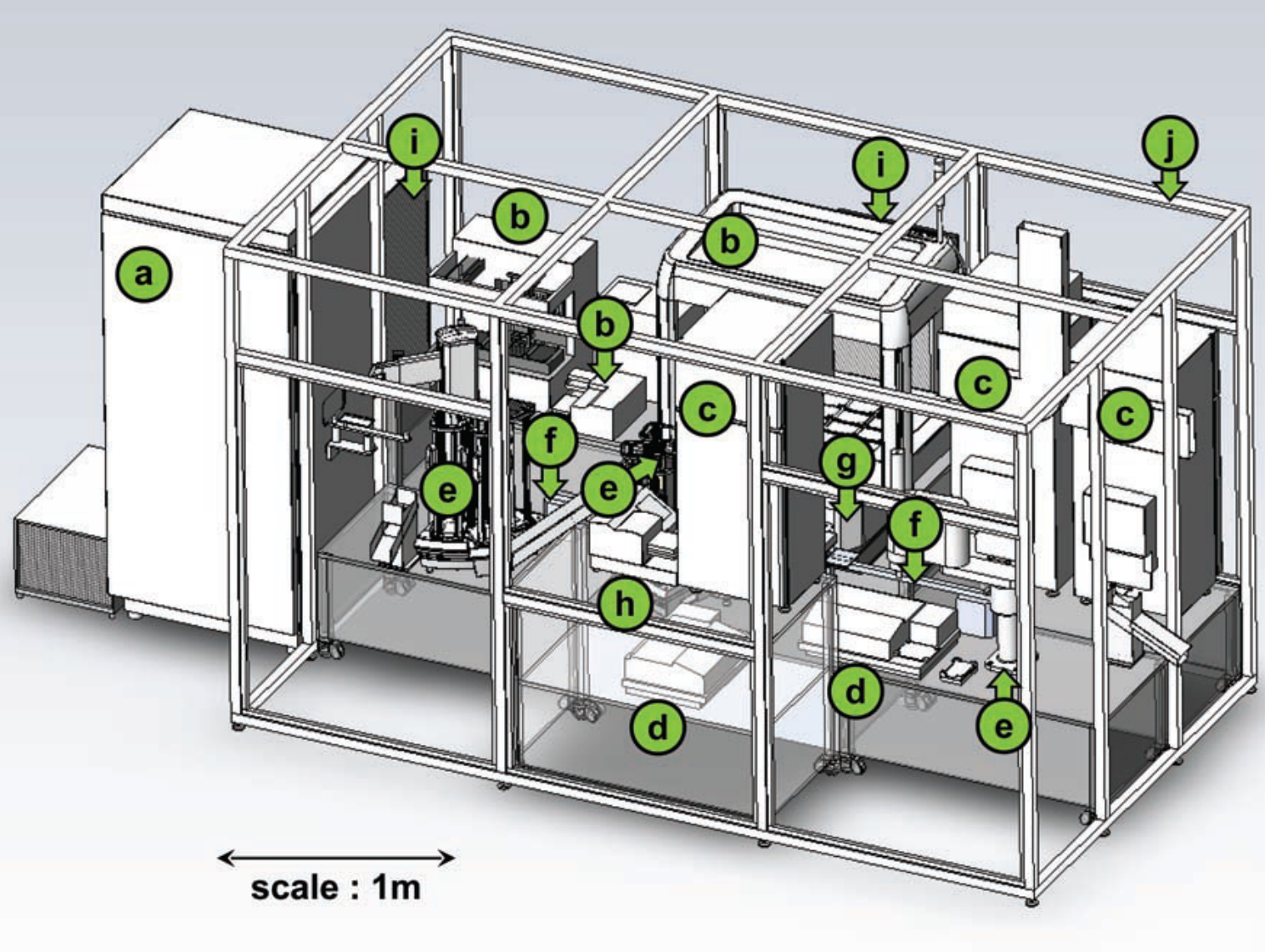


- AI-aided data collection and curation for scientific research
- Learning meaningful representations of scientific data
- AI-based generation of scientific hypotheses
- AI-driven experimentation and simulation

The Automation of Science

Ross D. King,^{1*} Jem Rowland,¹ Stephen G. Oliver,² Michael Young,³ Wayne Aubrey,¹ Emma Byrne,¹ Maria Liakata,¹ Magdalena Markham,¹ Pinar Pir,² Larisa N. Soldatova,¹ Andrew Sparkes,¹ Kenneth E. Whelan,¹ Amanda Clare¹

The basis of science is the hypothetico-deductive method and the recording of experiments in sufficient detail to enable reproducibility. We report the development of Robot Scientist “Adam,” which advances the automation of both. Adam has autonomously generated functional genomics hypotheses about the yeast *Saccharomyces cerevisiae* and experimentally tested these hypotheses



King, R. D. *et al.* The Automation of Science. *Science* **324**, 85–89 (2009)

Article

A mobile robotic chemist

<https://doi.org/10.1038/s41586-020-2442-2>

Received: 1 November 2019

Accepted: 25 March 2020

Published online: 8 July 2020

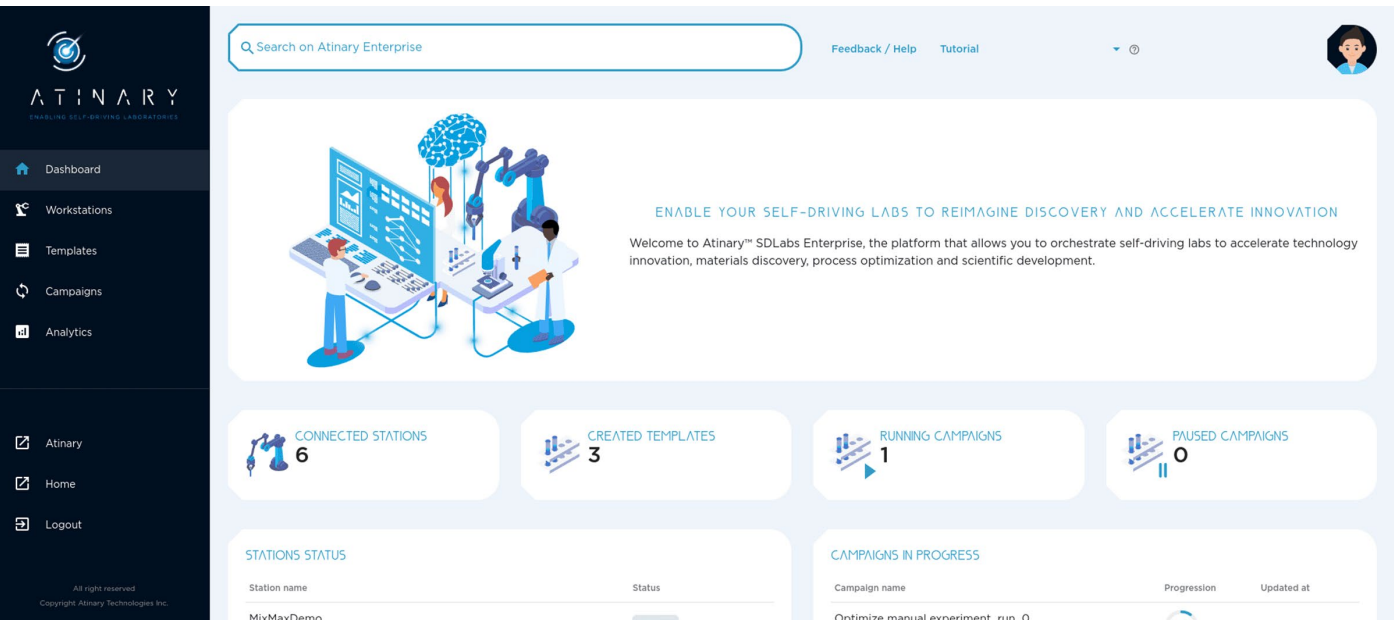
Check for updates

Benjamin Burger¹, Phillip M. Maffettone¹, Vladimir V. Gusev¹, Catherine M. Aitchison¹, Yang Bai¹, Xiaoyan Wang¹, Xiaobo Li¹, Ben M. Alston¹, Buyi Li¹, Rob Clowes¹, Nicola Rankin¹, Brandon Harris¹, Reiner Sebastian Sprick¹ & Andrew I. Cooper¹

Technologies such as batteries, biomaterials and heterogeneous catalysts have functions that are defined by mixtures of molecular and mesoscale components. As yet, this multi-length-scale complexity cannot be fully captured by atomistic simulations, and the design of such materials from first principles is still rare^{1–5}. Likewise, experimental complexity scales exponentially with the number of variables, restricting most searches to narrow areas of materials space. Robots can assist in experimental searches^{6–14} but their widespread adoption in materials research is challenging because of the diversity of sample types, operations, instruments and measurements required. Here we use a mobile robot to search for improved



Burger, B. *et al.* A mobile robotic chemist. *Nature* **583**, 237–241 (2020)



<https://youtu.be/SX26XRFx0U0>

Making AI for chemistry
available to everyone

<https://rxn.res.ibm.com/>

IBM **RXN** for Chemistry

Twitter @forRXN

Publications

Get Started! →

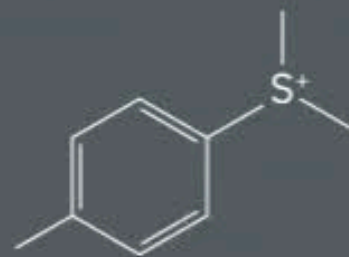
RXN for Chemistry

Use AI to predict outcomes of chemical reactions for optimized synthesis methods, and to automatically generate chemical procedures for use in manual or automated lab operations.

Start your Project Now



Synthesizing new molecule



Started: Nov 30 2020, 6:49am PT

Live from IBM RoboRXN

Action 2

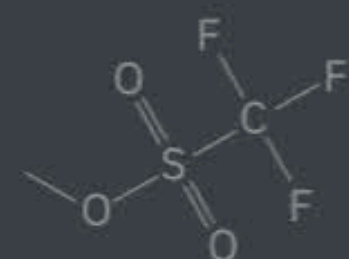
Overview

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

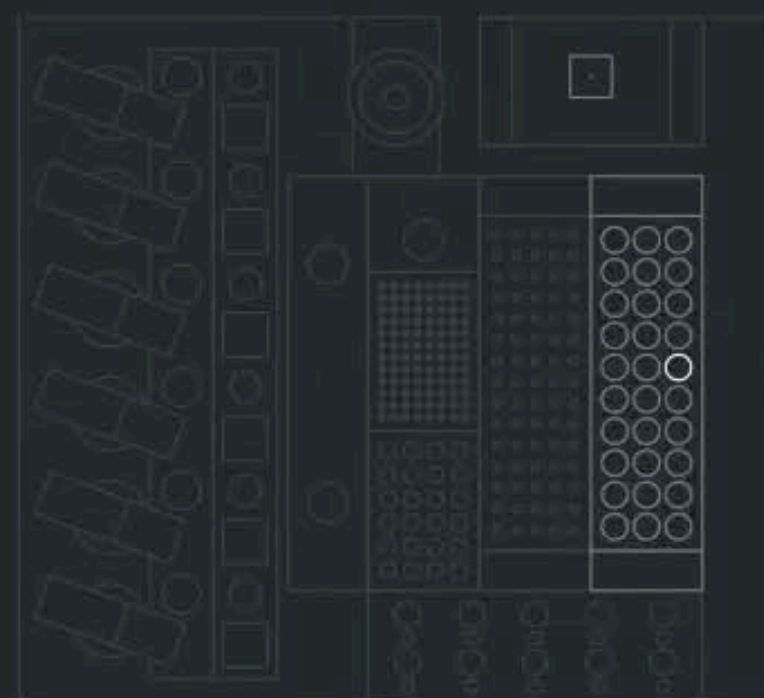


Methyl trifluoromethane sulfonate is a brown liquid. Insoluble in water. This material is a very reactive methylating agent, also known as methyl triflate.

● 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



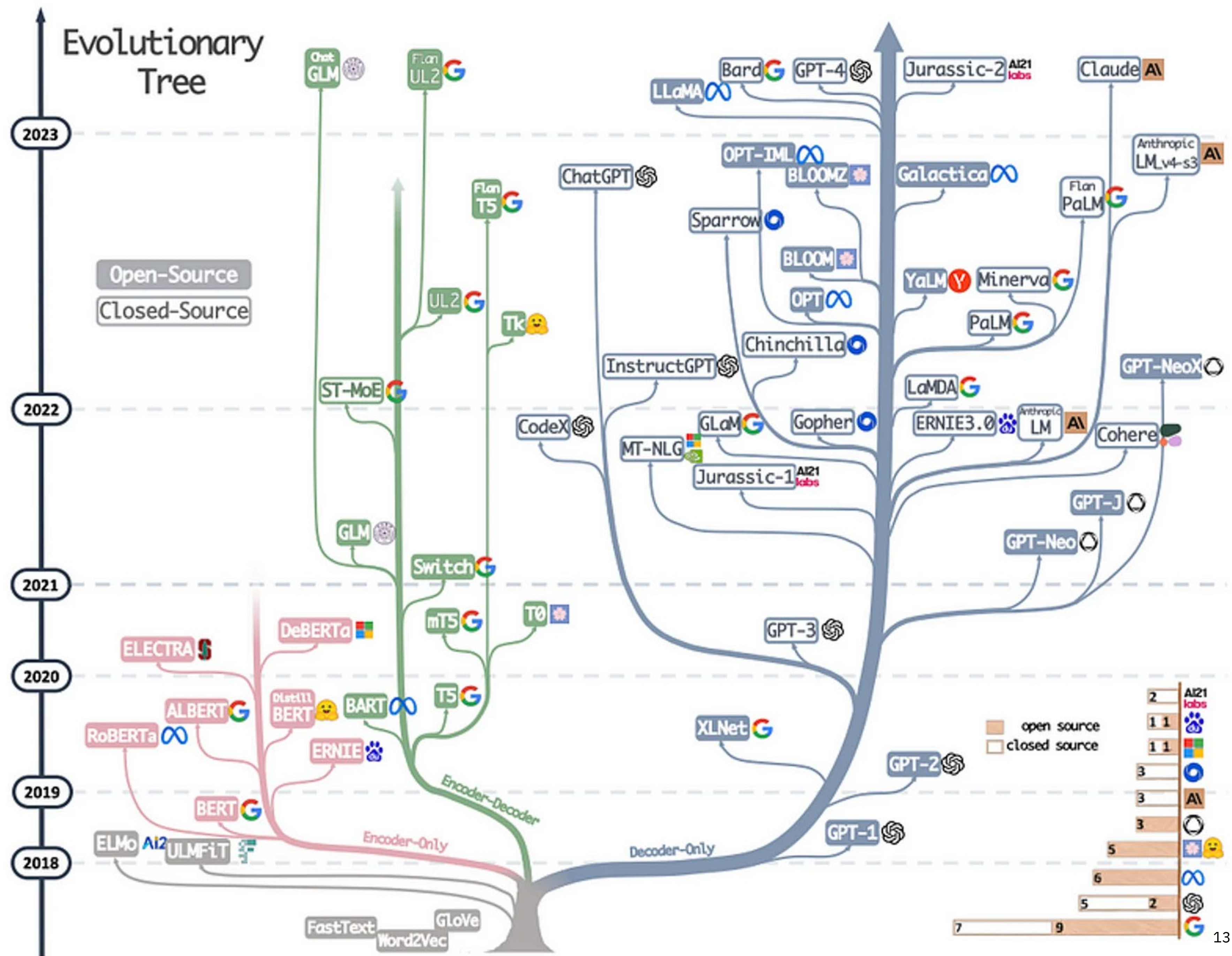
Live view module



Molecular generation

- Language models for molecular discovery
- Generative AI
- Graph neural networks

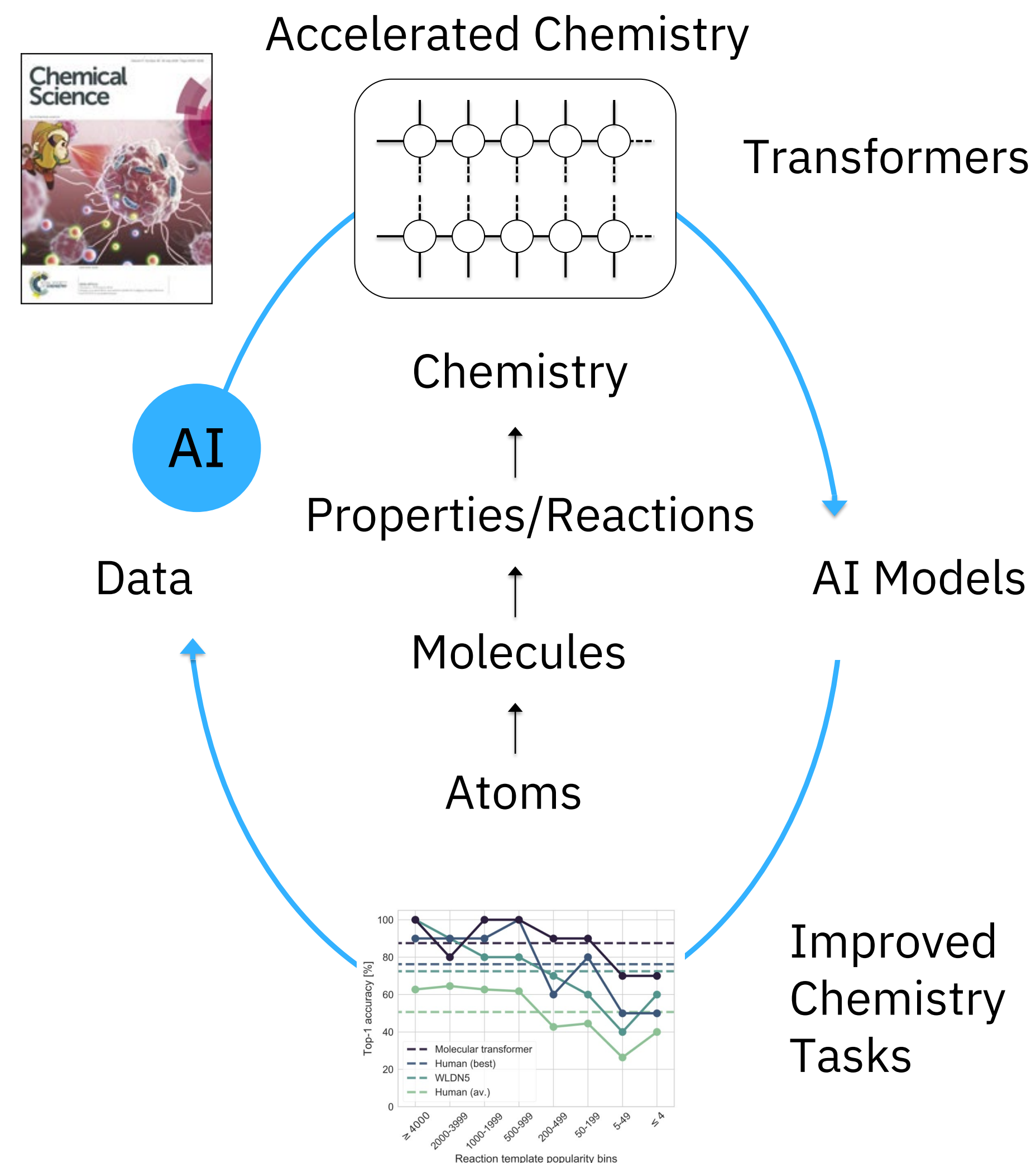
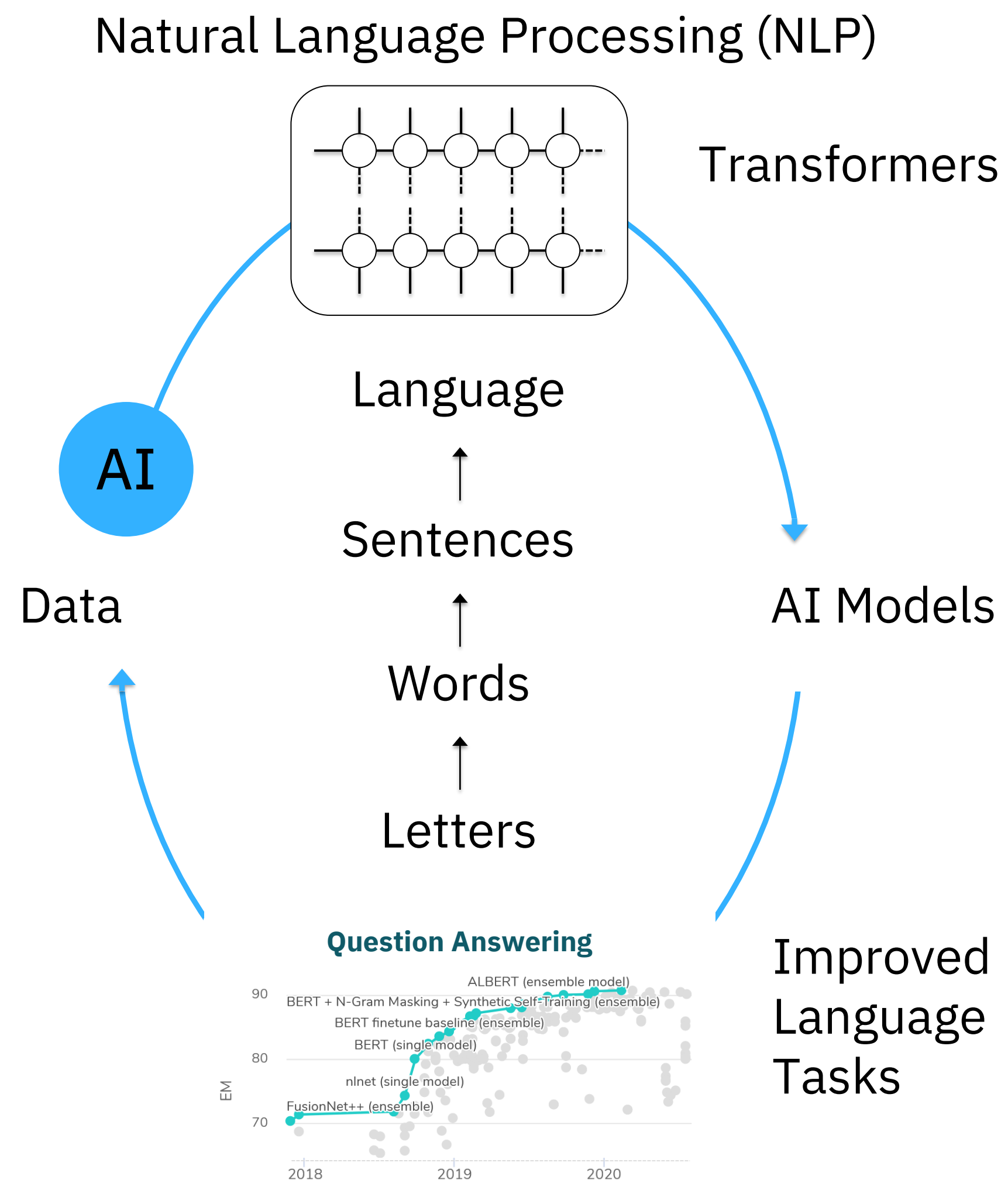
The world of large language models (LLMs)



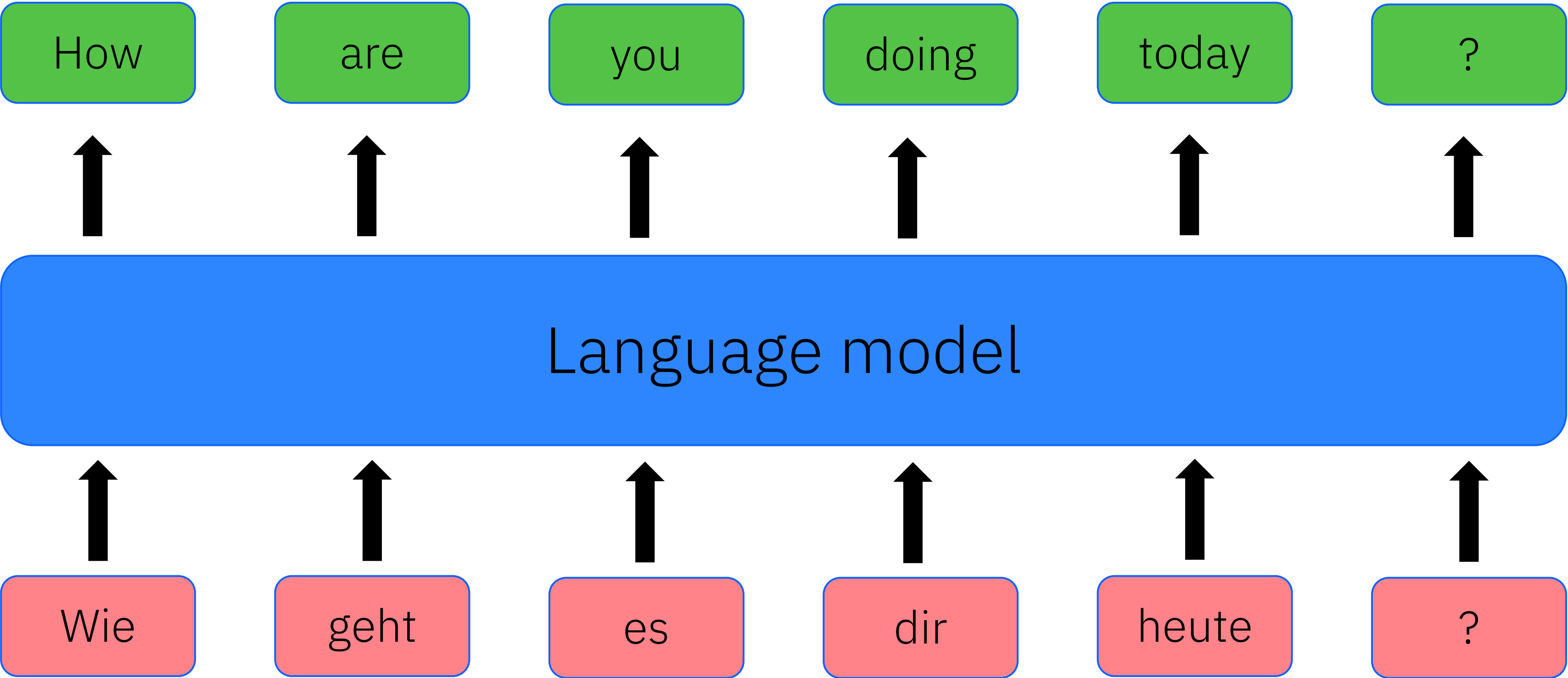
Credit: <https://blog.sylphai.com/introduction-to-large-language-models>

