# Autonomous chemical research with large language models

Daniil A. Boiko, Robert MacKnight, Ben Kline & Gabe Gomes ✉

**Student : 林子軒**

**ID：113023503**

**Date : 2024/09/16**

**Lab :楊自雄老師 智動化分子研發實驗室**

# Introduction

# SMILES format

➤ **3D, 2D -> 1D, how?**

➤ **Simplicity**: SMILES uses **ASCII characters** to represent molecules, making it easy to input and process by computers.

➤ **Readability**: For simple molecules, **SMILES strings can be relatively easy** for humans to read and interpret.

➤ **Compactness**: SMILES can represent complex molecular structures **in a compact string format**.

A

B

C

D

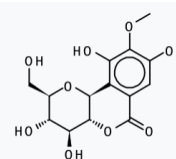N1CCN(CC1)C(C(F)=C2)=CC(=C2C4=O)N(C3CC3)C=C4C(=O)O

3

# SMILES format

- Atoms are represented by their **atomic symbols** (e.g., C for carbon, O for oxygen).

- **Single bonds** are implied and not explicitly written.

- **Double bonds** are represented by '=', triple bonds by '#'.

- **Branching** is shown using parentheses.

- **Rings** are indicated using numbers to show connection points.
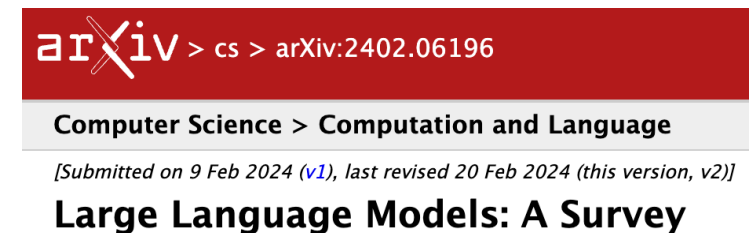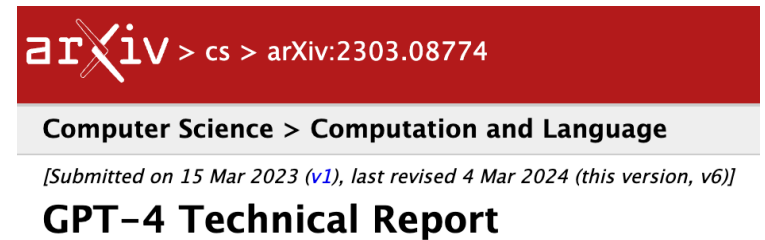
➤ Ethanol: CCO

➤ Benzene: c1ccccc1

Bergenin (cuscutin, a resin) ($C_{14}H_{16}O_9$)

`OC[C@@H](O1)[C@@H](O)[C@H](O)`
`[C@@H]2[C@@H]1c3c(O)c(OC)c(O)cc3C(=O)O2`

➤ CC1=CN(C2CC(O)C(CO[P](O)(=O)O[P](O)(=O)O[P](O)(O)=O)O2)C(=O)NC1=O

4

# Large language models (LLMs)



arXiv > cs > arXiv:2303.08774

**Computer Science > Computation and Language**

[Submitted on 15 Mar 2023 (v1), last revised 4 Mar 2024 (this version, v6)]

**GPT–4 Technical Report**

arXiv > math > arXiv:2404.03647

Search...
Help | Adv

**Mathematics > Optimization and Control**

[Submitted on 4 Apr 2024]

**Capabilities of Large Language Models in Control Engineering: A Benchmark Study on GPT–4, Claude 3 Opus, and Gemini 1.0 Ultra**

arXiv > cs > arXiv:2402.06196

**Computer Science > Computation and Language**

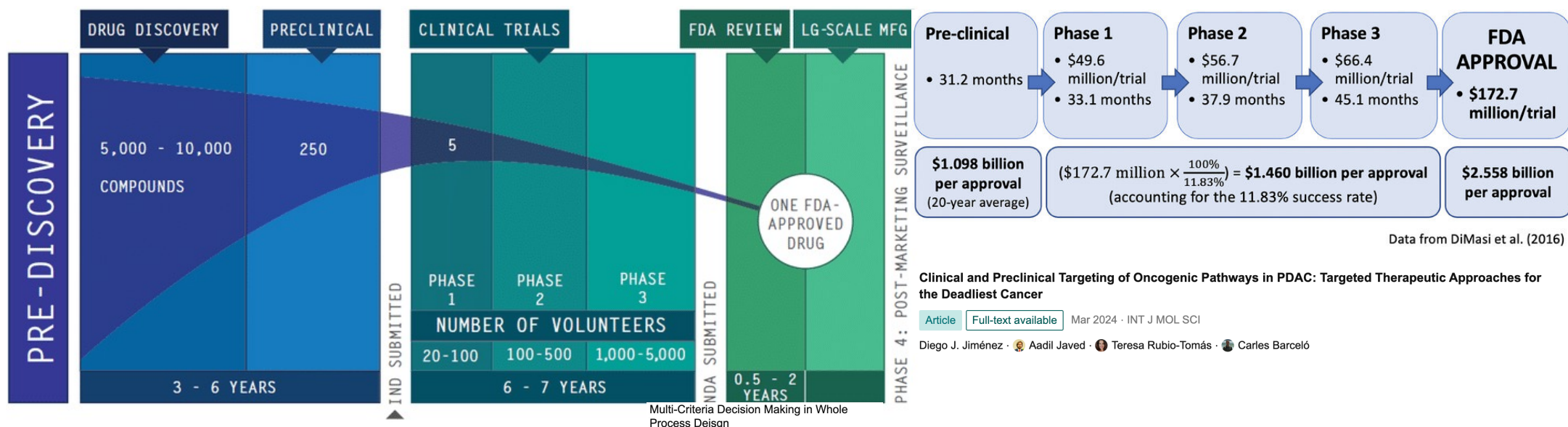[Submitted on 9 Feb 2024 (v1), last revised 20 Feb 2024 (this version, v2)]