The same AI breakthroughs for language are changing scientific discovery

Generative modeling and transformers are achieving new breakthroughs in scientific disciplines

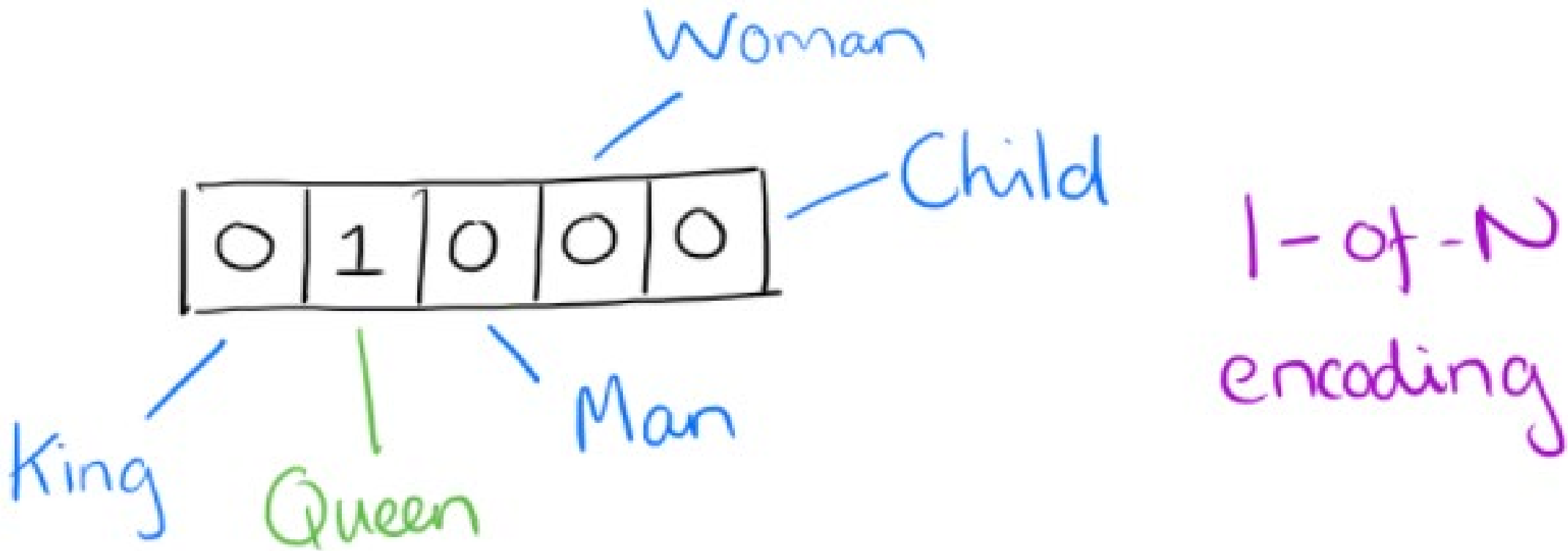


Translate from German to
English



Representing words as input to
a neural network

One-hot encoding?



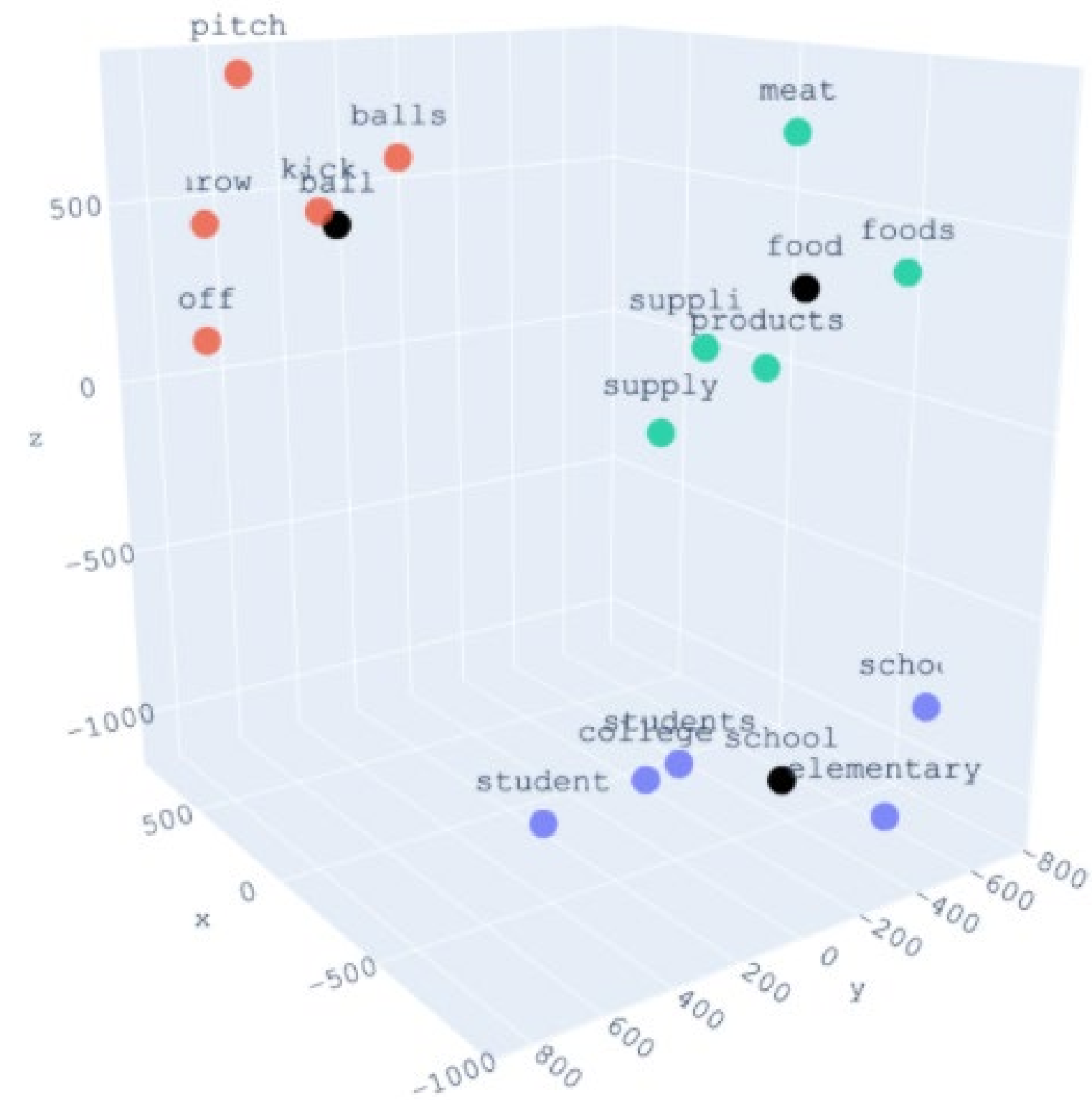
Learned encodings!

	King	Queen	Woman	Princess
Royalty	0.99	0.99	0.02	0.98
Masculinity	0.99	0.05	0.01	0.02
Femininity	0.05	0.93	0.999	0.94
Age	0.7	0.6	0.5	0.1
...	...			

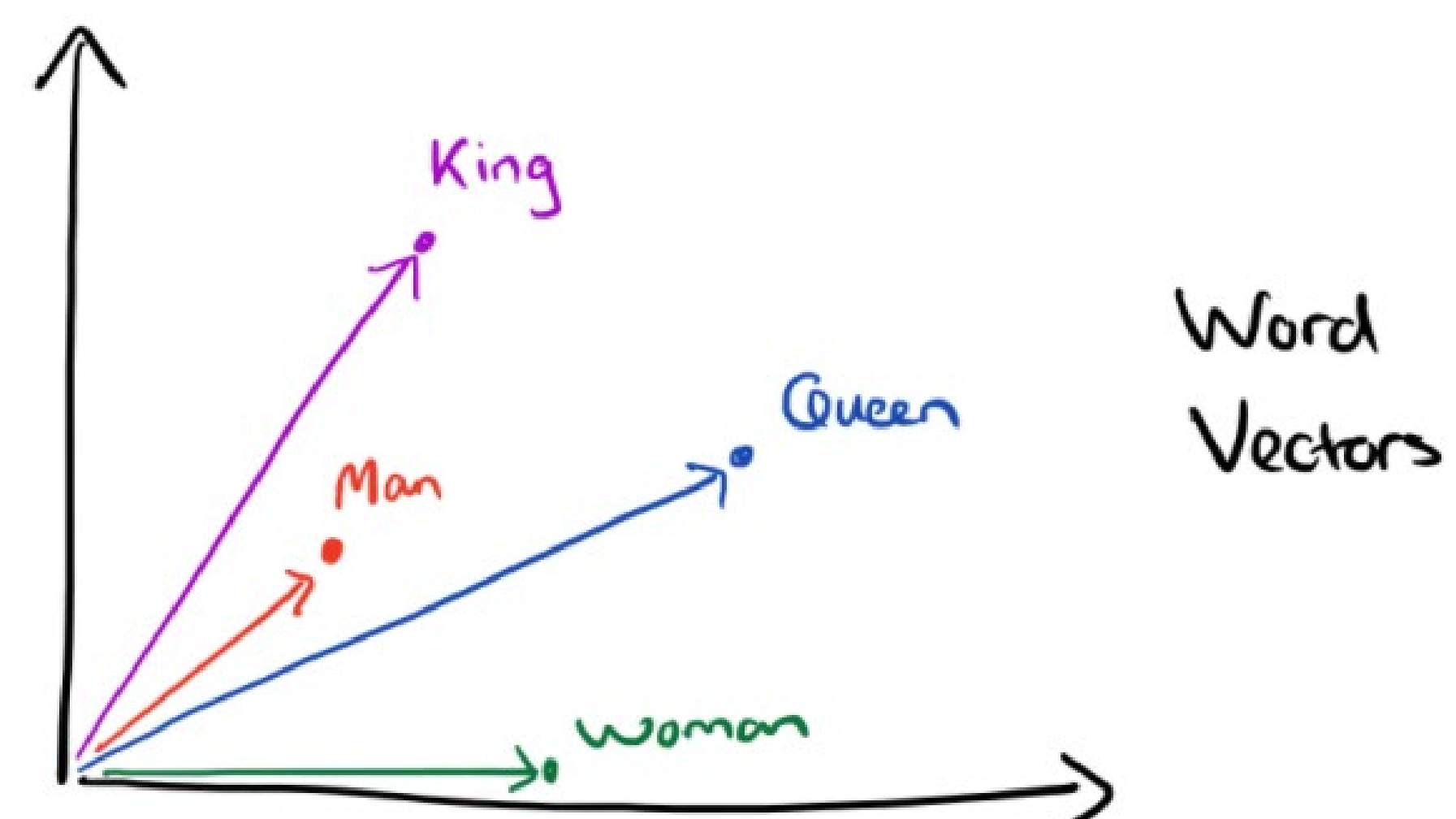
Word2Vec: Mikolov, T., Sutskever, I., Chen, K., Corrado, G. & Dean, J. Distributed Representations of Words and Phrases and their Compositionality. *arXiv:1310.4546* (2013)

Words as vectors

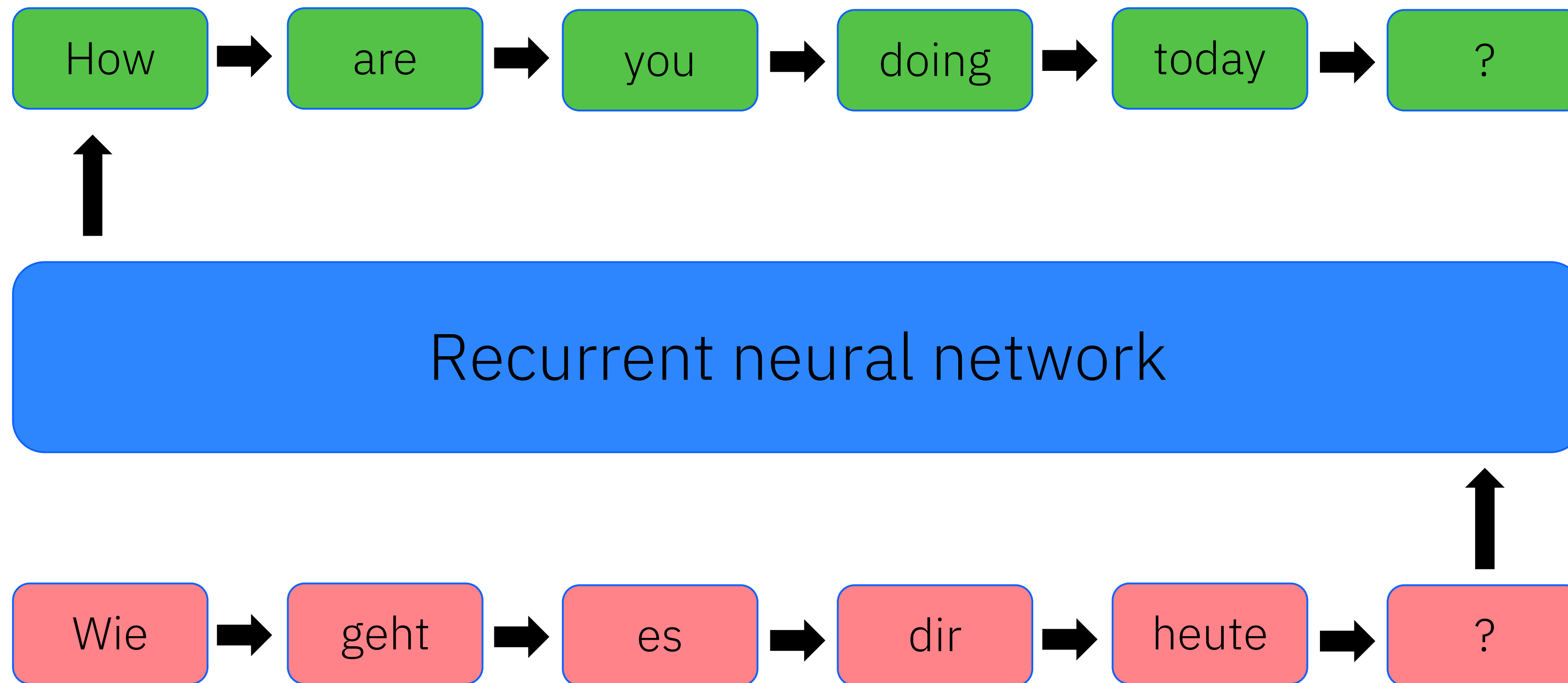
Similar words cluster



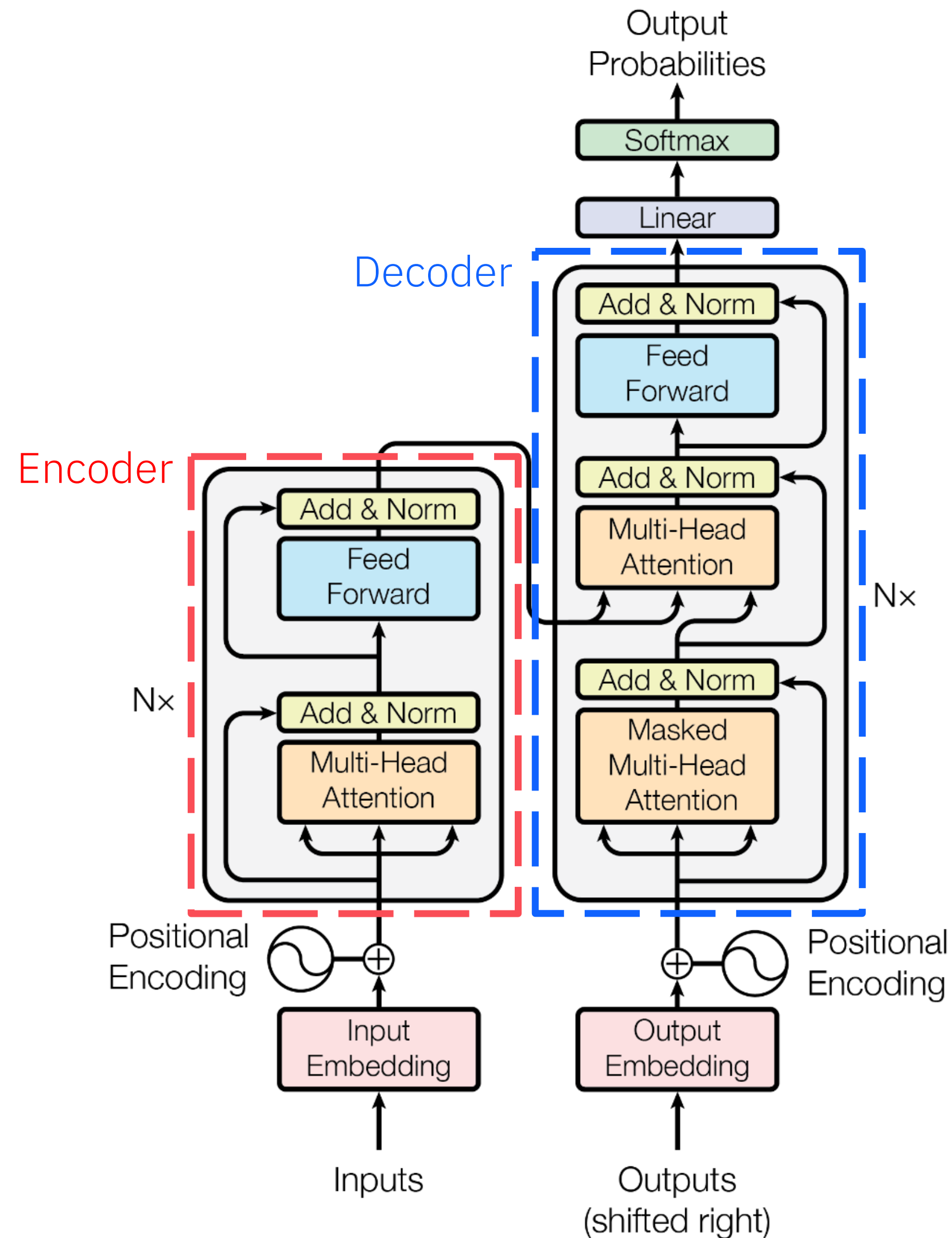
Word arithmetics



Until 2017, language models were based on recurrent neural networks (RNNs)



The transformer model



Key concept: **parallel multi-head attention mechanism**

- Tokenizers convert input (such as text) to tokens
- Embedding layer converts tokens and their positions to vector representations
- Transformer layers consisting of alternating attention and feedforward layers extract linguistic information

Superior translation quality and less training effort compared to previous state-of-the-art with recurrent neural networks (RNNs).

Encoder-only models:

Best suited for classification tasks
(e.g. sentence and word classification, entity recognition, extractive Q&A)

Decoder-only models (also called *auto-regressive models*)

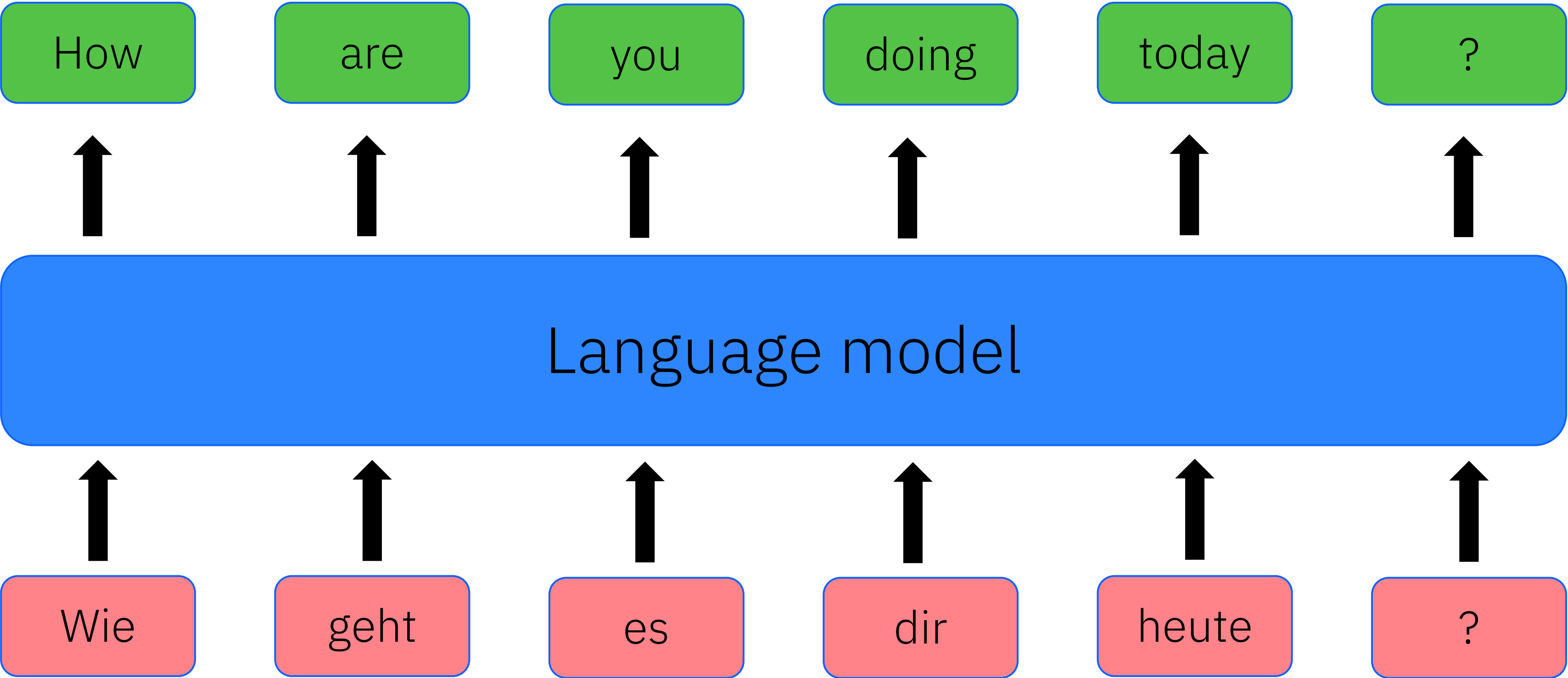
Best suited for generative tasks
(e.g. text generation)

Encoder-decoder models (also called *sequence-to-sequence models*)

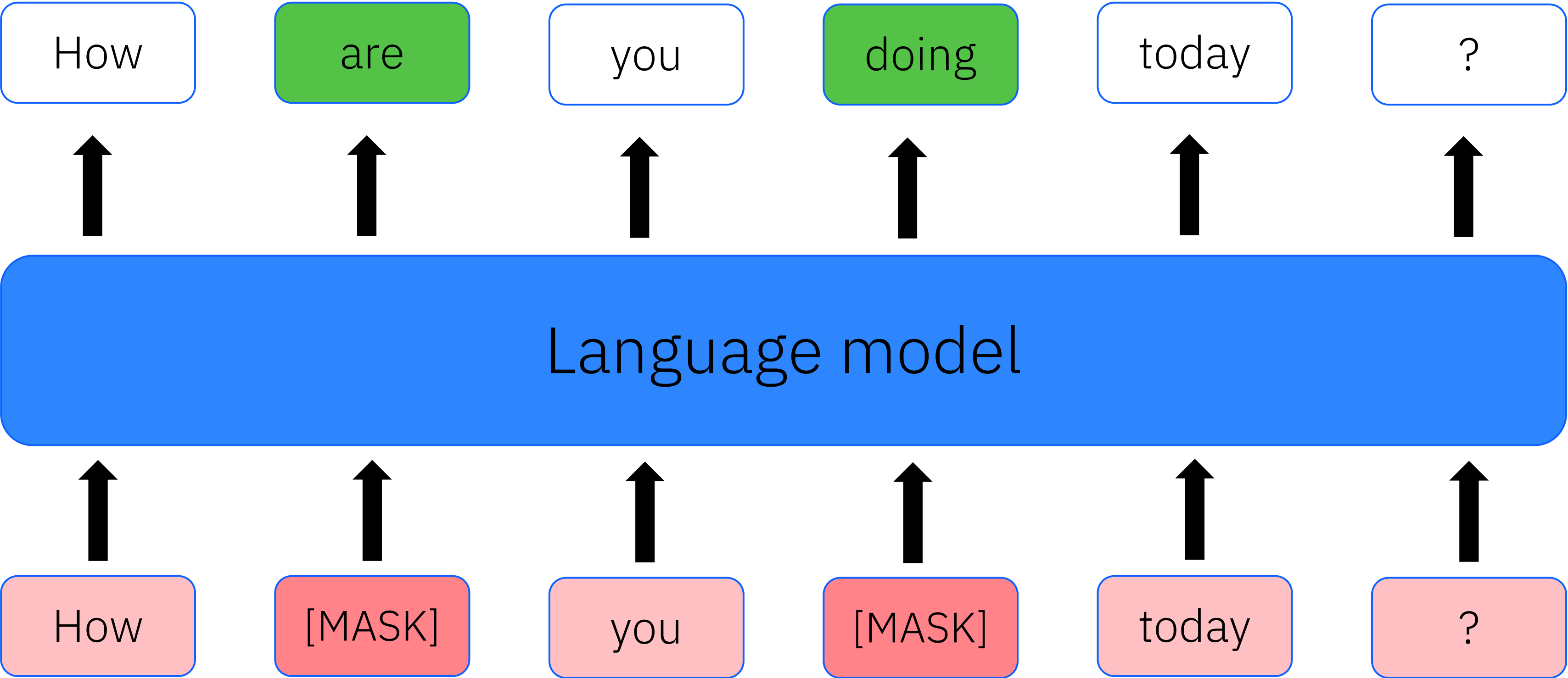
Best suited for generative tasks depending on a given input
(e.g. summarization, translation, generative Q&A)

📖 Vaswani, A. *et al.* Attention is All you Need. *NeurIPS* (2017)

Since 2017, transformers networks are the state-of-the-art for language models



Masked language modelling: training transformers with self-supervision



Transformers impacted multiple application domains

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.co

AN IMAGE IS WORTH 16X16 WORDS:
TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy*,†, Lucas Beyer*, Alexander Kolesnikov*, Dirk Weissenborn*,
Xiaohua Zhai*, Thomas
Georg Heigold, S
*equal t
{adosovit

RETURN TO ISSUE | < PREV RESEARCH ARTICLE NEXT >

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*

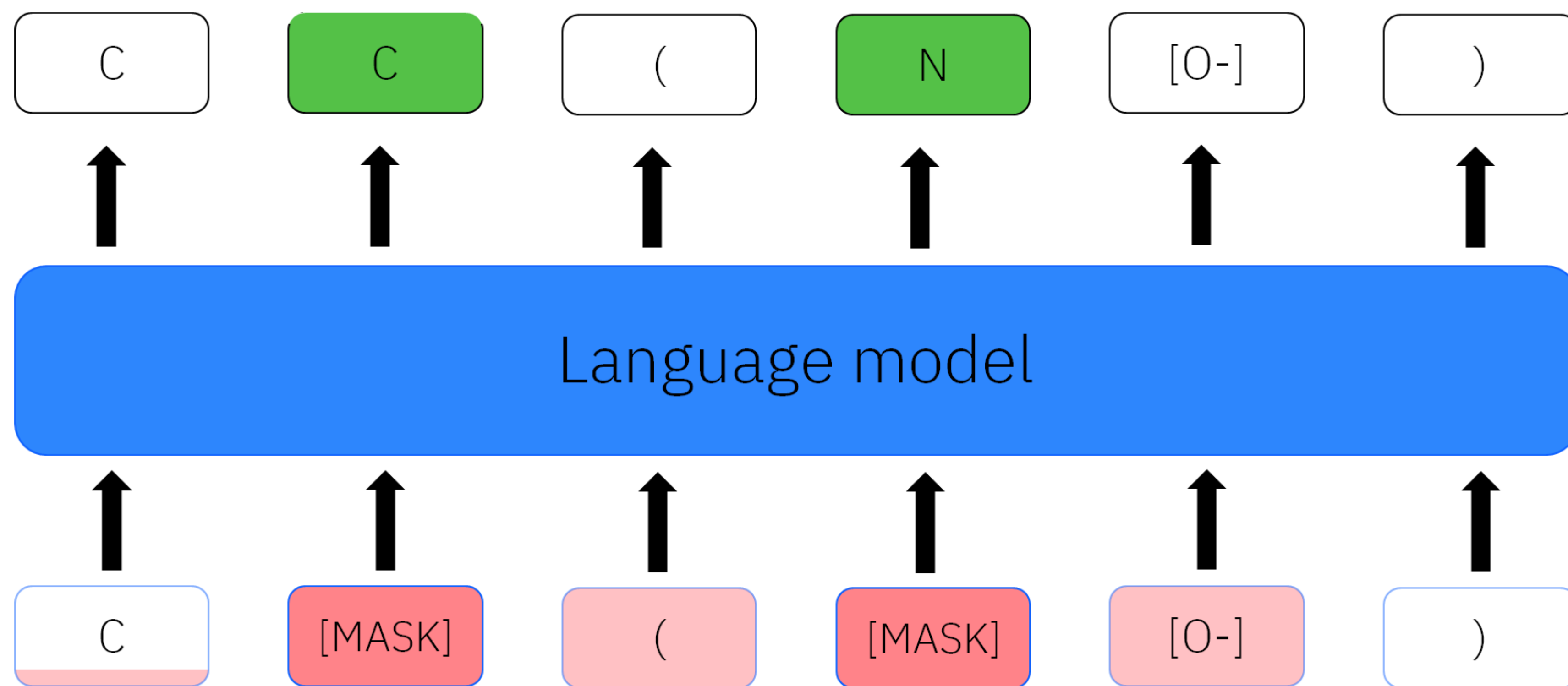
Cite this: ACS Cent. Sci. 2
Publication Date: August 30,
https://doi.org/10.1021/acsc
Copyright © 2019 American
RIGHTS & PERMISSIONS
PDF (1 MB)

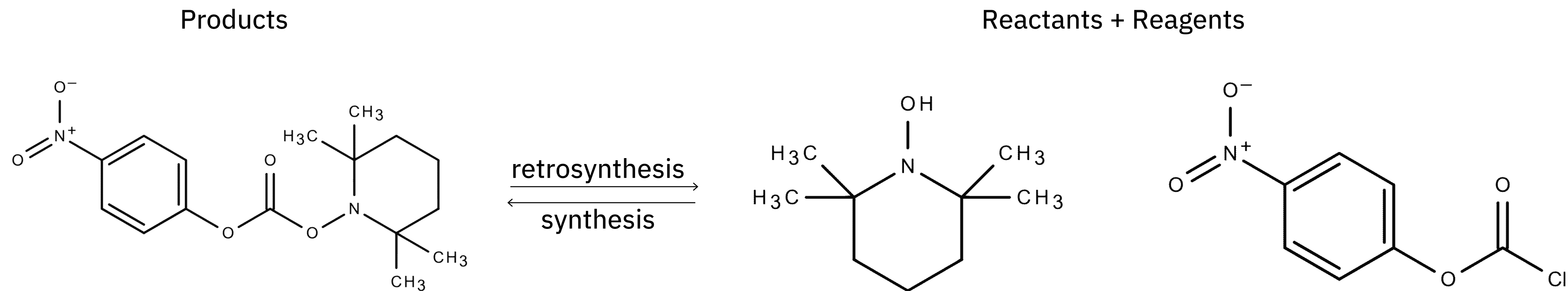
Pretrained Transformers as Universal Computation Engines

Kevin Lu,^{1,2}
¹ UC Ber

Decision Transformer: Reinforcement Learning via Sequence Modeling

Lili Chen*,¹, Kevin Lu*,¹, Aravind Rajeswaran², Kimin Lee¹,
Aditya Grover^{2,3}, Michael Laskin¹, Pieter Abbeel¹, Aravind Srinivas^{†,4}, Igor Mordatch^{†,5}





Textual representation (SMILES)

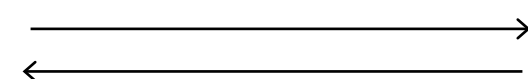
CC1(CCCC(N1OC(OC1=CC=C([N+](=O)[O-])C=C1)=O)(C)C)C

CC1(CCCC(N1O)(C)C)C

O(C(=O)Cl)C1=CC=C([N+](=O)[O-])C=C1

“Sentence of atoms”

C C 1 (C C C C (N 1 O C (O C 1
= C C = C ([N+] (= O) [O-]) C
= C 1) = O) (C) C) C



C C 1 (C C C C (N 1 O) (C) C) C . O (C (= O) Cl)
C 1 = C C = C ([N+] (= O) [O-]) C = C 1

Translation

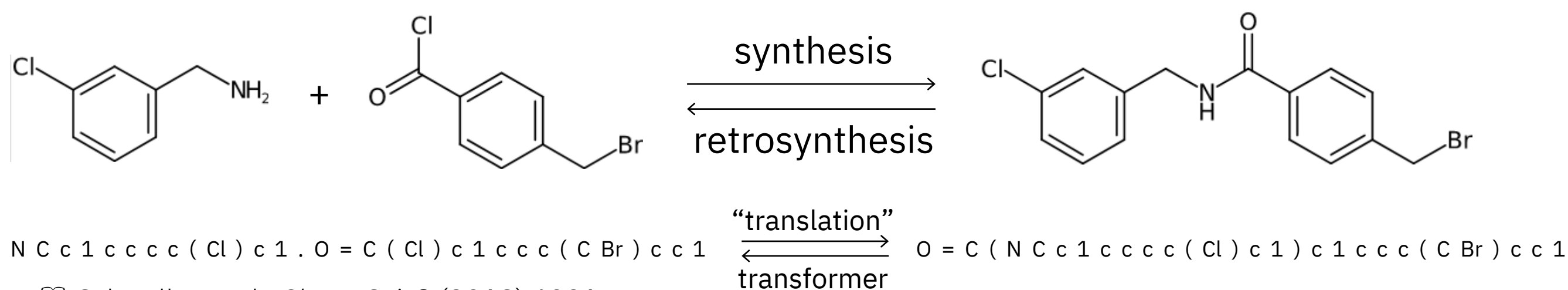
English



Spanish

Automating lab synthesis and experimentation
with help of language models

1. AI-based chemical reaction prediction



📖 Schwaller et al., *Chem. Sci.* **9** (2018) 6091

📖 Schwaller et al., *ACS Cent. Sci.* **5** (2019) 1572

📖 Schwaller et al., *Chem. Sci.* **11** (2020) 3316

2. Chemical procedures from text (Paragraph2Actions)

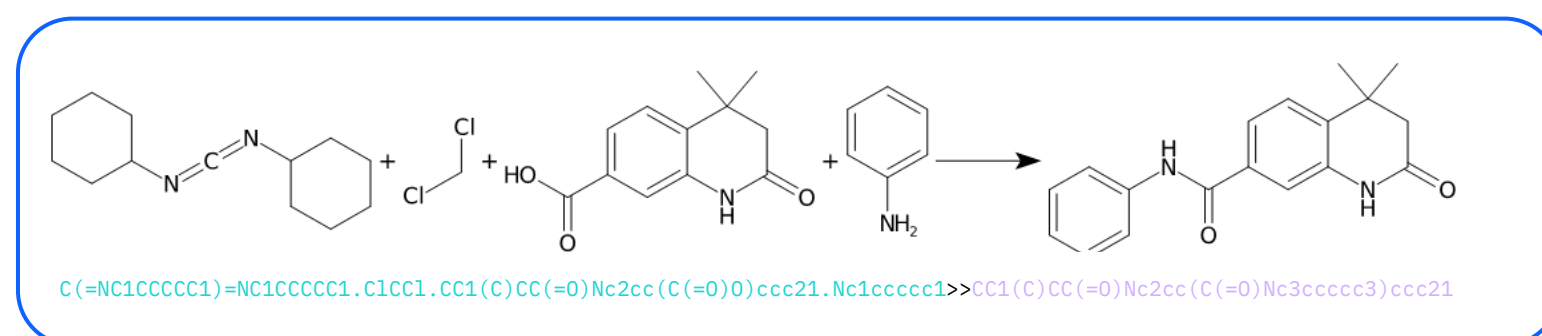
The TFA was removed in vacuo and a saturated solution of NaHCO₃ was added.

translation

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

📖 Vaucher et al., *Nat. Comm.* **11** (2020) 3601

3. Chemical procedures from reactions (Smiles2Actions)

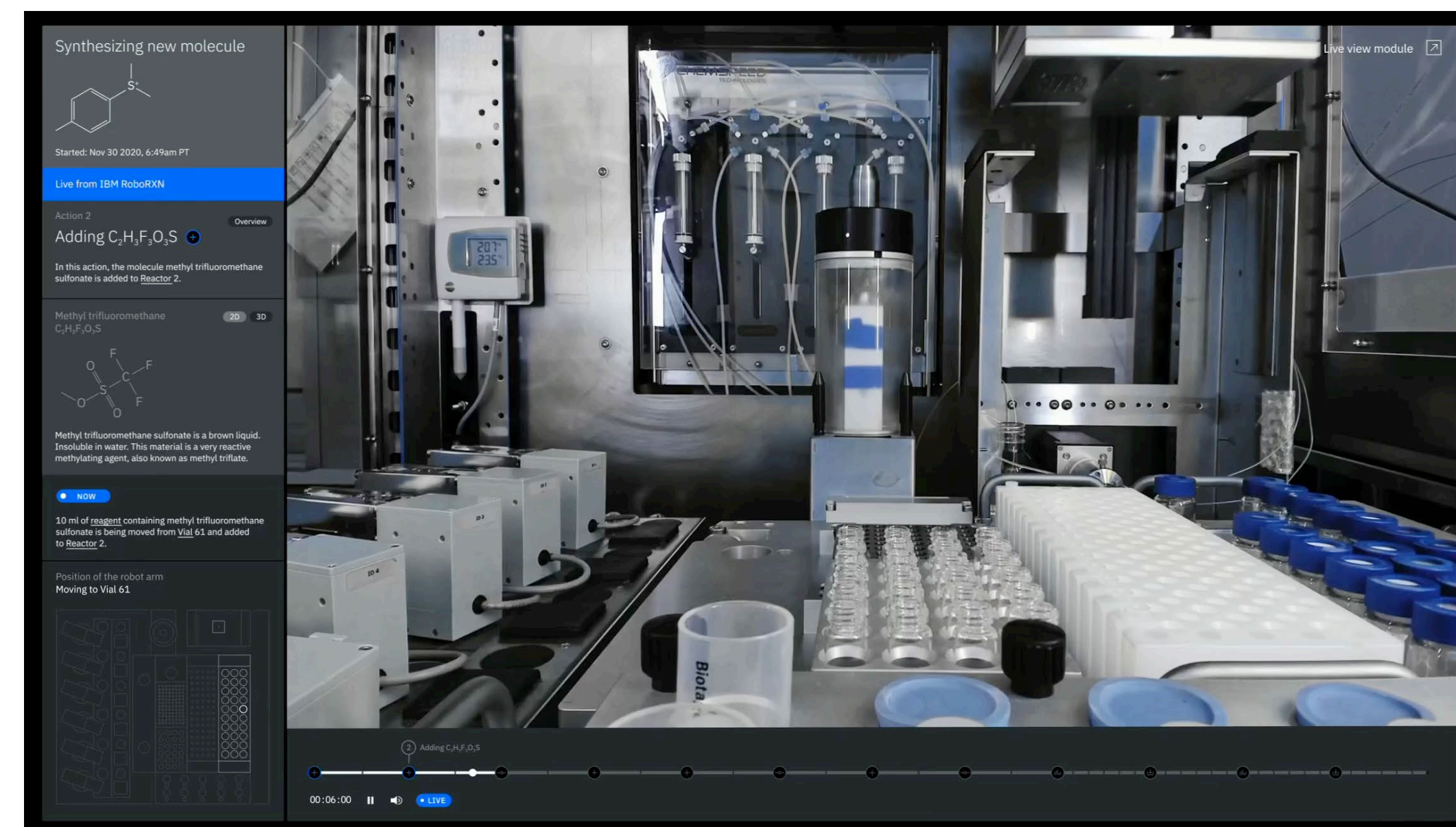


📖 Vaucher et al., *Nat. Comm.* **12** (2022) 2573

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\$

Automated synthesis via cloud + robotic lab

<https://rxn.res.ibm.com/rxn/robo-rxn>



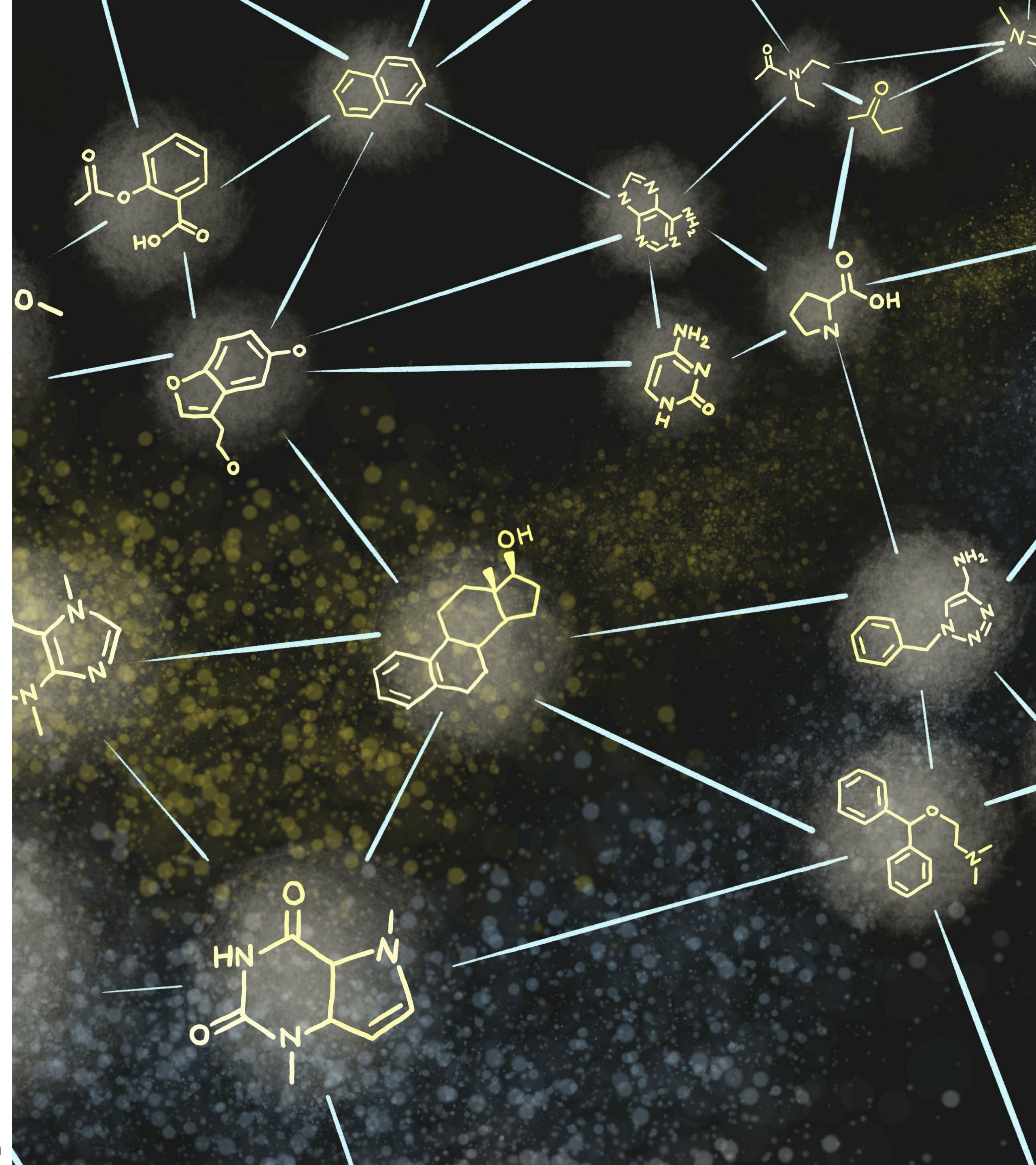
30k+
global users
via cloud

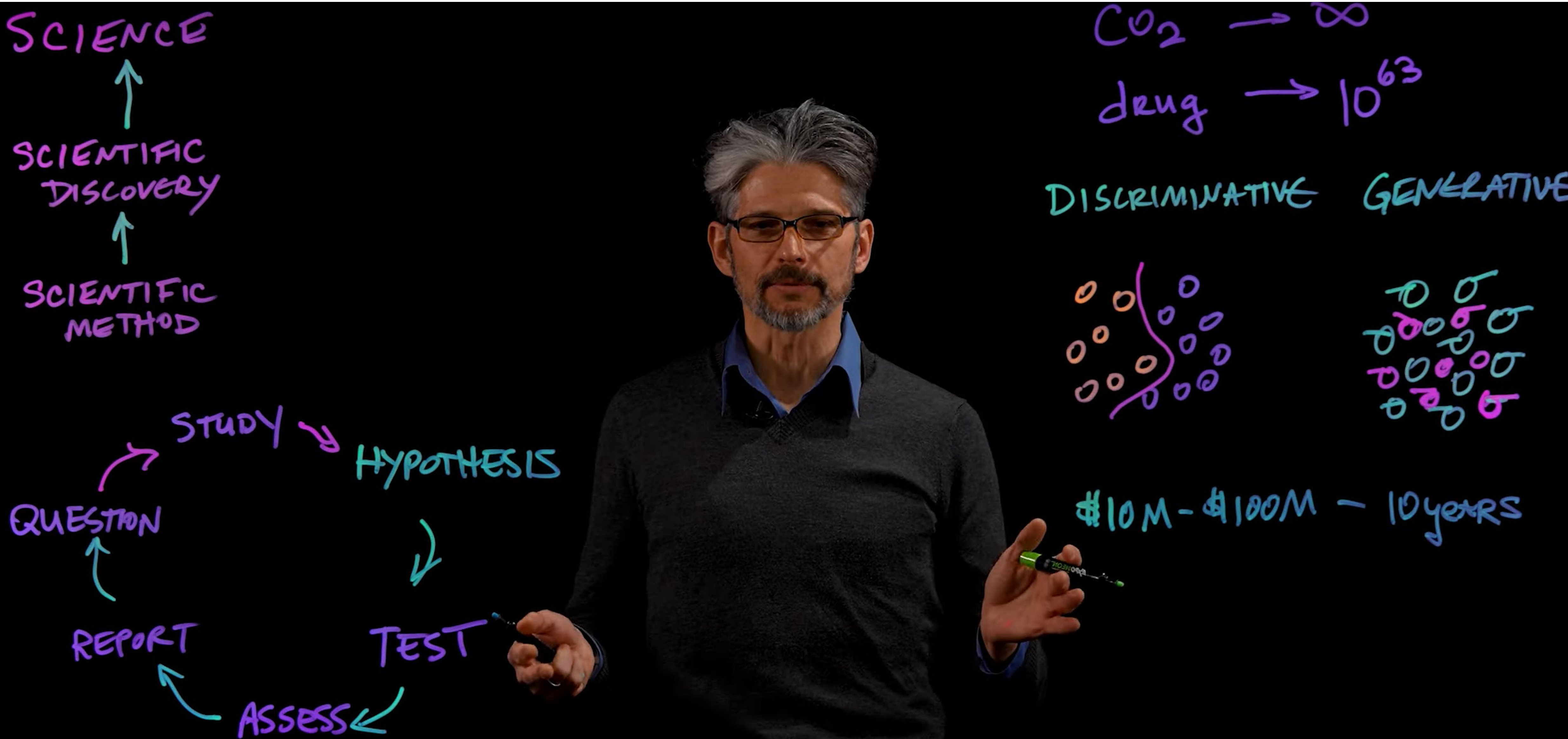
10+ million
reaction
predictions

The size of chemical space

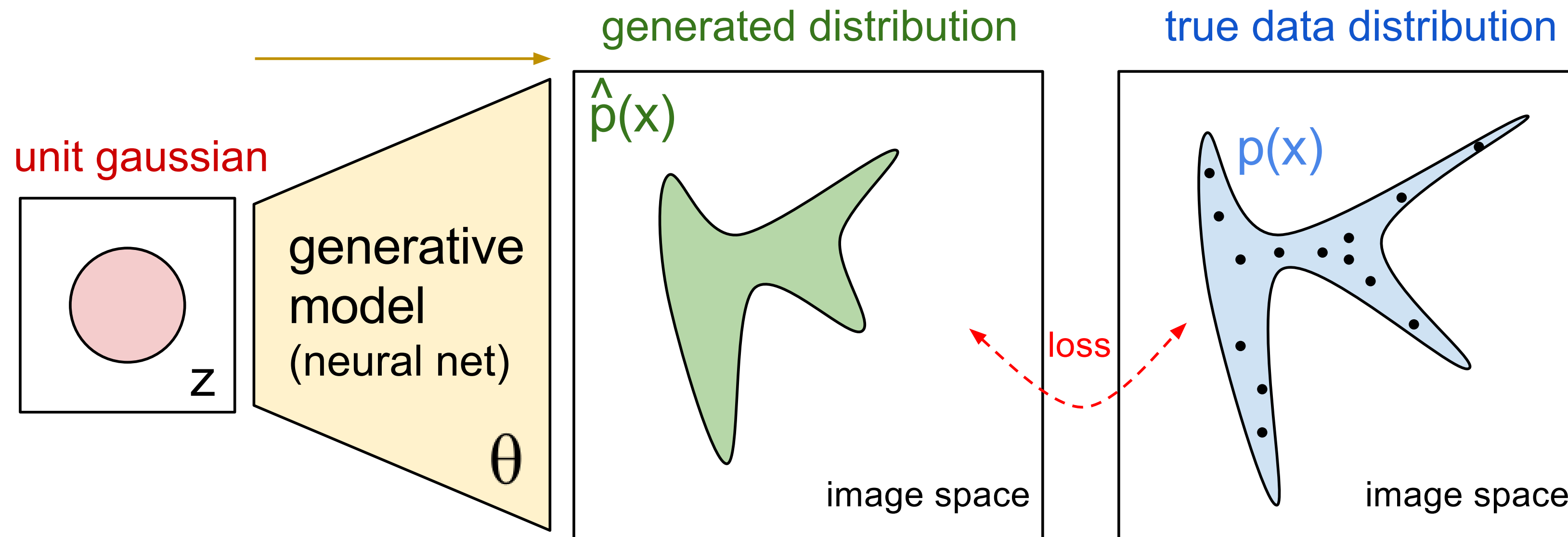
Number of compounds	Reference
$6,2 \times 10^{13}$	Henze and Blair [4]
$1,3 \times 10^{15}$	Blair and Henze [5]
10^{21}	Weaver and Weaver [8]
10^{23}	Ertl [7]
10^{26}	Ogata et al. [24]
10^{33}	Weininger [23]
10^{33}	This work (see ref. below)
10^{60}	Bohacek et al. [6]
10^{100}	Walters et al. [26]
10^{180}	Weininger [27]

📖 Polishchuk, P. G., Madzhidov, T. I. & Varnek, A. Estimation of the size of drug-like chemical space based on GDB-17 data. *J Comput Aided Mol Des* **27**, 675–679 (2013)





Principle of generative models



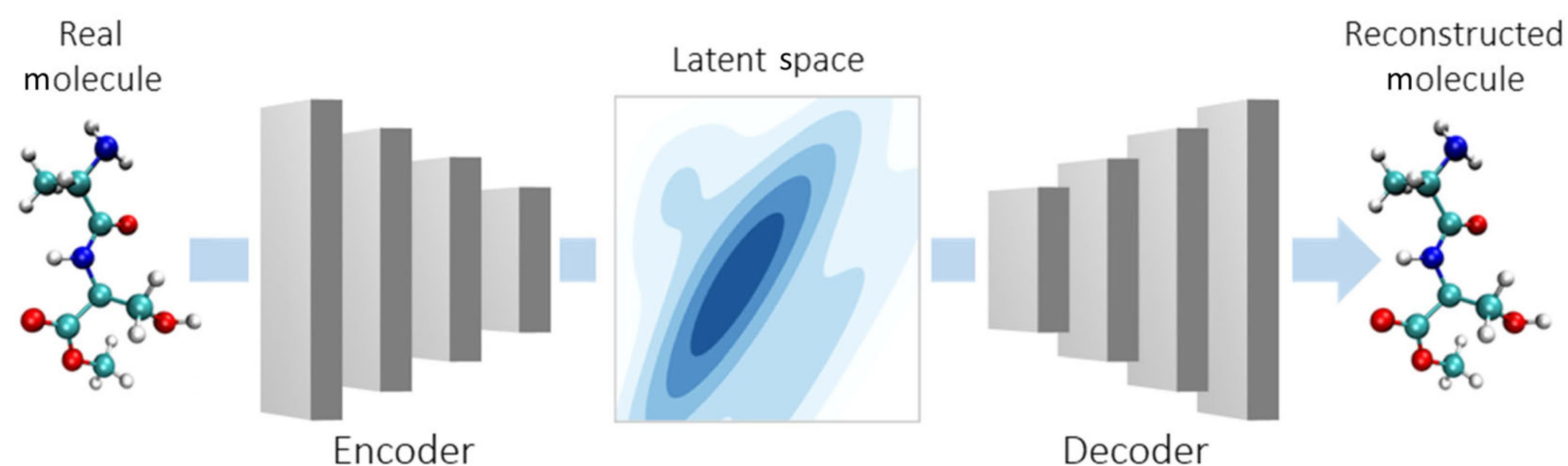
Credit: <https://openai.com/research/generative-models>

Example approaches to **generative models**:

- Variational autoencoders (VAEs)
- Generative adversarial networks (GANs)
- Generative flow networks (GFNs)
- Diffusion models (DMs)

Examples of generative models

Variational autoencoders

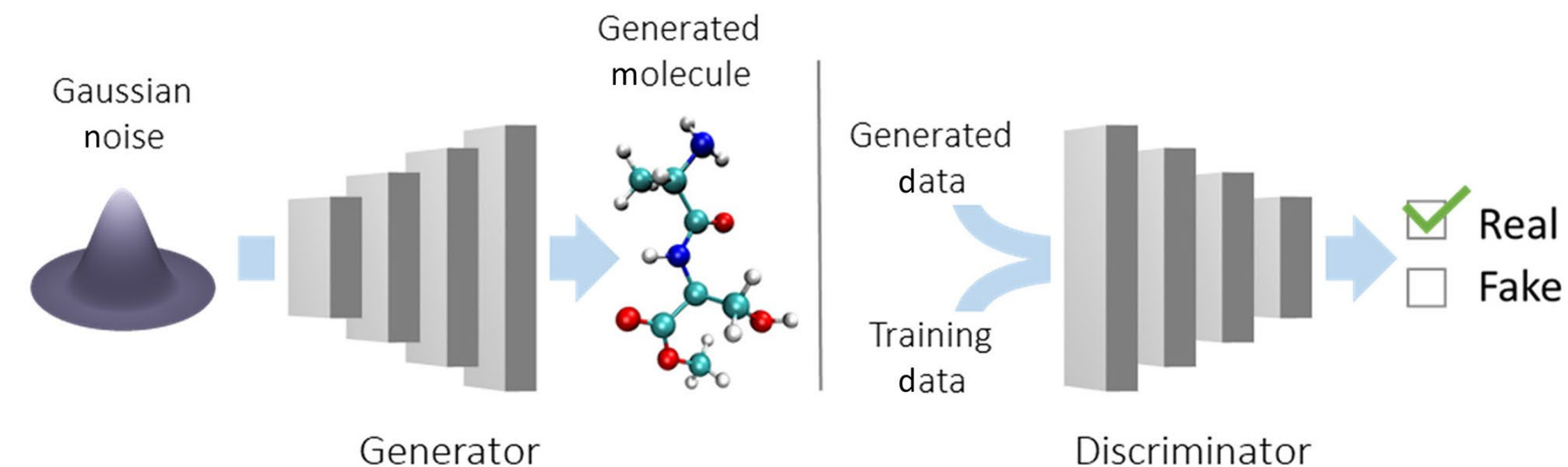


Credit: Bilodeau *et al.* (2022)

- Consist of an encoder and a decoder
- Learn the best encoding-decoding scheme through iterative optimization
- Loss function with two terms: reconstruction loss and regularization term

Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. *arXiv:1312.6114v11* (2022)

Generative adversarial networks

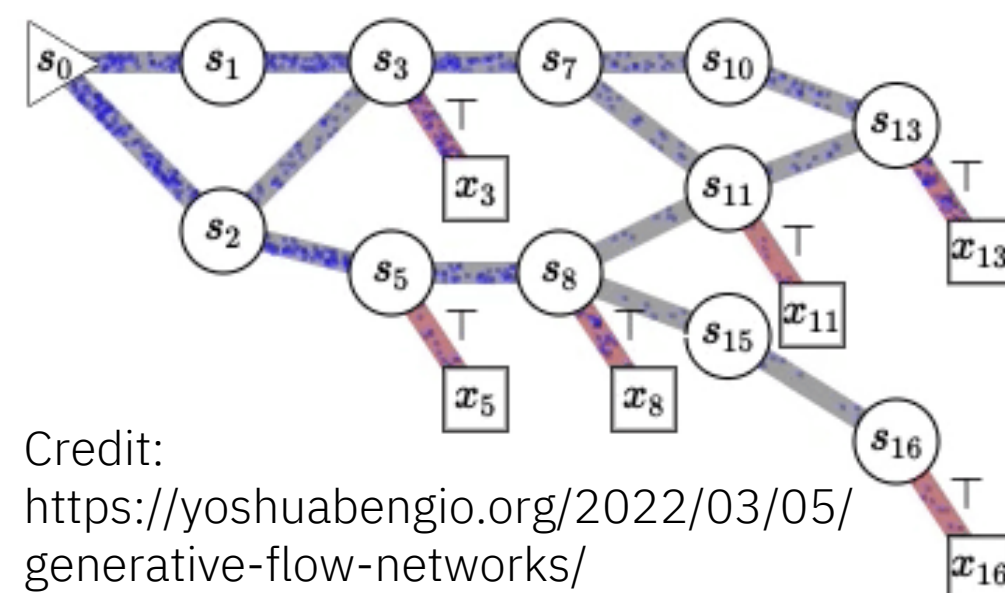


Credit: Bilodeau *et al.* (2022)

- Consist of a generator and a discriminator
- Adversarial training of generator to trick discriminator
- Tends to produce high fidelity, risk of collapsing to low diversity

Goodfellow, I. *et al.* Generative Adversarial Nets. *NeurIPS* (2014)

Generative flow networks

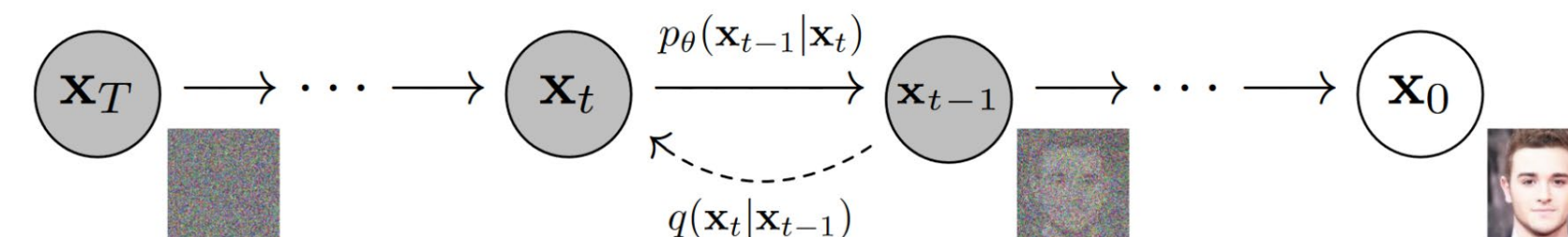


Credit:
<https://yoshuabengio.org/2022/03/05/generative-flow-networks/>

- Allow neural nets to model distributions over data structures like graphs
- Generates a series of actions at a frequency proportional to their reward
- Improves samples diversity and provides non-iterative sampling mechanism

Bengio, E., Jain, M., Korablyov, M., Precup, D. & Bengio, Y. Flow Network based Generative Models for Non-Iterative Diverse Candidate Generation. *NeurIPS* (2021)

Diffusion models

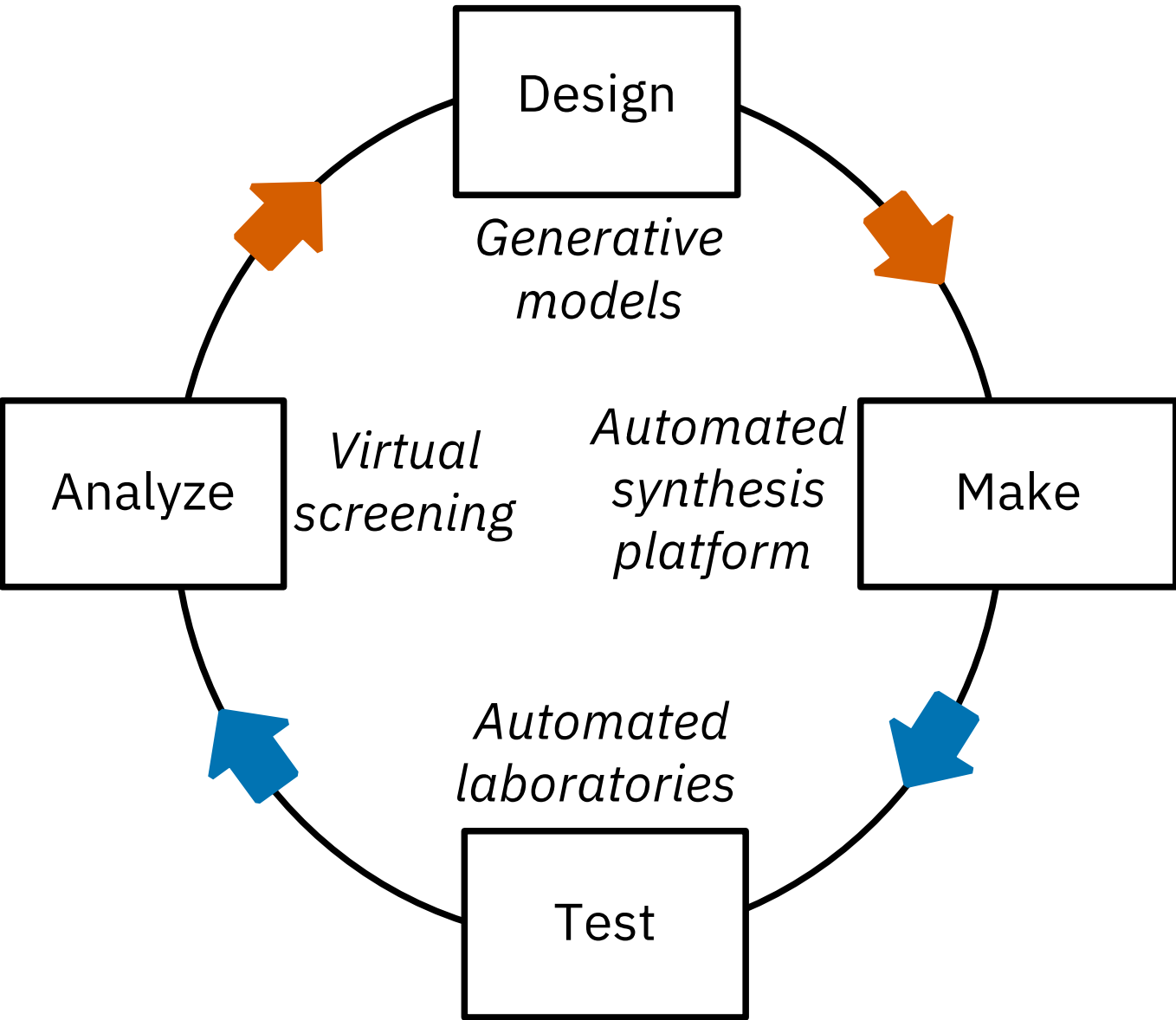


Credit: Ho *et al.* (2022)

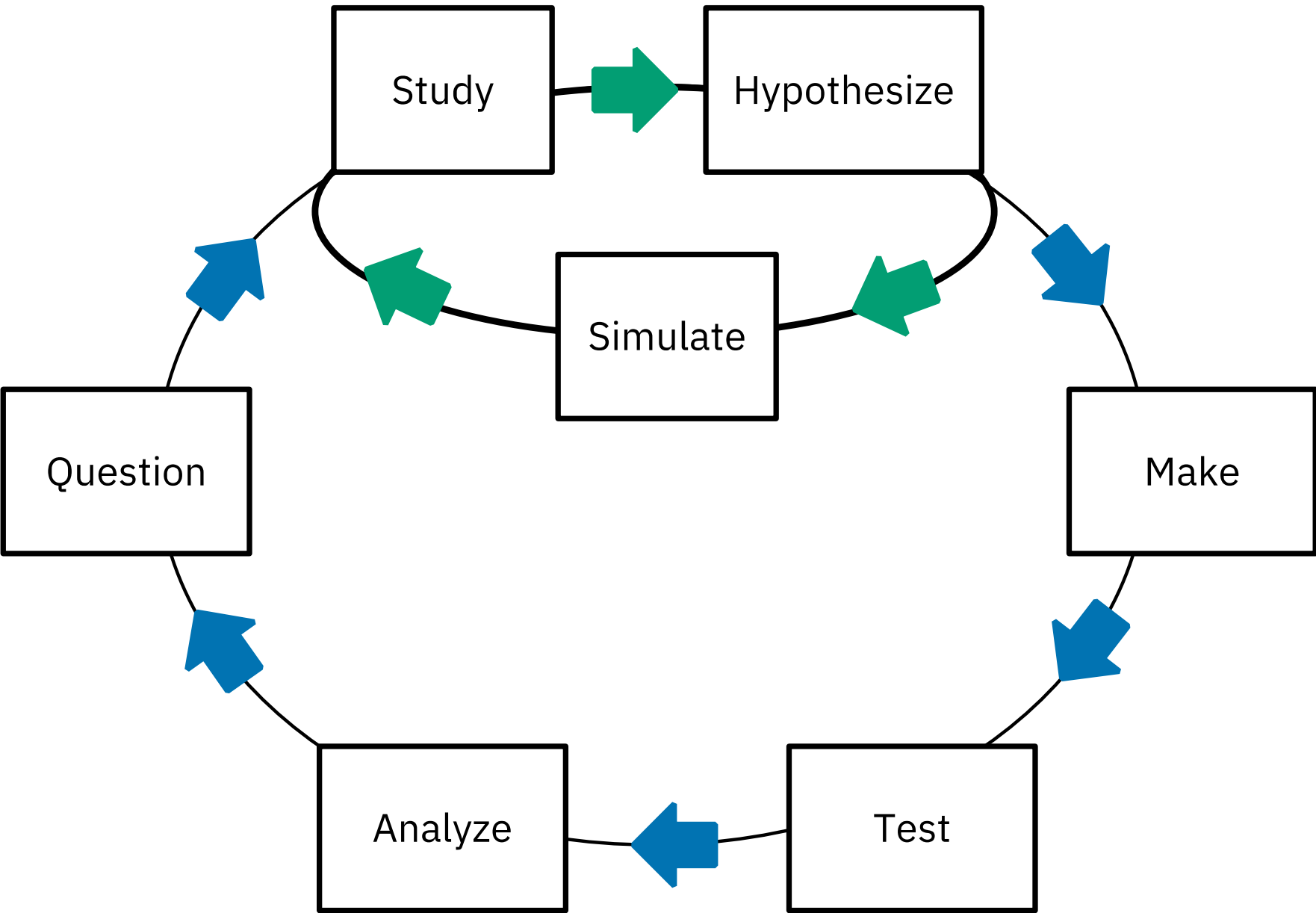
- Learn high-dimensional distributions by denoising data at multiple scales
- Fixed forward diffusion process and learnable reverse diffusion process
- High fidelity and diversity samples, but slow sample generation

Bengio, E., Jain, M., Korablyov, M., Precup, D. & Bengio, Y. Flow Network based Generative Models for Non-Iterative Diverse Candidate Generation. *NeurIPS* (2021)

Design-Make-Test-Analyze cycle in chemistry



Accelerated molecular discovery



Generative Toolkit for Scientific Discovery
(GT4SD): an open-source library to accelerate
hypothesis generation in scientific discovery



GT4SD makes generative AI algorithms and
models easier to use in scientific discovery
<https://github.com/GT4SD/gt4sd-core>

1. Train generative models

```
gt4sd-trainer --training_pipeline_name paccmann-vae-trainer --epochs 250
```

2. Create inference pipelines

```
gt4sd-saving --training_pipeline_name paccmann-vae-trainer --model_path /t
```

3. Run inference pipelines

```
gt4sd-inference --algorithm_name PaccMannGP --algorithm_application PaccMa
```

Manica *et al.*, *npj Comput. Mater.* **9**, 69 (2023)

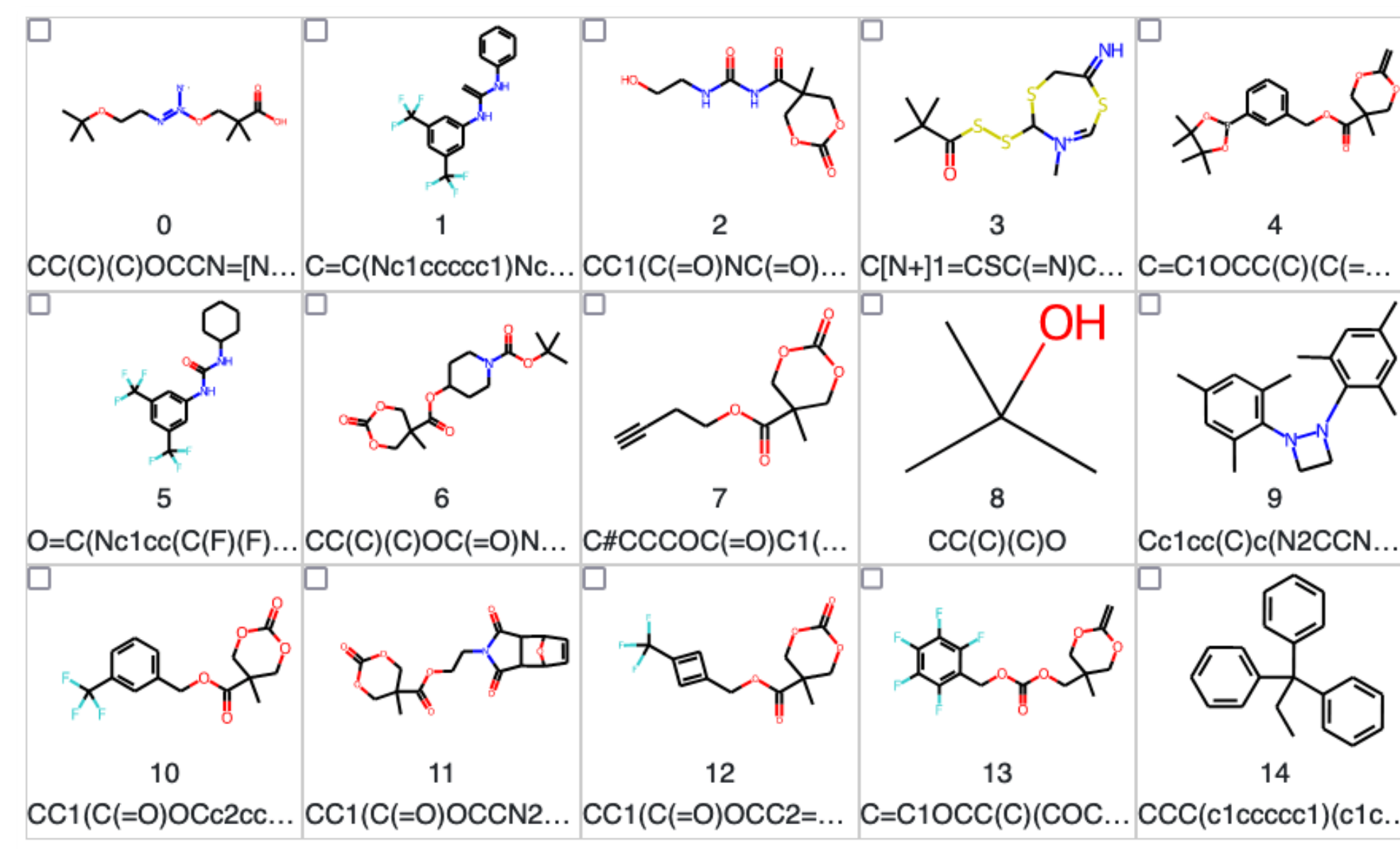
Study

Hypothesize

Test

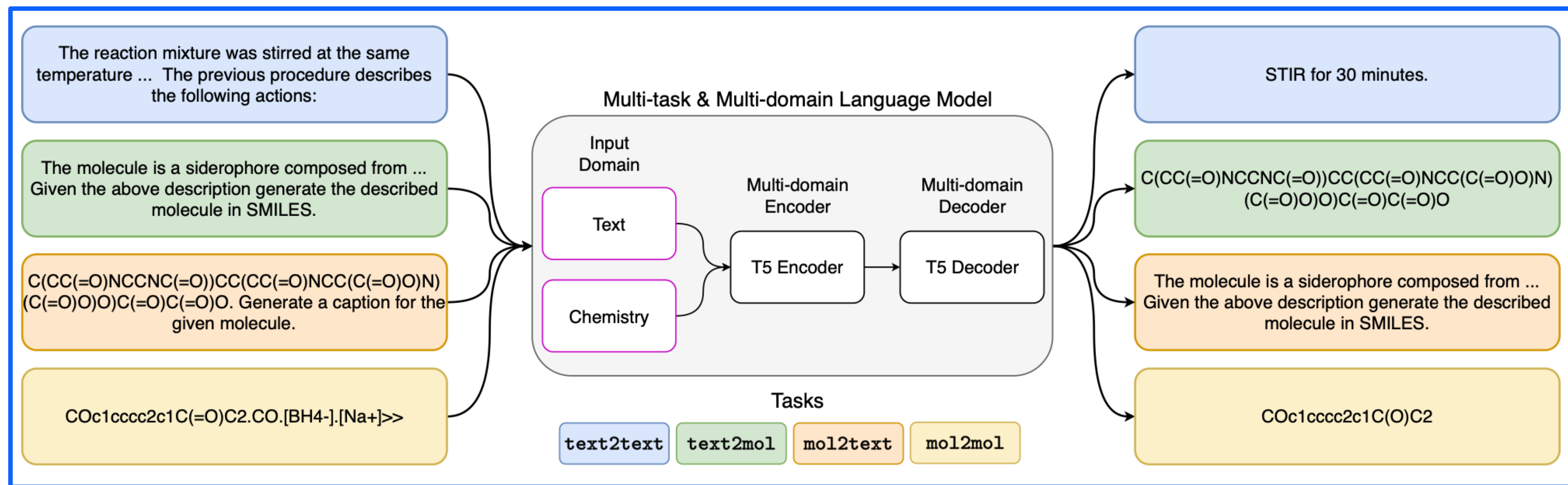
Applications include hypothesis generation for
inverse design and discovery of materials and
therapeutics like antivirals and antimicrobials

Example molecules generated using GT4SD



Multi-modal foundation models will accelerate fundamental research tasks

An example from chemistry...



Christofidellis *et al.*, ICML (2023)

End-to-end discovery workflow
using multi-task text and
chemistry T5 model

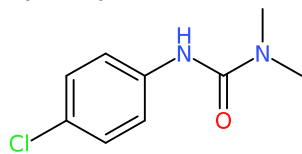
De novo design —→ Retrosynthesis —→ Actions

Input

Write in SMILES the described molecule: Give me a member of the class of phenylureas that is urea in which one of the nitrogens is substituted by a p-chlorophenyl group while the other is substituted by two methyl groups. It has a role as a herbicide, a xenobiotic and an environmental contaminant. It should be a member of monochlorobenzenes and a member of phenylureas.

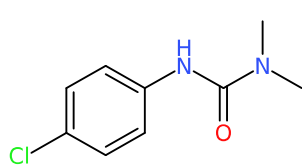
Target

CN(C)C(=O)NC1=CC=C(C=C1)Cl



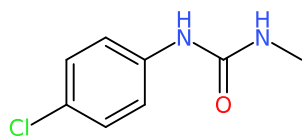
Text + chemistry T5

CN(C)C(=O)NC1=CC=C(C=C1)Cl



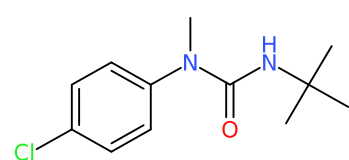
Galactica

1-(4-Chlorophenyl)-3-methylurea



ChatGPT

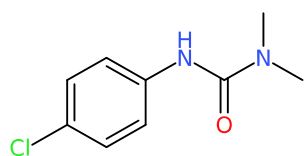
CC(NC(=O)N(C)C1=CC=C(Cl)C=C1)(C)C



Input

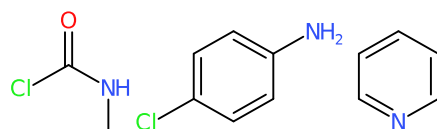
Predict the reaction that produces the following product:

CN(C)C(=O)NC1=CC=C(C=C1)Cl



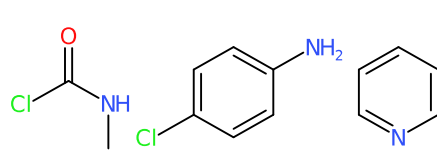
Target

CN(C)C(=O)Cl.Nc1ccc(Cl)cc1.c1ccncc1



Text + chemistry T5

CN(C)C(=O)Cl.NC1=CC=C(C=C1)Cl.c1ccncc1



RXN confidence = 1.0

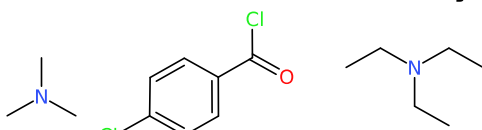
Galactica

(a) (b) (c) (d) (e) (f) (g) (h) (i) (j) (k) (l) (m) (n) (o) (p)
(q) (r) (s) (t) (u) (v) (w)

Invalid reaction

ChatGPT

CN(C)C + 4-chlorobenzoyl chloride + base



RXN confidence = 0.6

Input

Which actions are described in the following paragraph:

4-Chloro-aniline (1.0 g) was dissolved in pyridine and cooled to 0°C before adding dimethyl carbamyl chloride (0.683 g, 6.35 mmol). The reaction mixture was stirred at room temperature for 2 h and then heated overnight at 60°C under argon. The solution was cooled to room temperature, poured into ice water and extracted with EtOAc. The organic extract was dried over MgSO₄, filtered and concentrated to a residue to afford 3-(4-Chlorophenyl)-1,1-dimethylurea.

Target

- | | |
|--|---|
| 1. ADD 4-Chloroaniline (1.0 g) | 8. ADD ice water |
| 2. ADD Pyridine | 9. EXTRACT with EtOAc |
| 3. SET TEMPERATURE 0°C | 10. COLLECTLAYER organic |
| 4. ADD dimethyl carbamyl chloride (0.683 g, 6.35 mmol) | 11. DRY SOLUTION over MgSO ₄ |
| 5. STIR for 2 h at room temperature | 12. FILTER keep filtrate |
| 6. STIR for overnight at 60°C under argon | 13. CONCENTRATE |
| 7. SETTEMPERATURE room temperature | 14. YIELD 3-(4-Chlorophenyl)-1,1-dimethylurea |

Text + chemistry T5

- | | |
|------------------------------------|--|
| 1. ADD 4-Chloroaniline (1.0 g); | 7. ADD water; |
| 2. ADD pyridine; | 8. EXTRACT with EtOAc; |
| 3. SETTEMPERATURE 0°C; | 9. DRY SOLUTION over MGSO ₄ |
| 4. ADD dimethyl carbamyl chloride; | 10. FILTER keep filtrate; |
| 5. STIR for overnight at 60°C; | 11. CONCENTRATE; |
| 6. SETTEMPERATURE 60°C; | 12. YIELD product |

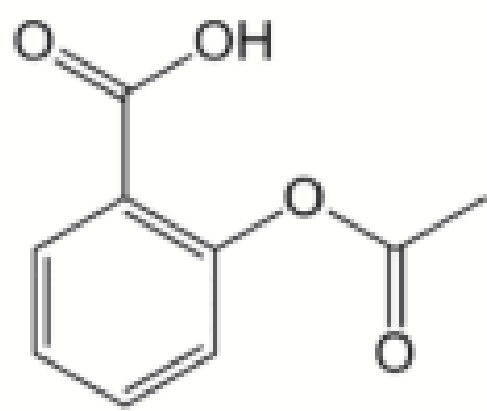
Galactica

- 3-(4-Chlorophenyl)-1,1-dimethylurea (1.0 g, 4.5 mmol) was dissolved in dry THF (10 mL) and cooled to 0° C. before adding n-BuLi (1.6 M in hexane, 1.6 mL, 2.5 mmol).
- The reaction mixture was stirred at 0° C. for 10 min and then at room temperature for 1 h.
- The reaction mixture was quenched with water and extracted with EtOAc.
- The organic extract was dried over MgSO₄, filtered and concentrated to a residue.
- The residue was purified by flash chromatography (silica gel, 10% EtOAc in hexane) to afford 3-(4-Chlorophenyl)-1,1-dimethylurea (1.0 g, 4.5 mmol) as a white solid (90%)

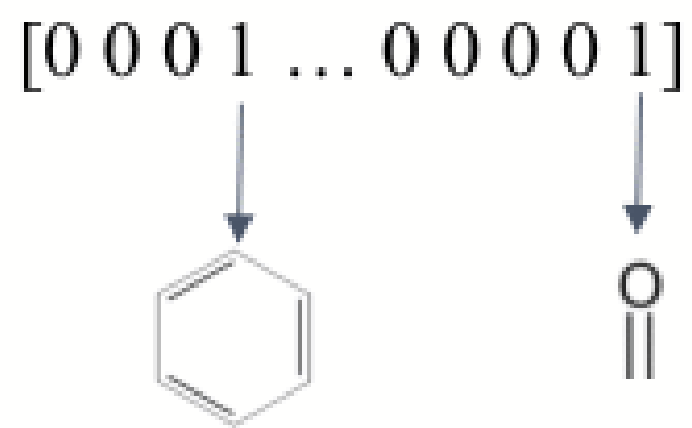
ChatGPT

- Dissolving 4-chloro-aniline in pyridine and cooling the mixture to 0°C.
- Adding dimethyl carbamyl chloride to the mixture.
- Stirring the reaction mixture at room temperature for 2 hours.
- Heating the reaction mixture overnight at 60°C under argon.
- Cooling the solution to room temperature and pouring it into ice water.
- Extracting the mixture with ethyl acetate.
- Drying the organic extract over magnesium sulfate.
- Filtering the extract and concentrating it to a residue.

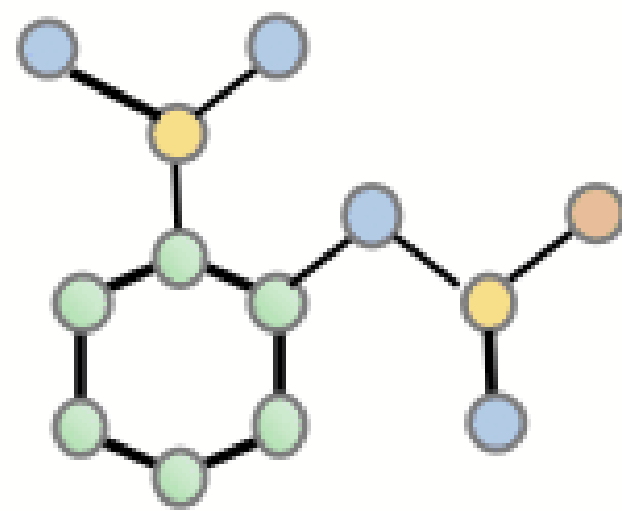
Different molecular representations for machine learning



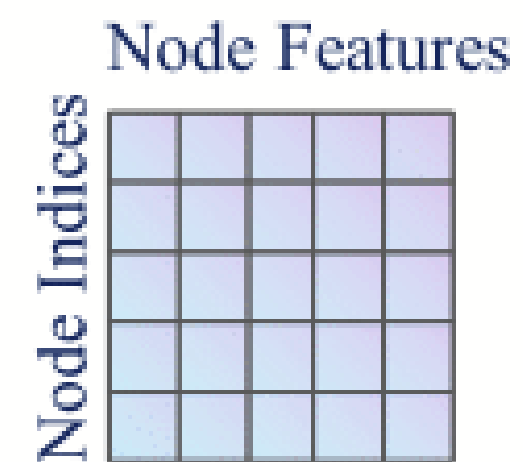
A. Kekulé Diagram



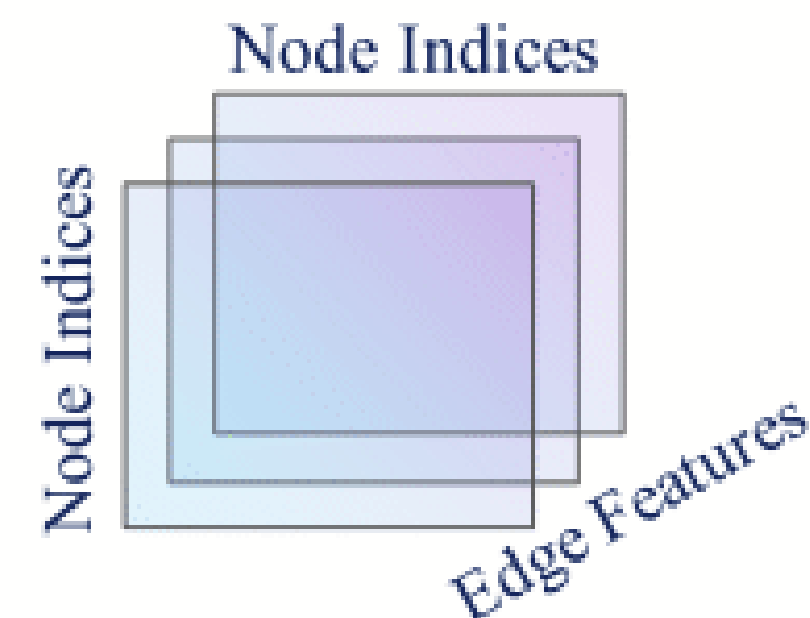
B. Fingerprints



C. Molecular Graph



Node Feature Matrix

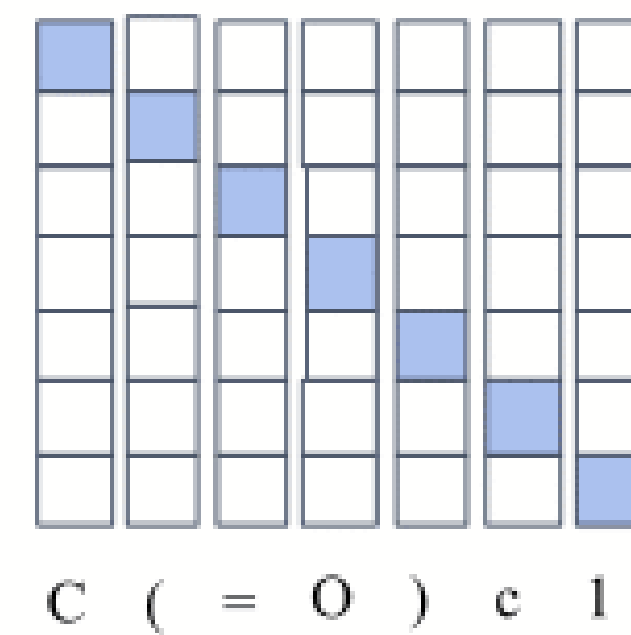


Adjacency Tensor



Tokenization

One-Hot Encoding



D. SMILES String

📖 Deng, J., Yang, Z., Ojima, I., Samaras, D. & Wang, F. Artificial intelligence in drug discovery: applications and techniques. *Briefings in Bioinformatics* **23**, bbab430 (2022)

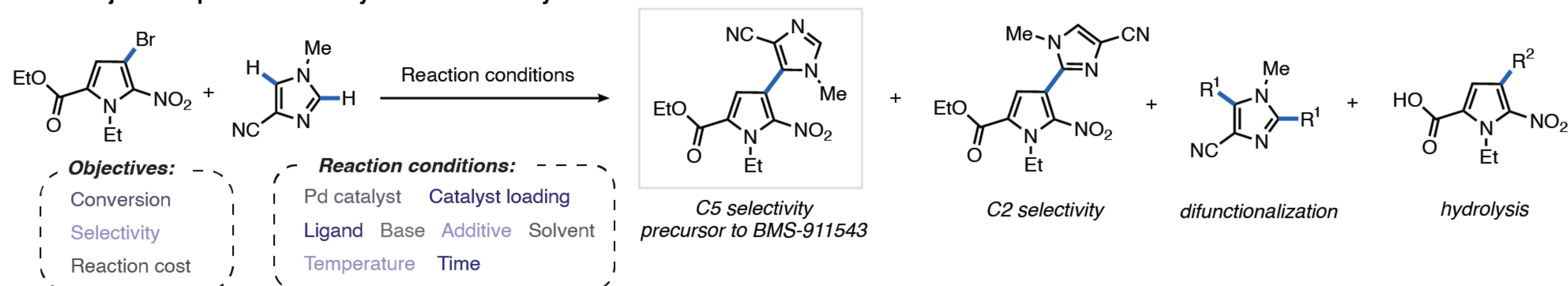
📖 Wigh, D. S., Goodman, J. M. & Lapkin, A. A. A review of molecular representation in the age of machine learning. *WIREs Computational Molecular Science* **12**, e1603 (2022)

Optimization strategies

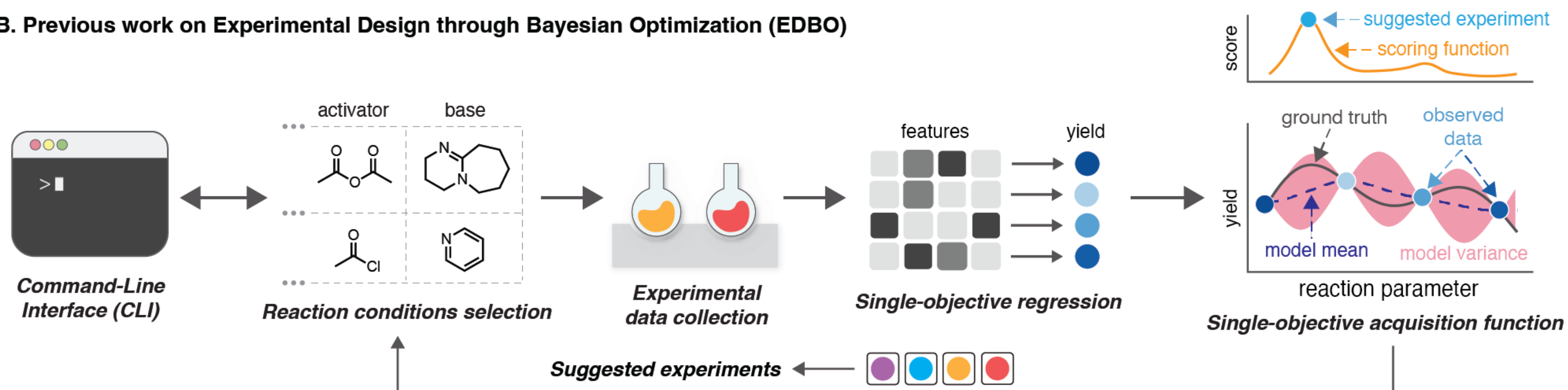
- Design of Experiment
- Bayesian Optimization
- Reinforcement Learning

Optimization in chemical reaction planning

A. Multi-objective optimization in synthetic chemistry



B. Previous work on Experimental Design through Bayesian Optimization (EDBO)



Torres, J. A. G. *et al.* A Multi-Objective Active Learning Platform and Web App for Reaction Optimization. *J. Am. Chem. Soc.* **144**, 19999–20007 (2022)

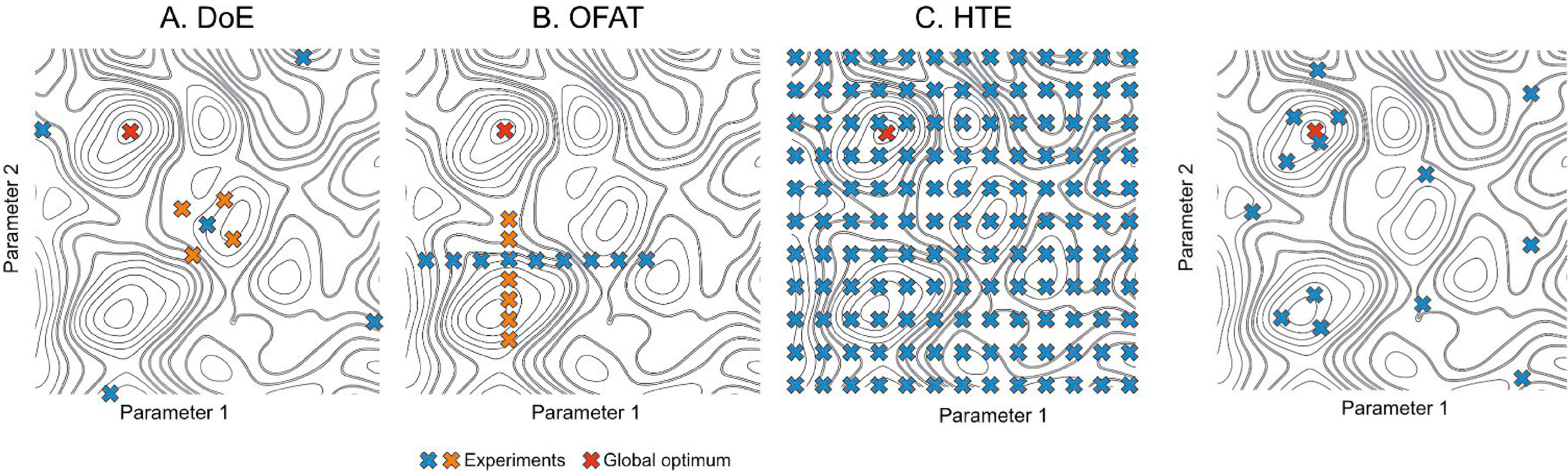
Comparison of approaches to experiment planning

Design of experiment

One-factor-at-a-time

High-throughput experimentation

Machine learning driven



- | | | | |
|--|--|---|--|
| <ul style="list-style-type: none">▪ Identify factors affecting response▪ Relies heavily on initial design▪ Number of experiments grows exponentially with number of parameters | <ul style="list-style-type: none">▪ Identify effect of single factor▪ Time consuming for many factors▪ Difficult to isolate confounding factors and potential interactions between factors | <ul style="list-style-type: none">▪ Grid of possible values▪ Exhaustive exploration▪ Experimentally expensive▪ Sensitive to choice of parameters | <ul style="list-style-type: none">▪ Adaptive methods that learn from experiments in real time▪ Recommend experiments where there is little information (exploration) or where better results are likely to occur (exploitation) |
|--|--|---|--|

📖 Gutierrez, D. P., Folkmann, L. M., Tribukait, H. & Roch, L. M. How to Accelerate R&D and Optimize Experiment Planning with Machine Learning and Data Science. *Chimie* **77**, 7–16 (2023)

Key concepts in Bayesian optimization

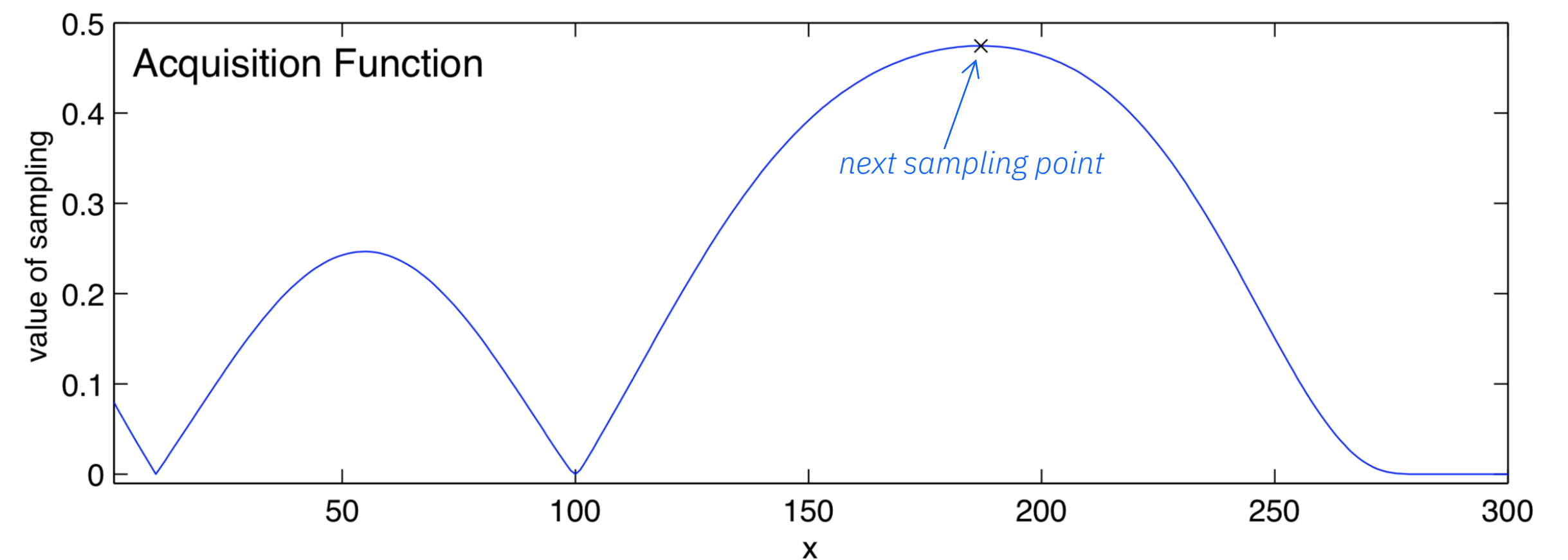
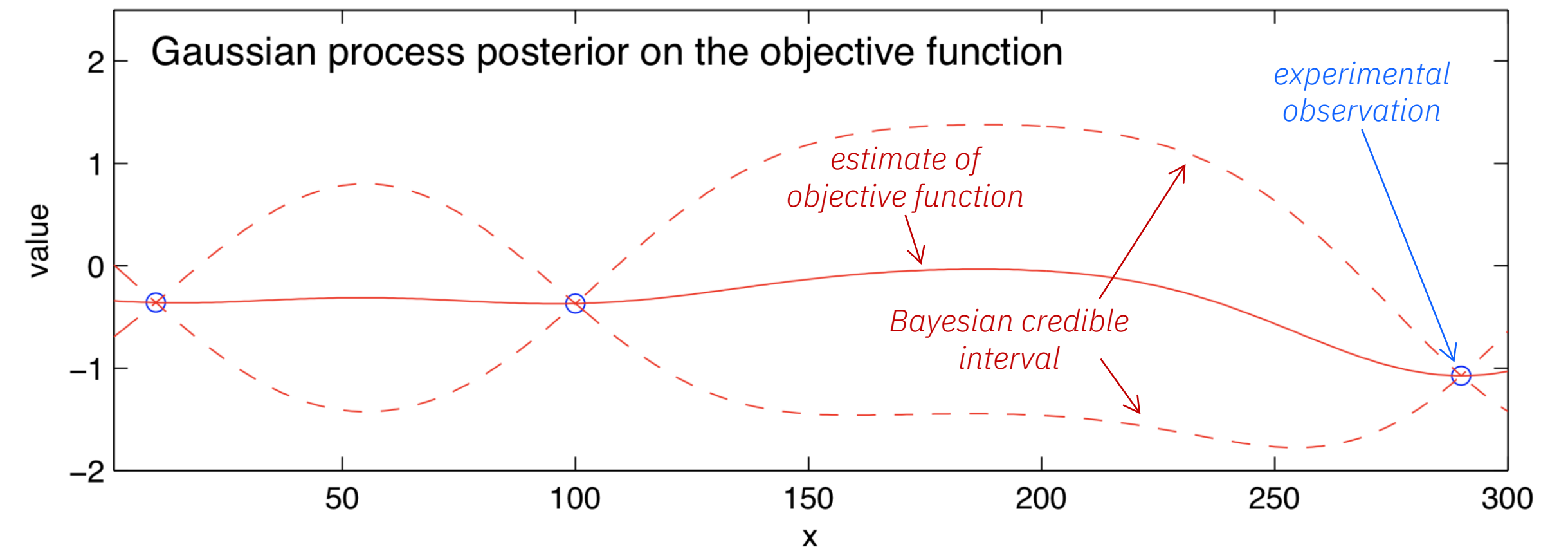
Bayesian optimization is a global optimization algorithm that reduces the need for many experiments.

Surrogate model:

- Probabilistic model to past observations
- Example function: Gaussian process

Acquisition model:

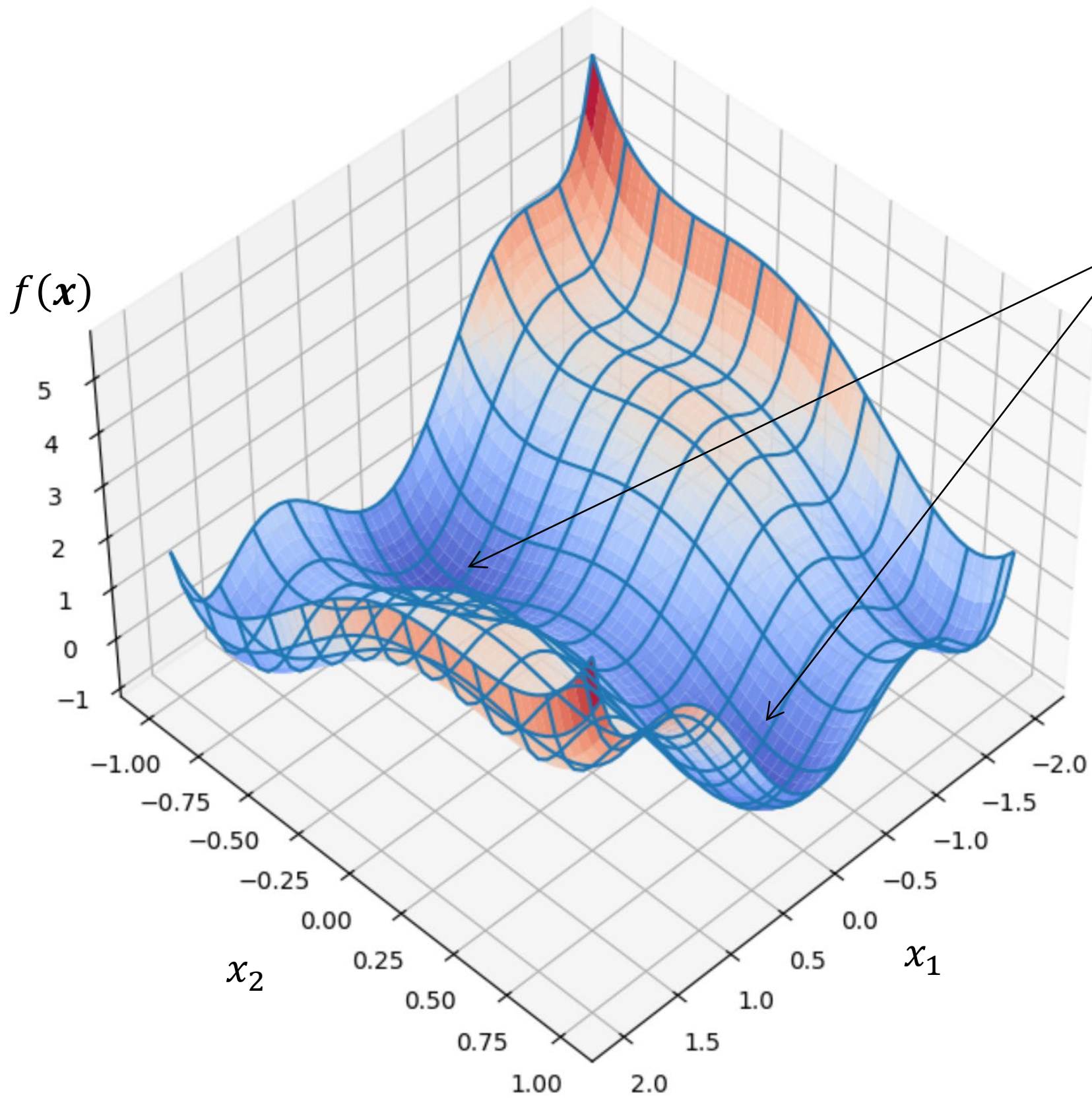
- Selection policy to decide which point to evaluate next
- Determines tradeoff between exploration and exploitation
- Example function: Expected improvement



 Frazier, P. I. A Tutorial on Bayesian Optimization. arXiv:1807.02811 (2018)

Six-hump camelback function

$$f(\mathbf{x}) = \left(4 - 2.1x_1^2 + \frac{x_1^4}{3}\right)x_1^2 + x_1x_2 + (-4 + 4x_2^2)x_2^2$$



Global minima:

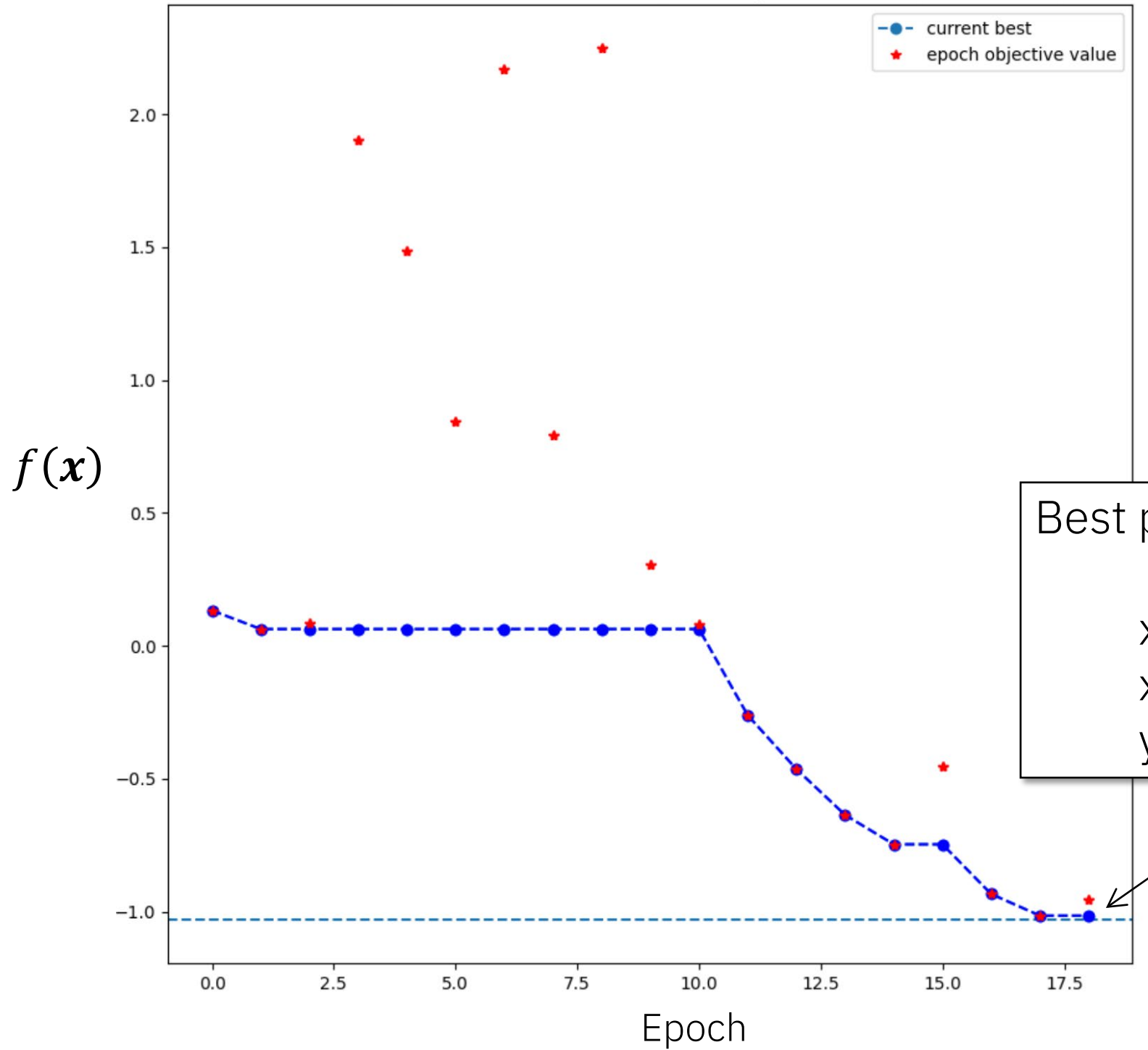
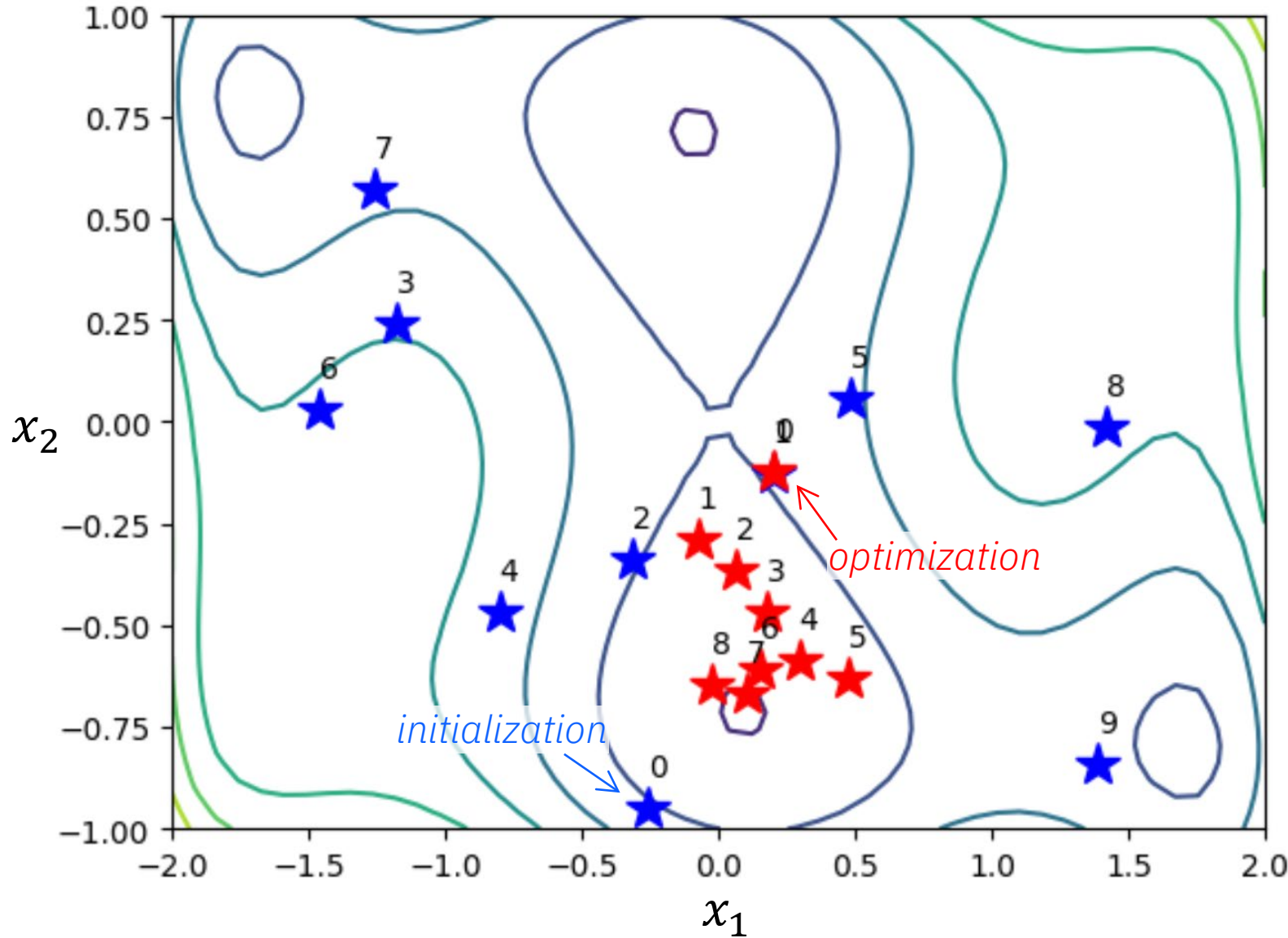
$x_1 = 0.0898$
 $x_2 = -0.7126$
 $y = -1.0316$

and

$x_1 = -0.0898$
 $x_2 = 0.7126$
 $y = -1.0316$

Other interactive examples online:
<https://distill.pub/2020/bayesian-optimization/>

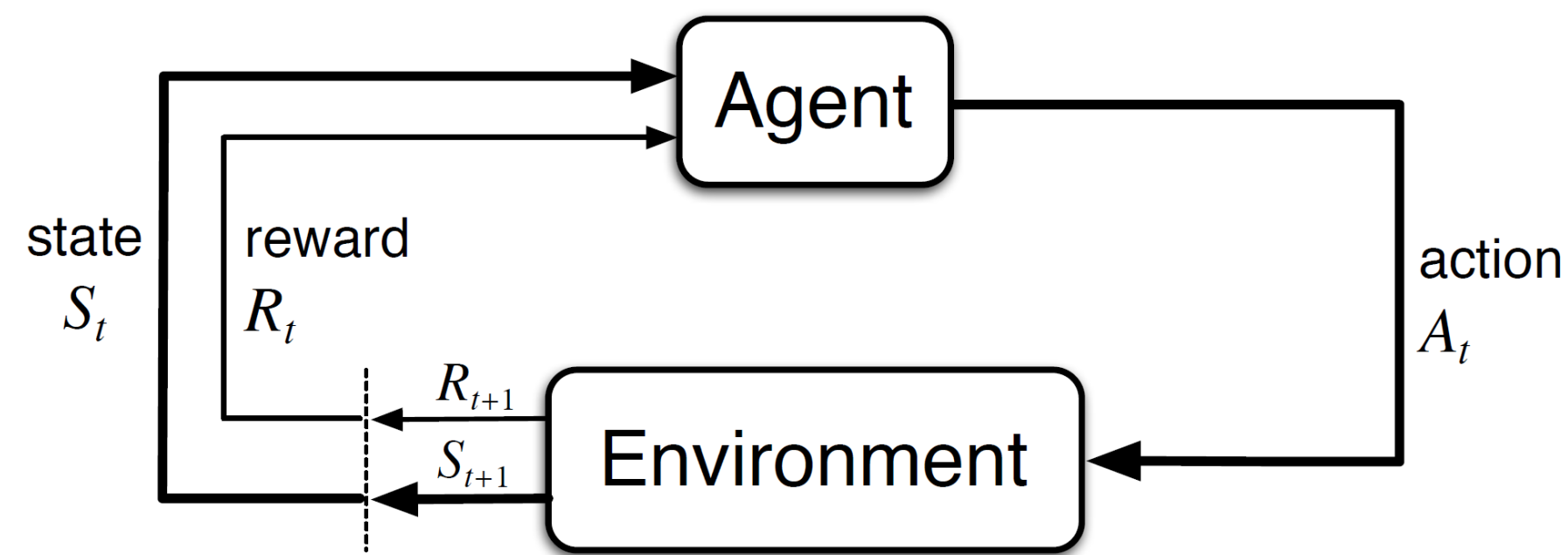
Exploration and exploitation in 10 epochs



Best parameters found:

$x_1 = 0.1020$
 $x_2 = -0.6682$
 $y = -1.0153$

Reinforcement learning (RL) to learn complex behaviors through trial-and-error interactions



📖 Sutton & Barton, Reinforcement Learning: An Introduction (2018)



PacMan agent using deep reinforcement learning

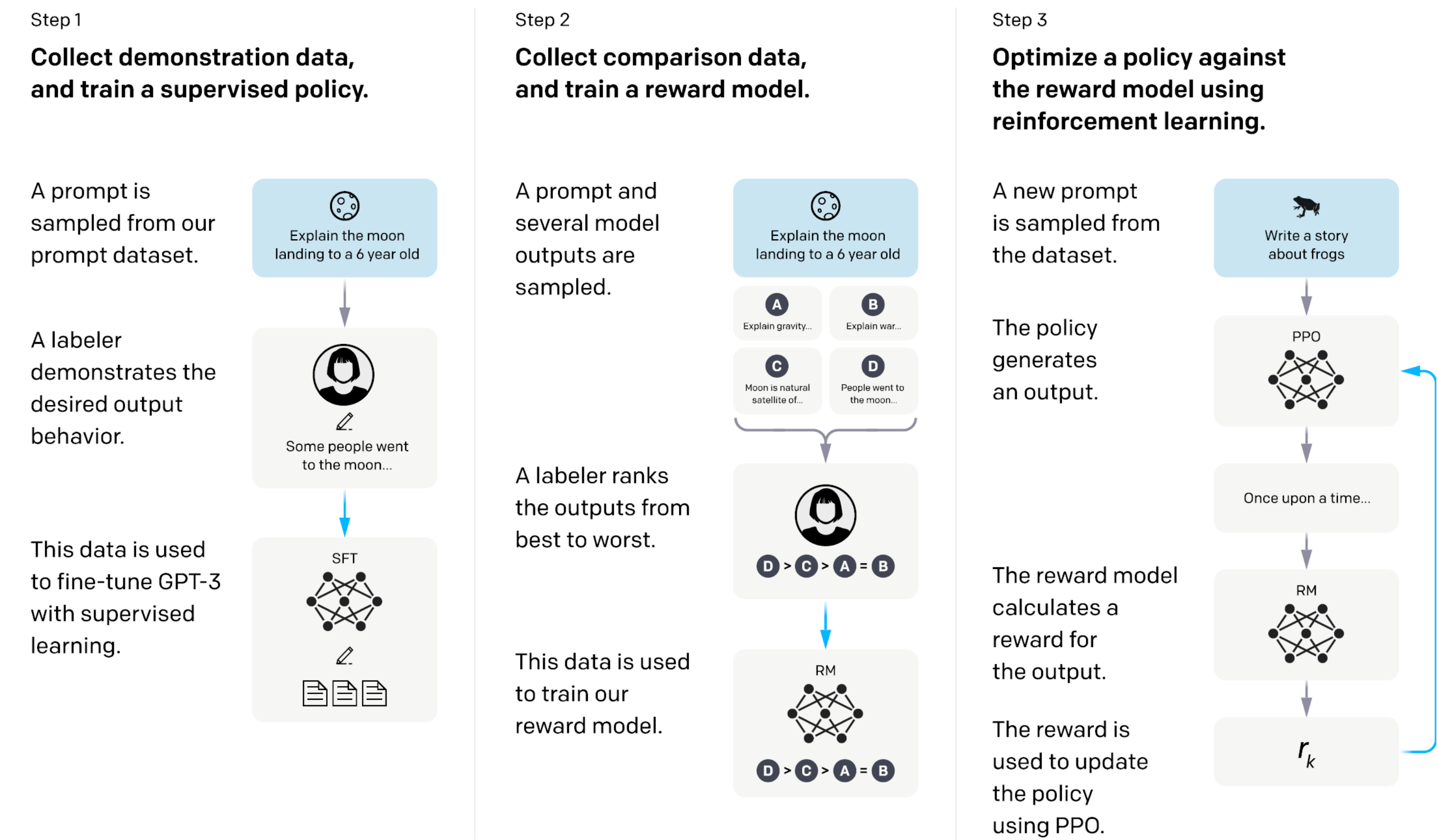
<https://youtu.be/QilHGSYbjDQ?si=wyP1NDdrymngUb7j>

- Reinforcement learning uses punishments and rewards as signals for positive and negative behavior
- Find suitable action model to maximize the total cumulative reward
- Exploration vs. exploitation trade-off
- Two main approaches in deep reinforcement learning:
 - Value functions (estimate expected return, e.g. through neural networks)
 - Policy search (directly find policies, e.g. through gradient-free or gradient-based methods)

📖 Arulkumaran, K., Deisenroth, M. P., Brundage, M. & Bharath, A. A. Deep Reinforcement Learning: A Brief Survey. *IEEE Signal Processing Magazine* **34**, 26–38 (2017)

📖 Azar, M. G. *et al.* A General Theoretical Paradigm to Understand Learning from Human Preferences. arXiv:2310.12036 (2023)

Reinforcement learning with human feedback (RLHF) to align model with user intent



- Common approach to RLHF is based on training a reward model (classifier) and then fine-tuning a policy to maximize the reward
- Recent approaches bypass reward models to directly train policies through Direct Preference Optimization (DPO) or pairwise preferences

📖 Rafailov, R. *et al.* Direct Preference Optimization: Your Language Model is Secretly a Reward Model. arXiv:2305.18290 (2023)

📖 Azar, M. G. *et al.* A General Theoretical Paradigm to Understand Learning from Human Preferences. arXiv:2310.12036 (2023)

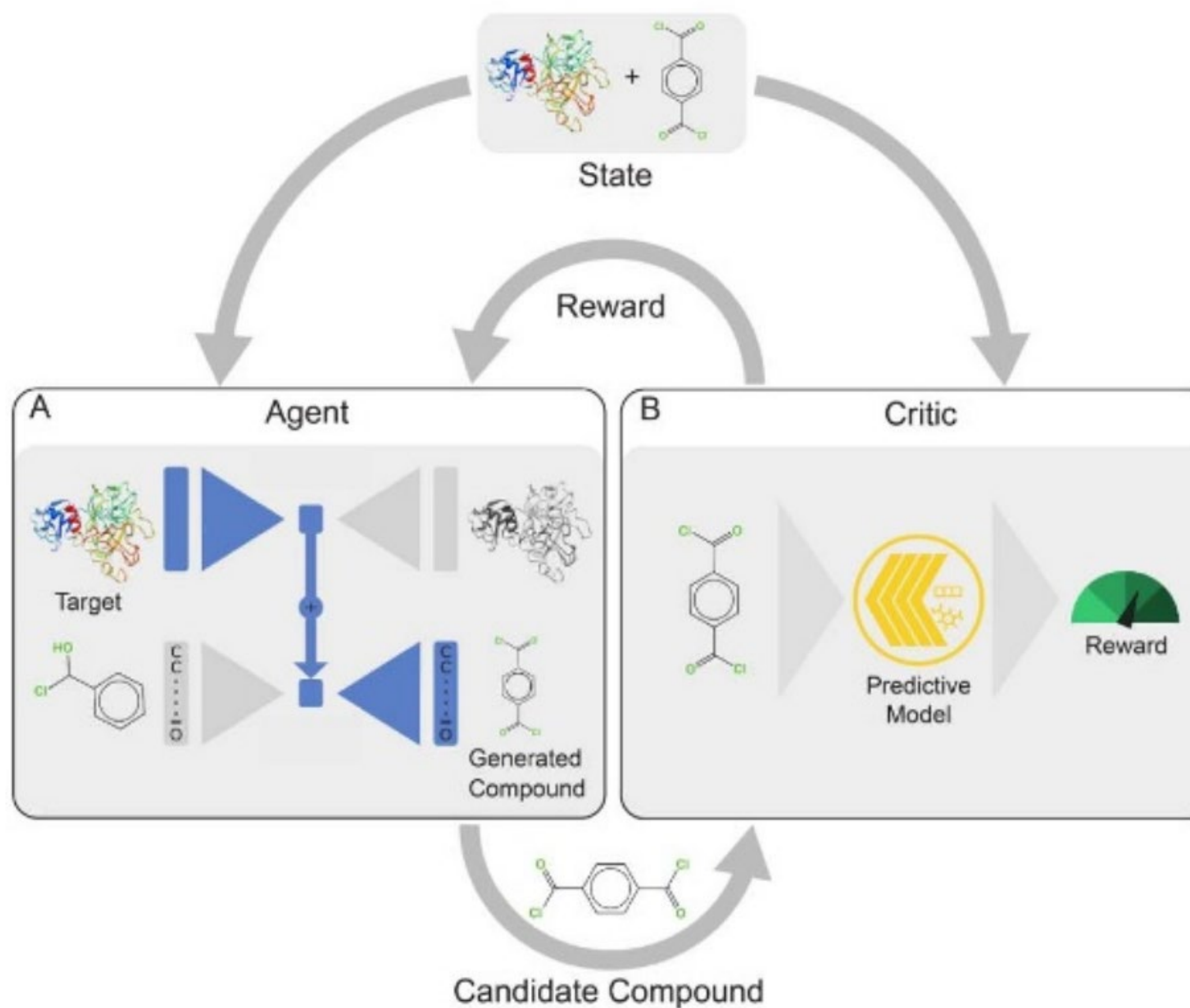
📖 Ouyang, L. *et al.* Training language models to follow instructions with human feedback. arXiv:2203.02155 (2022)

Closing the loop

- Case studies

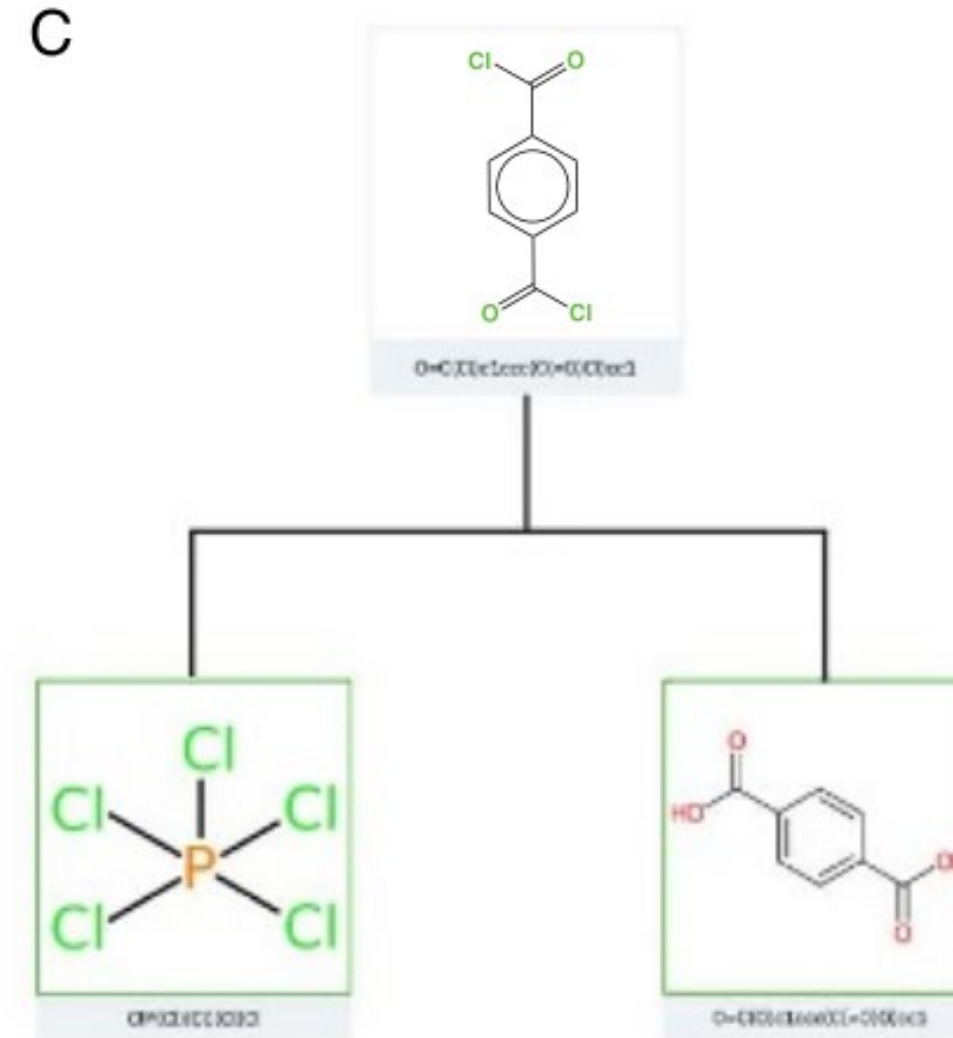
Toward closed-loop
autonomous molecular
discovery: designing novel
antiviral candidates against
SARS-CoV-2

Conditional generative model



Retrosynthesis prediction model

C



Predicted synthesis plan

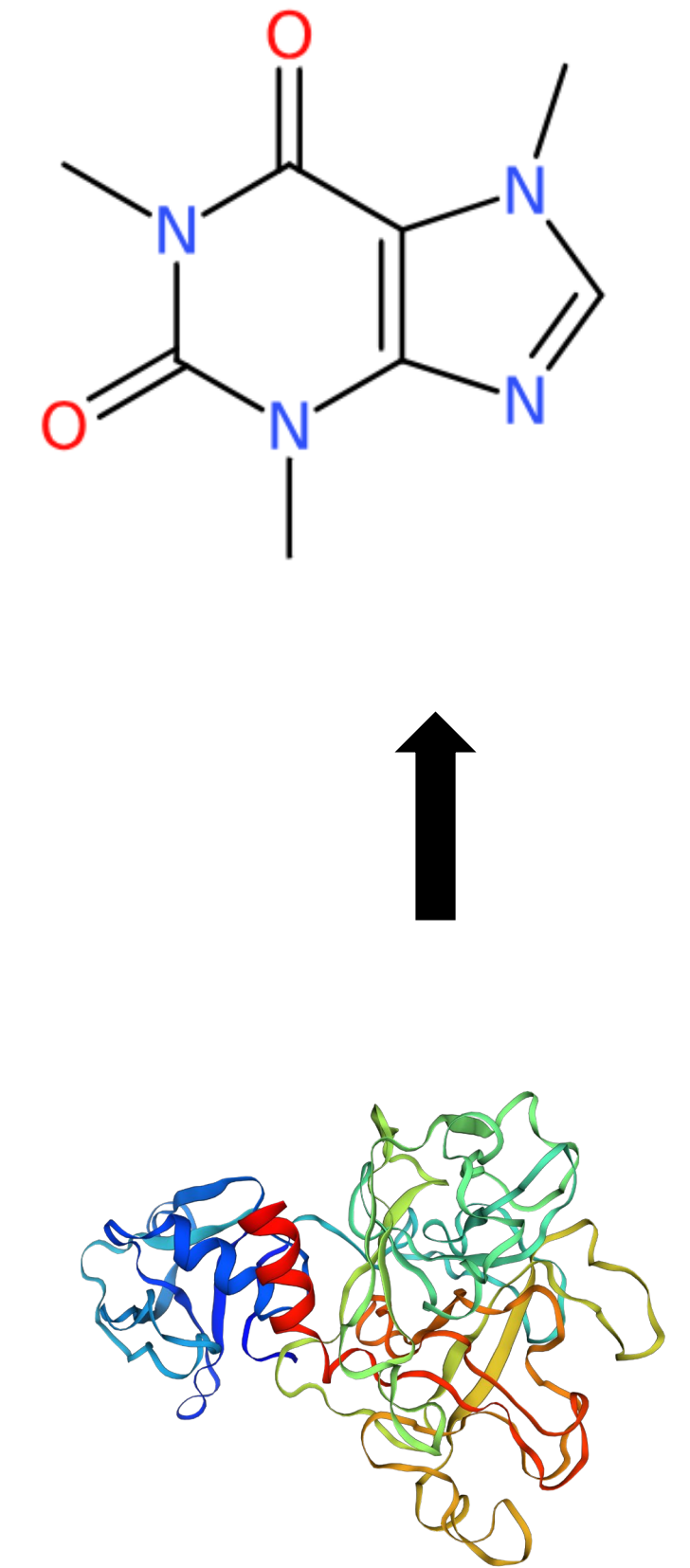
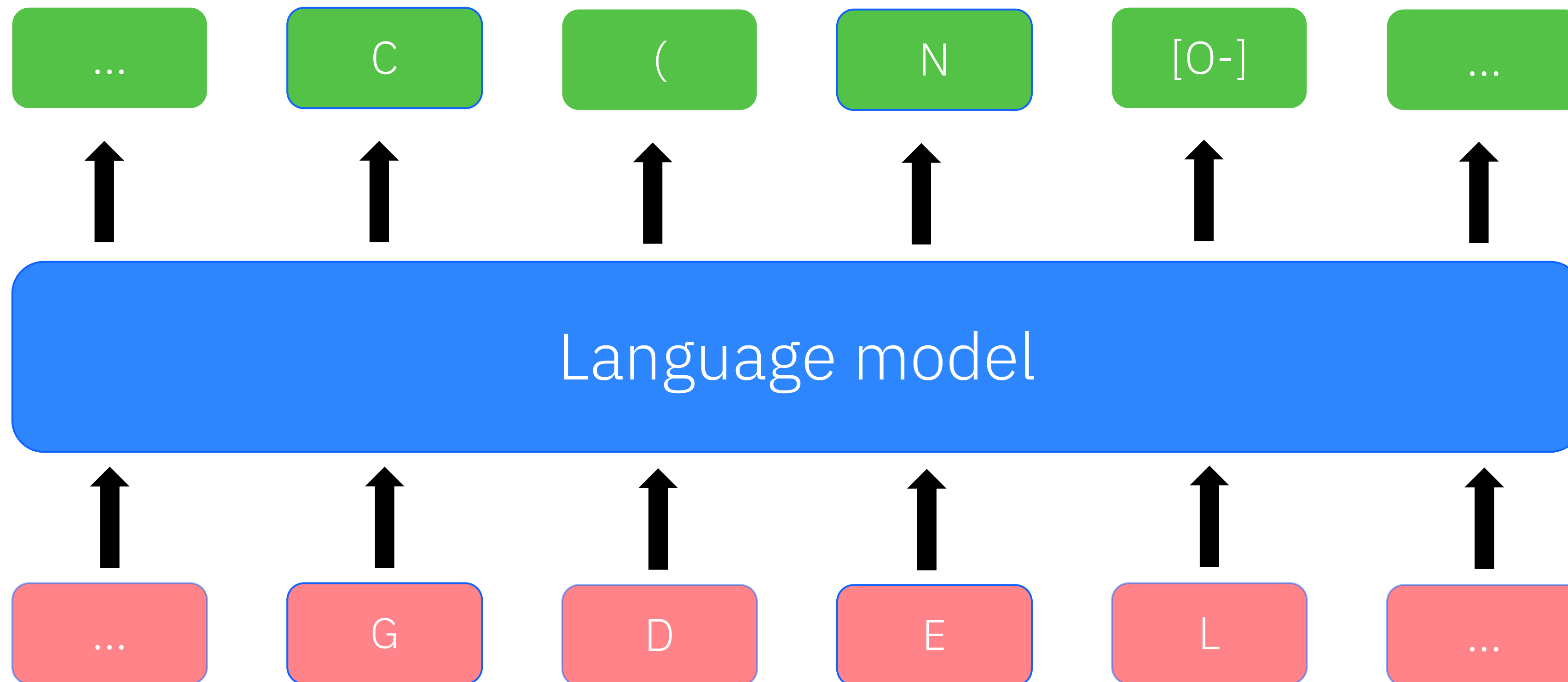
D

- 1) ADD ClP(Cl)(Cl)(Cl)Cl
- 2) ADD O=C(O)c1ccc(C(=O)O)cc1
- 3) STIR at 60 °C for 70 minutes
- 4) QUENCH
- 5) EXTRACT with Et2O
- 6) EVAPORATE
- 7) PURIFY



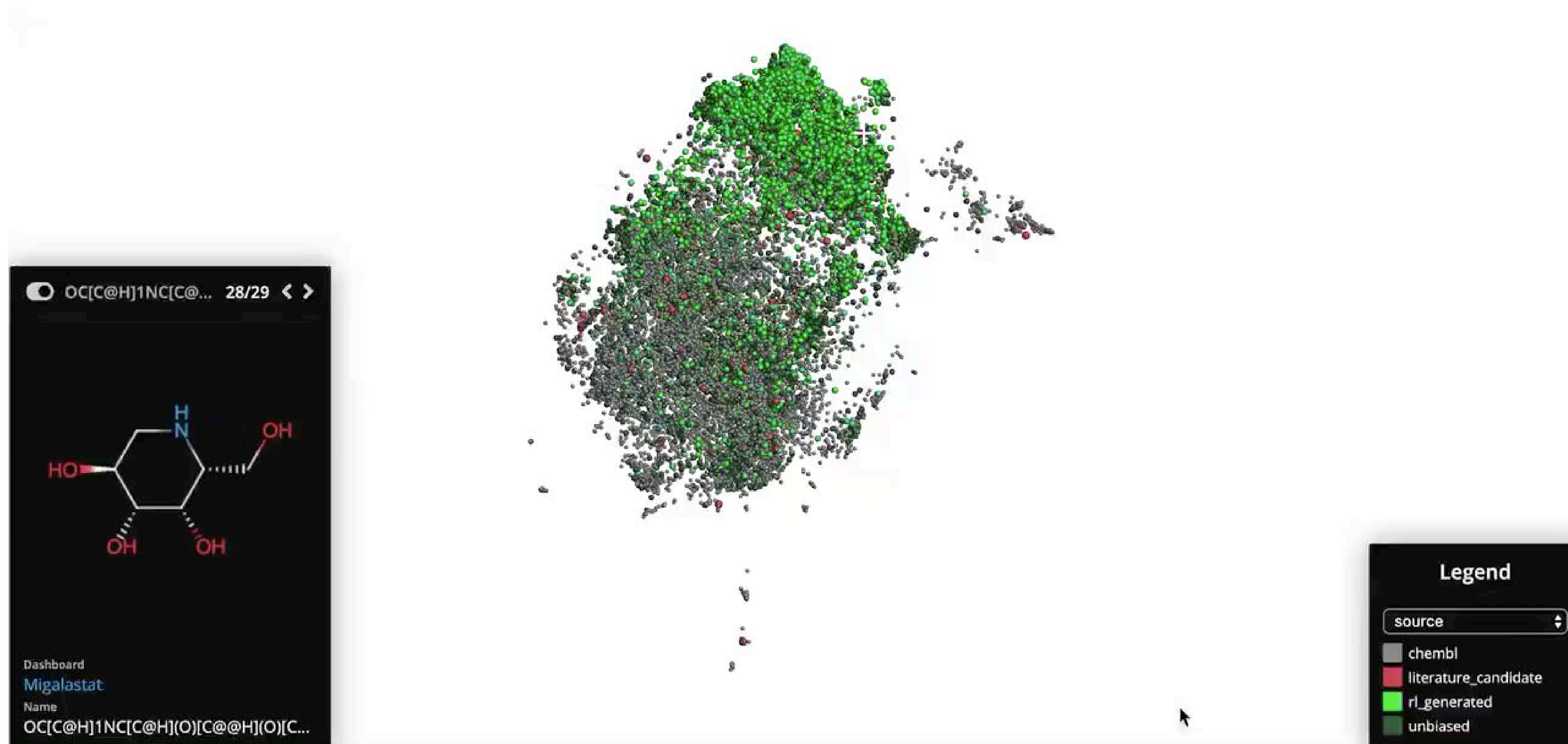
Candidate compound synthesis

Translate from amino acid
sequence to potentially binding
ligands



 Born, J. *et al.* Data-driven molecular design for discovery and synthesis of novel ligands: a case study on SARS-CoV-2. *Mach. Learn.: Sci. Technol.* **2**, 025024 (2021)

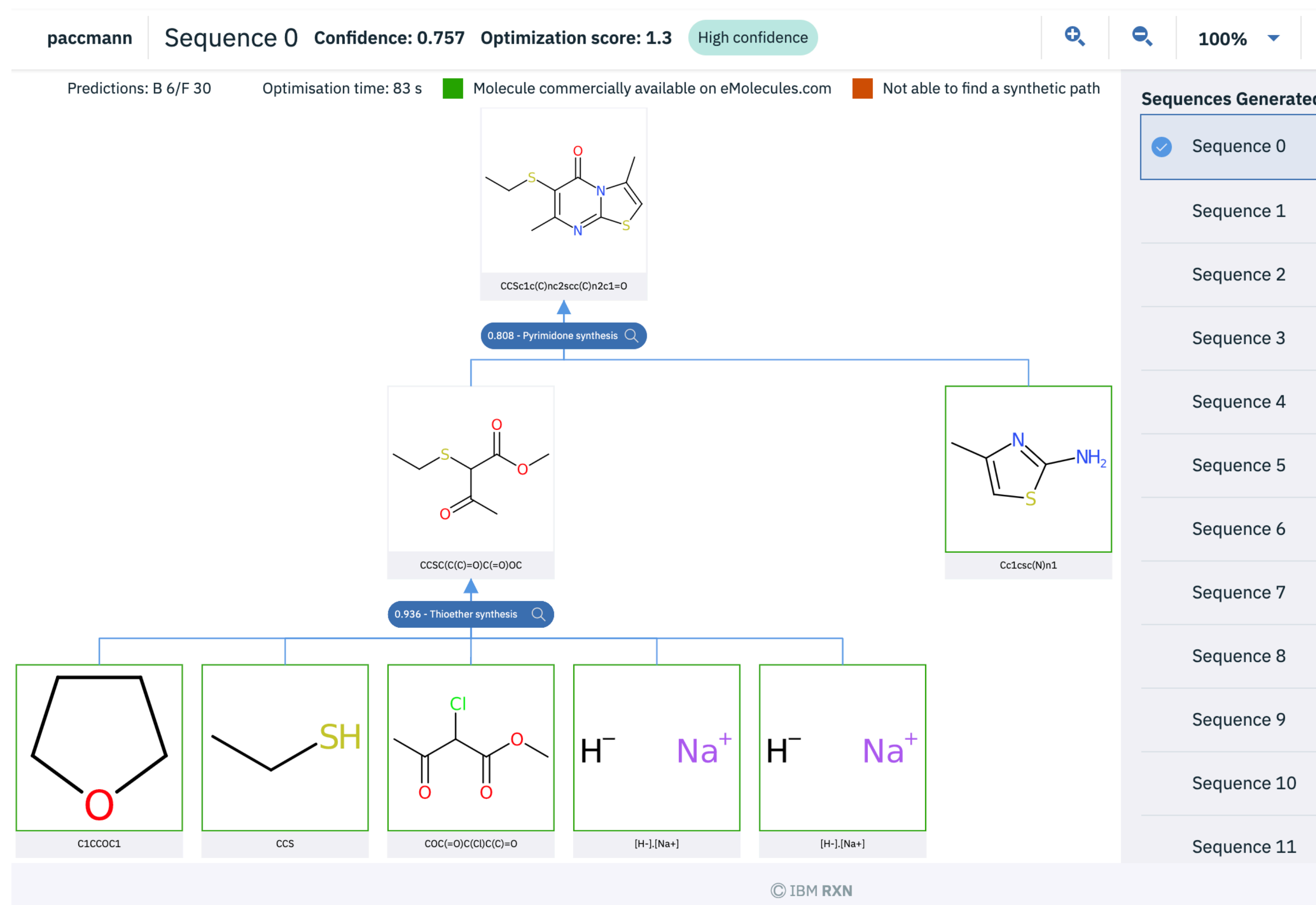
Visualizing the learned chemical space



📖 Born, J. *et al.* Data-driven molecular design for discovery and synthesis of novel ligands: a case study on SARS-CoV-2. *Mach. Learn.: Sci. Technol.* **2**, 025024 (2021)

Find a viable synthetic route

- Web service:
rxn.res.ibm.com
- Retrosynthesis route for most promising compounds
- Reaction prediction and retro-synthesis models (Transformer)
- Synthesis routes found for ca. 30% of the best candidates



📖 Born, J. *et al.* Data-driven molecular design for discovery and synthesis of novel ligands: a case study on SARS-CoV-2. *Mach. Learn.: Sci. Technol.* **2**, 025024 (2021)

Selection of synthesis candidate

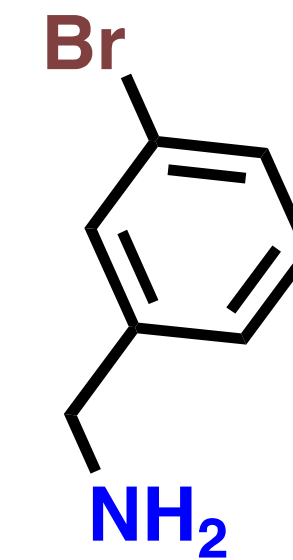
3-Bromobenzylamine:

- Full substructure of Arbidol
- Generated by our model to target ACE2 receptor
- Presence of bromine is key for efficacy of Arbidol

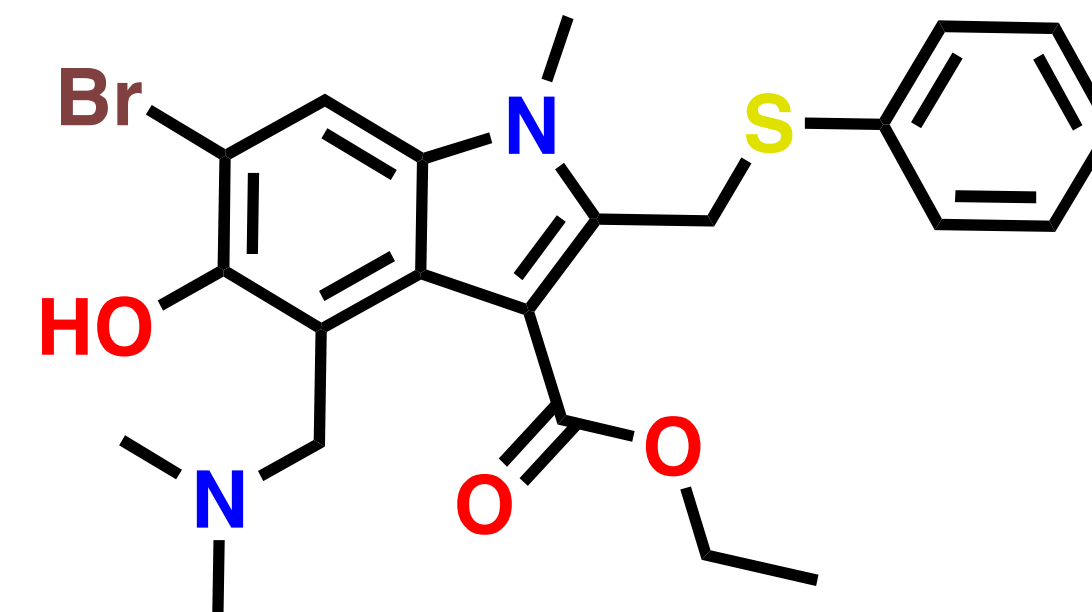
Arbidol:

- Approved, broad-spectrum antiviral drug
- Positive evidence for COVID-19 treatment
- Interacts with the ACE2 receptor

3-Bromobenzylamine



Arbidol

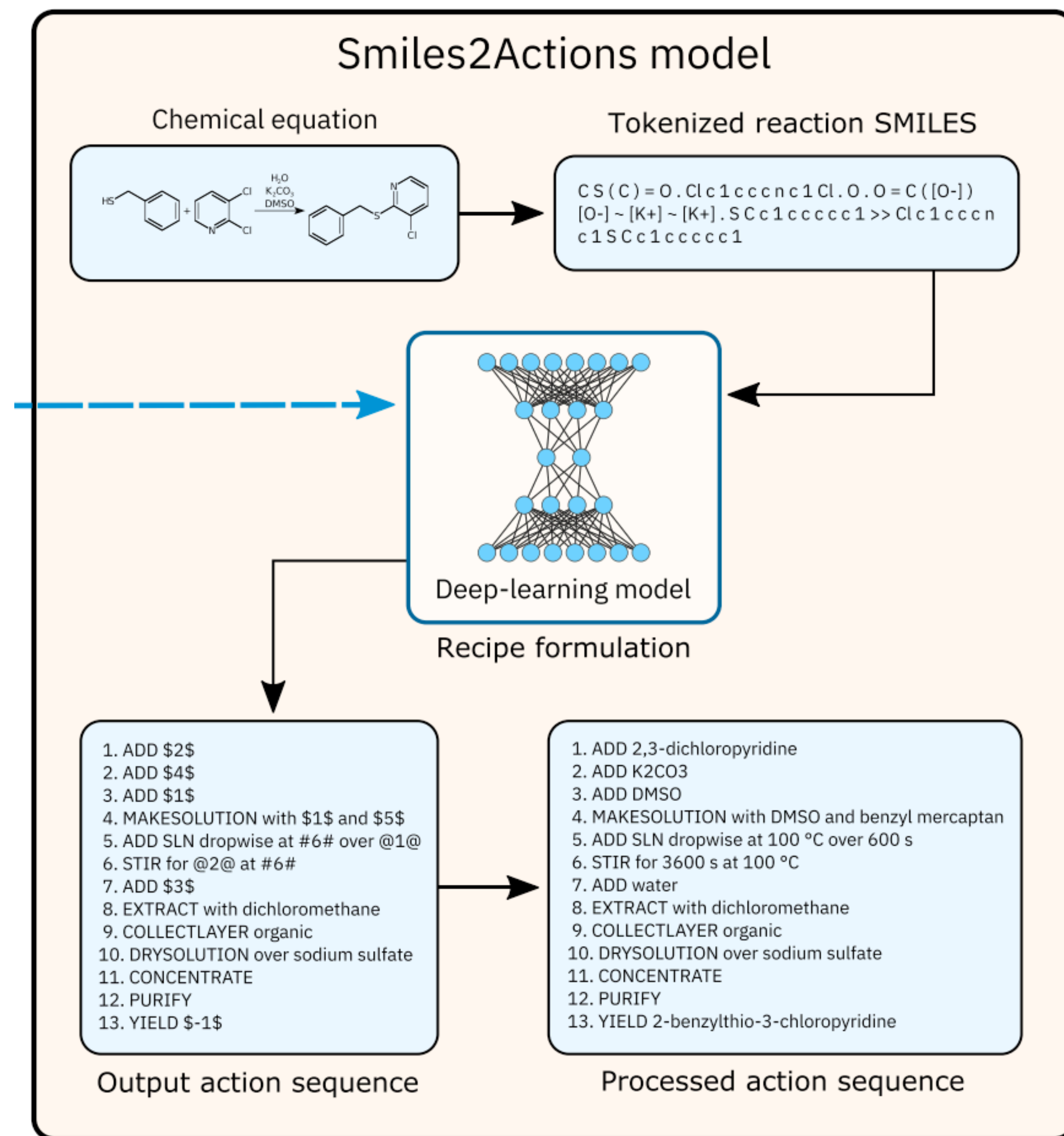
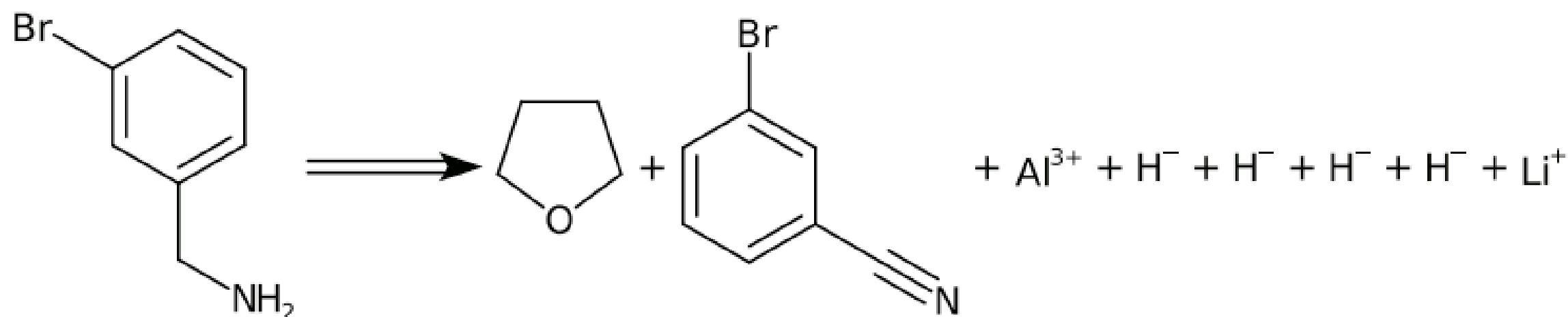


 Born, J. *et al.* Data-driven molecular design for discovery and synthesis of novel ligands: a case study on SARS-CoV-2. *Mach. Learn.: Sci. Technol.* **2**, 025024 (2021)

Stepwise synthesis execution plan

Type: Nitrile reduction, Confidence: 0.985

C1CCOC1.N#Cc1cccc(Br)c1.[Al+3].[H-].[H-].[H-].[H-].[Li+]>>BrC1C=CC=C(CN)C=1



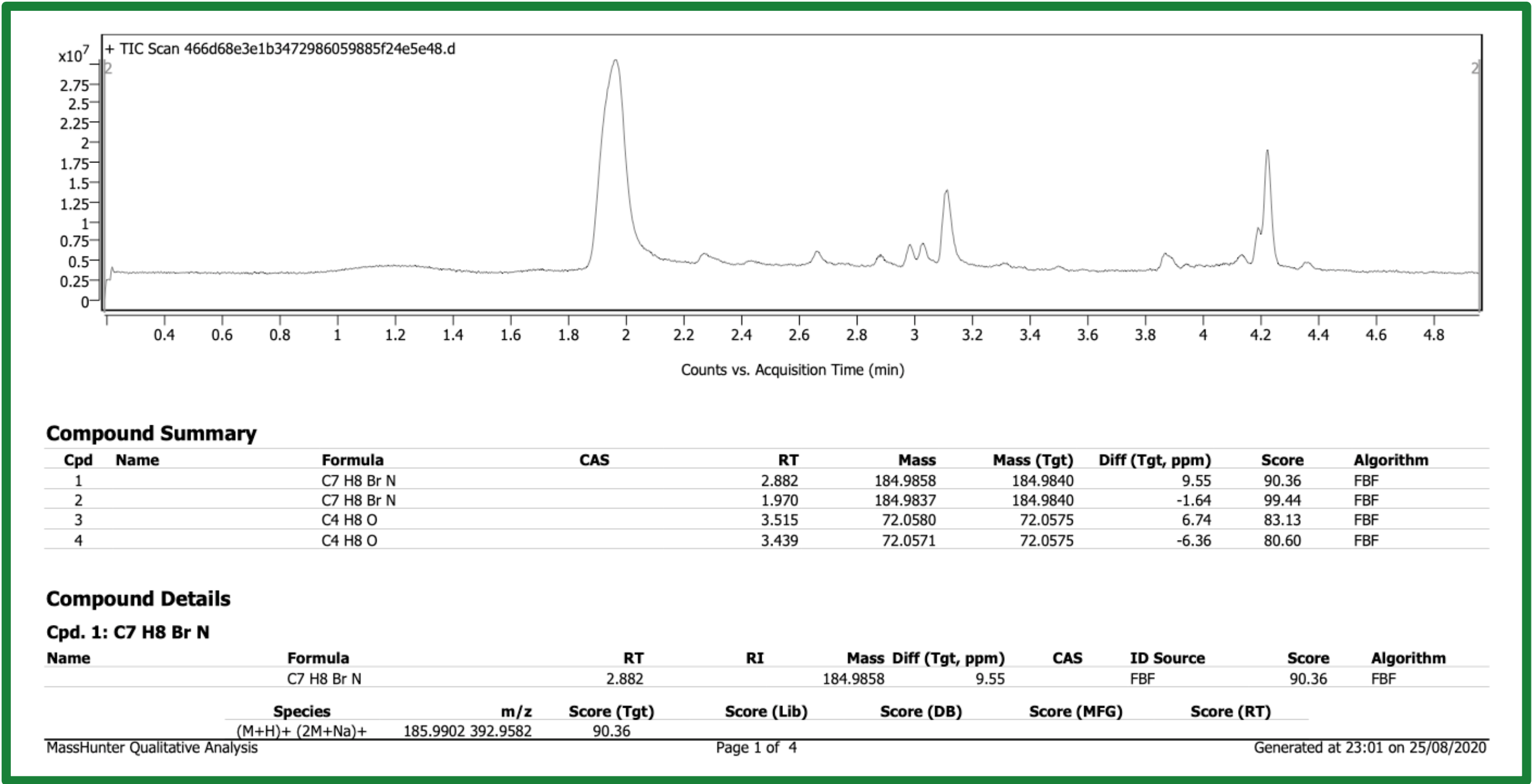
📖 Born, J. *et al.* Data-driven molecular design for discovery and synthesis of novel ligands: a case study on SARS-CoV-2. *Mach. Learn.: Sci. Technol.* **2**, 025024 (2021)

Molecular synthesis using IBM RoboRXN platform

Automated synthesis



LC/MS chromatogram



📖 Born, J. *et al.* Data-driven molecular design for discovery and synthesis of novel ligands: a case study on SARS-CoV-2. *Mach. Learn.: Sci. Technol.* **2**, 025024 (2021)

Google AI and robots join forces to build new materials

Tool from Google DeepMind predicts nearly 400,000 stable substances, and an autonomous system learns to make them in the lab.

Article

Scaling deep learning for materials discovery


<https://doi.org/10.1038/s41586-023-06735-9>

Received: 8 May 2023

Accepted: 10 October 2023

Published online: 29 November 2023

Open access

 Check for updates

Amil Merchant^{1,3}, Simon Batzner^{1,3}, Samuel S. Schoenholz^{1,3}, Muratahan Aykol¹, Gowoon Cheon² & Ekin Dogus Cubuk^{1,3}

Novel functional materials enable fundamental breakthroughs across technological applications from clean energy to information processing^{1–11}. From microchips to batteries and photovoltaics, discovery of inorganic crystals has been bottlenecked by expensive trial-and-error approaches. Concurrently, deep-learning models for

Article

An autonomous laboratory for the accelerated synthesis of novel materials


<https://doi.org/10.1038/s41586-023-06734-w>


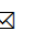
Received: 16 May 2023

Accepted: 10 October 2023


Published online: 29 November 2023

Open access

 Check for updates

Nathan J. Szymanski^{1,2,5}, Bernardus Rendy^{1,2,5}, Yuxing Fei^{1,2,5}, Rishi E. Kumar^{3,5}, Tanjin He^{1,2}, David Milsted², Matthew J. McDermott^{1,2}, Max Gallant^{1,2}, Ekin Dogus Cubuk⁴, Amil Merchant⁴, Haegyeom Kim², Anubhav Jain³, Christopher J. Bartel², Kristin Persson^{1,2}, Yan Zeng² & Gerbrand Ceder^{1,2}

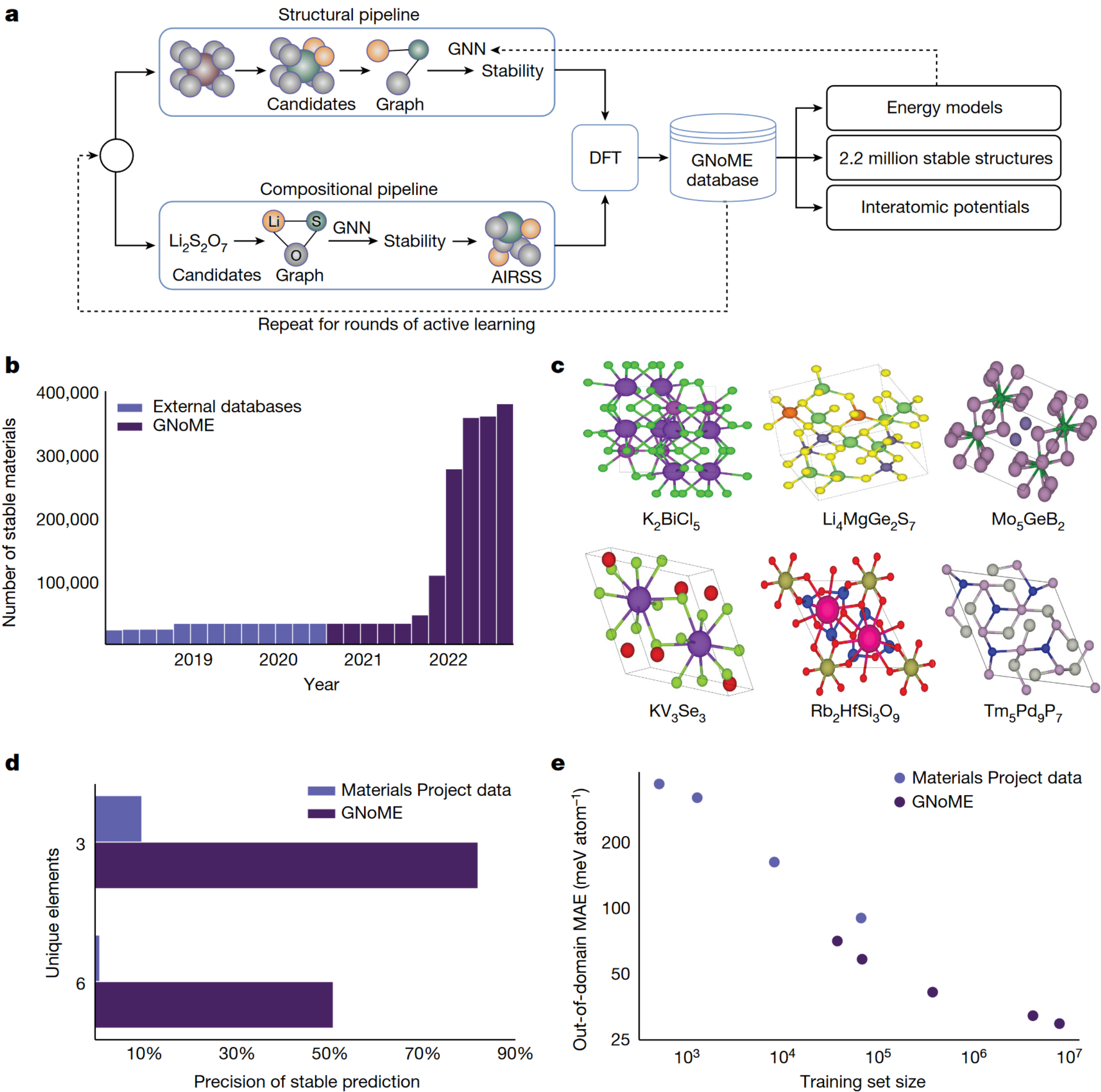
To close the gap between the rates of computational screening and experimental realization of novel materials^{1,2}, we introduce the A-Lab, an autonomous laboratory for the solid-state synthesis of inorganic powders. This platform uses computations,

 Merchant, A. *et al.* Scaling deep learning for materials discovery. *Nature* **624**, 80–85 (2023)

 Szymanski, N. J. *et al.* An autonomous laboratory for the accelerated synthesis of novel materials. *Nature* **624**, 86–91 (2023)

Generating a large of database stable crystal structures by deep learning

- Generate diverse candidate structures leveraging symmetry-aware partial substitutions and random structure search
- Graph neural networks (GNNs) used to model materials properties and filter candidate materials → graph networks for materials exploration (GNoME)
- Energy of filtered candidates is computed using density functional theory (DFT) and used to improve GNoME models through active learning

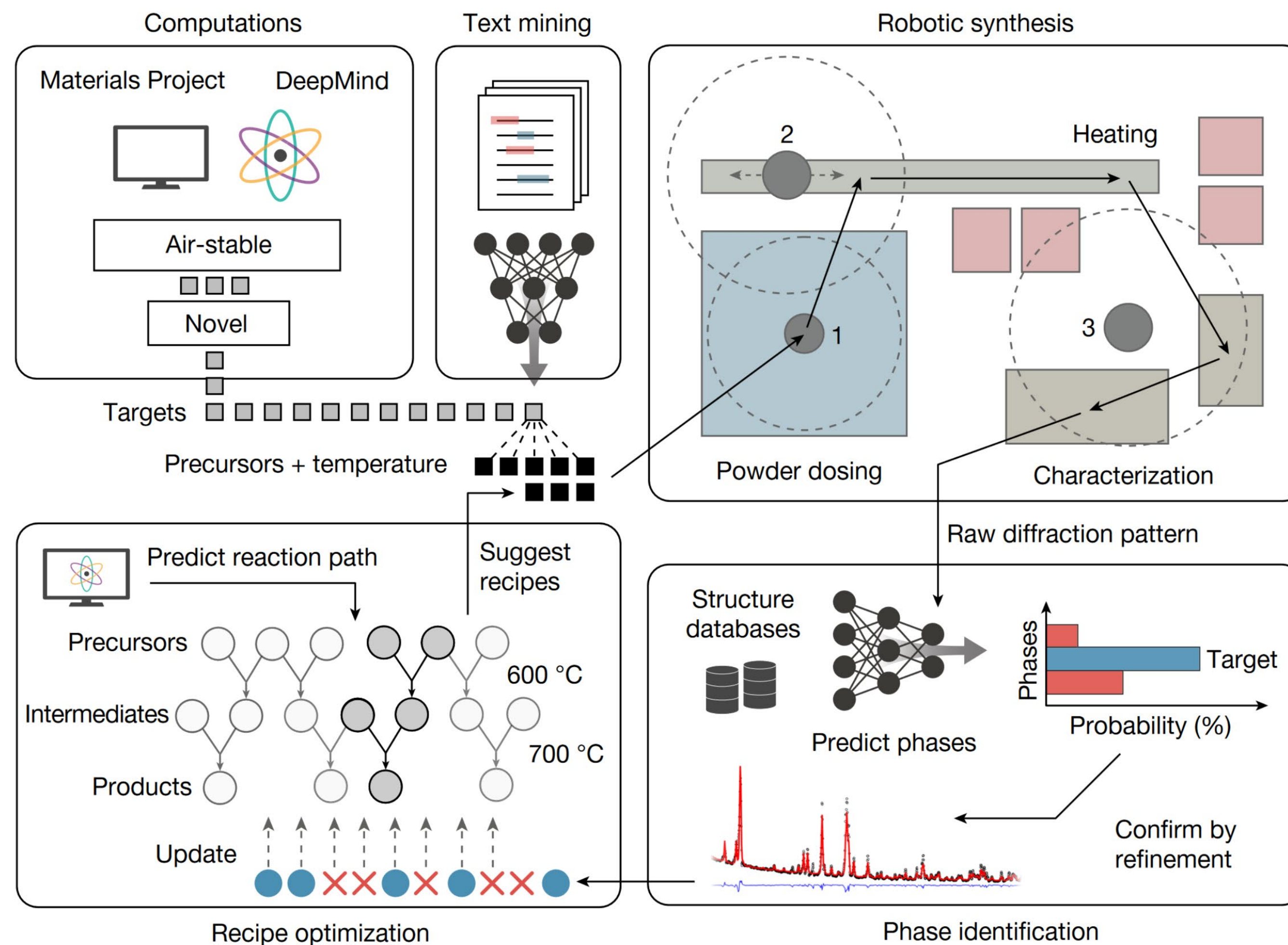


📖 Merchant, A. *et al.* Scaling deep learning for materials discovery. *Nature* **624**, 80–85 (2023)

📖 Szymanski, N. J. *et al.* An autonomous laboratory for the accelerated synthesis of novel materials. *Nature* **624**, 86–91 (2023)

Performing automated experimental validation of new crystal structures

- Synthesis recipes based on similar materials are generated by a language model trained on solid-state synthesis conditions from journal data
- Continuous experimentation with active learning to learn reaction pathways based on fixed set of precursors
- Phase and weight fractions extracted from X-ray diffraction (XRD) patterns using a convolutional neural network (CNN), automated Rietveld refinement, and manual analysis



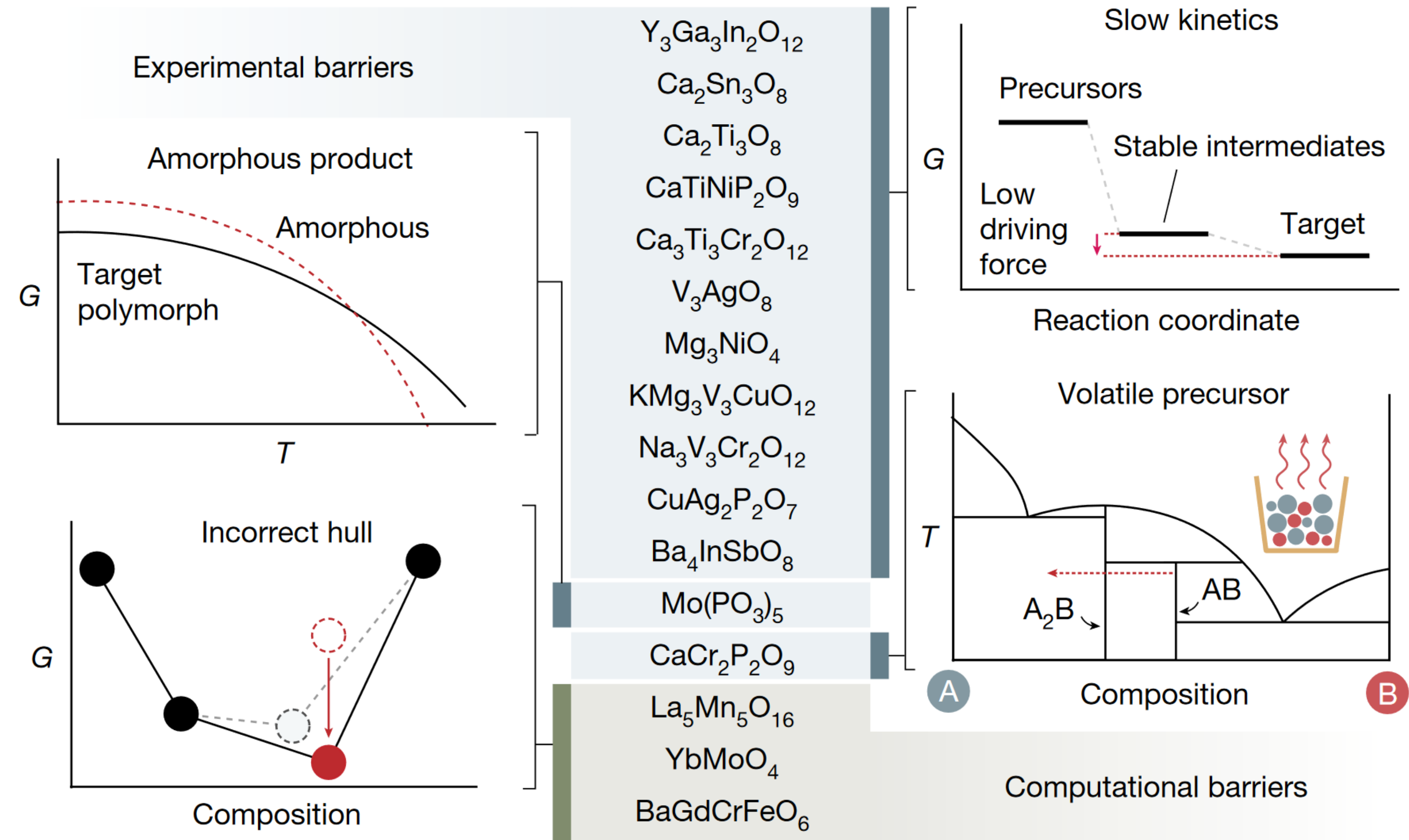
📖 Merchant, A. *et al.* Scaling deep learning for materials discovery. *Nature* **624**, 80–85 (2023)

📖 Szymanski, N. J. *et al.* An autonomous laboratory for the accelerated synthesis of novel materials. *Nature* **624**, 86–91 (2023)

Lessons learned from failures
to synthesize target product

17 out of 58 targets were not
obtained at the end of the
active learning cycle:

- 1. 11 failures due to sluggish
reaction kinetics
- 2. Decomposition and
evaporation of ammonium
phosphate precursors
- 3. Crystallization may be
inhibited upon melting of
samples at high
temperatures
- 4. Underestimation of energy
by DFT calculation



Merchant, A. *et al.* Scaling deep learning for materials discovery. *Nature* **624**, 80–85 (2023)

Szymanski, N. J. *et al.* An autonomous laboratory for the accelerated synthesis of novel materials. *Nature* **624**, 86–91 (2023)

Hands-on lab

- Design and synthesize your own molecule

Design and synthesize your
own molecule

1. Generate novel molecules and pick your favorite one!
https://huggingface.co/spaces/GT4SD/regression_transformer
2. Predict its molecular properties:
https://huggingface.co/spaces/GT4SD/regression_transformer
3. Find synthesis routes with IBM RXN:
<https://rxn.res.ibm.com>
4. Press a button to synthesize your own compound 😊



Summary

1. Scientific discovery as a loop
 - Evolution of the scientific method
 - AI for science
 - State of the art
2. Molecular generation
 - Language models for molecular discovery
 - Generative AI
 - Graph neural networks
3. Optimization strategies
 - Design of Experiment
 - Bayesian Optimization
 - Reinforcement Learning
4. Closing the loop: case studies
5. Hands-on lab