**Large Language Models: A Survey**

➢Large Language Models (LLMs) are AI systems that process and generate **human language**. They are "large" due to their billions to trillions of parameters.

➢**Massive Training Data**: LLMs are trained on **vast amounts** of text data from diverse sources, including books, websites, and articles.

# Large language models (LLMs)

➤ But how is LLMs important in chemistry?

➤ What if it can predict a new drug or a new structure within only a few seconds and dollars?



Data from DiMasi et al. (2016)

Clinical and Preclinical Targeting of Oncogenic Pathways in PDAC: Targeted Therapeutic Approaches for the Deadliest Cancer

Article · Full-text available · Mar 2024 · INT J MOL SCI

Diego J. Jiménez · Aadil Javed · Teresa Rubio-Tomás · Carles Barceló

Multi-Criteria Decision Making in Whole Process Deisgn
February 2013
Thesis for: PhD · Advisor: Gary Montague, Elaine B Martin OBE FReng
Authors:
Richard Edgar Hodgett
University of Leeds

Article: Multi-Criteria Decision Making in Whole Proce
Author: Gary Montague, Elaine B Martin OBE FReng

6

# Coscientist

# Coscientist 6 tasks

✓planning chemical syntheses of known compounds using **publicly available data**

✓**efficiently searching** and navigating through extensive hardware documentation

✓using documentation to execute high-level commands in a **cloud laboratory**

✓precisely controlling **liquid handling instruments** with low-level instructions

✓tackling **complex scientific tasks** that demand simultaneous use of multiple hardware modules and integration of diverse data sources

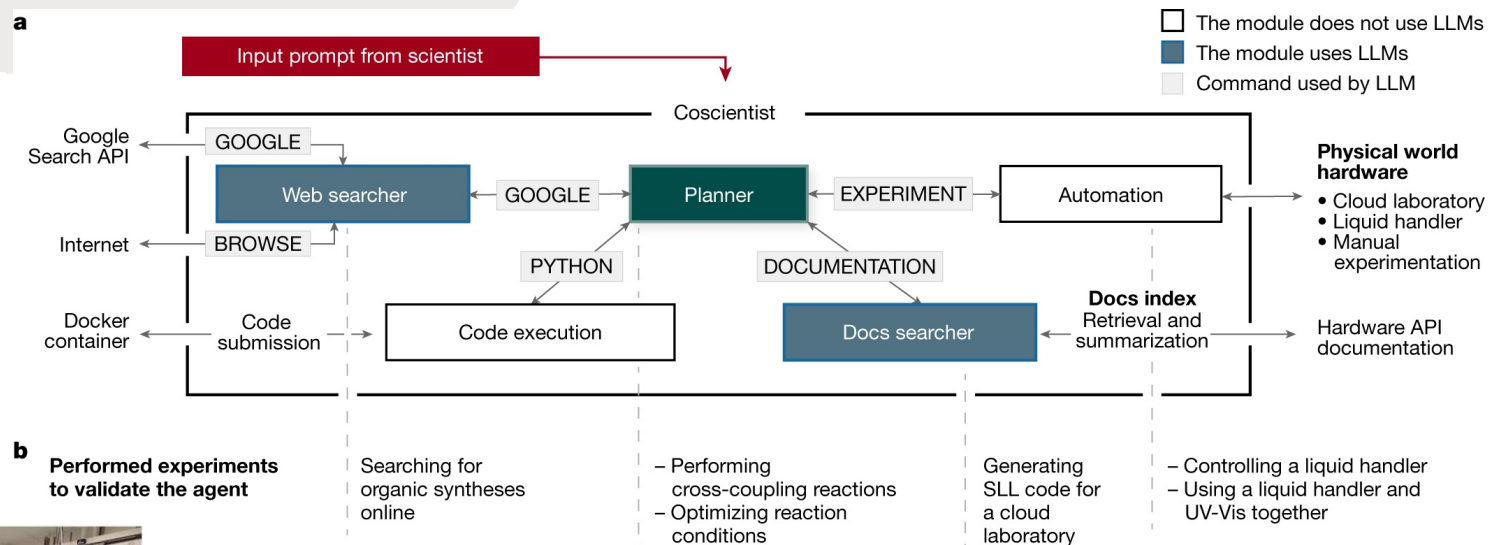✓**solving optimization problems** requiring analyses of previously collected **experimental data**

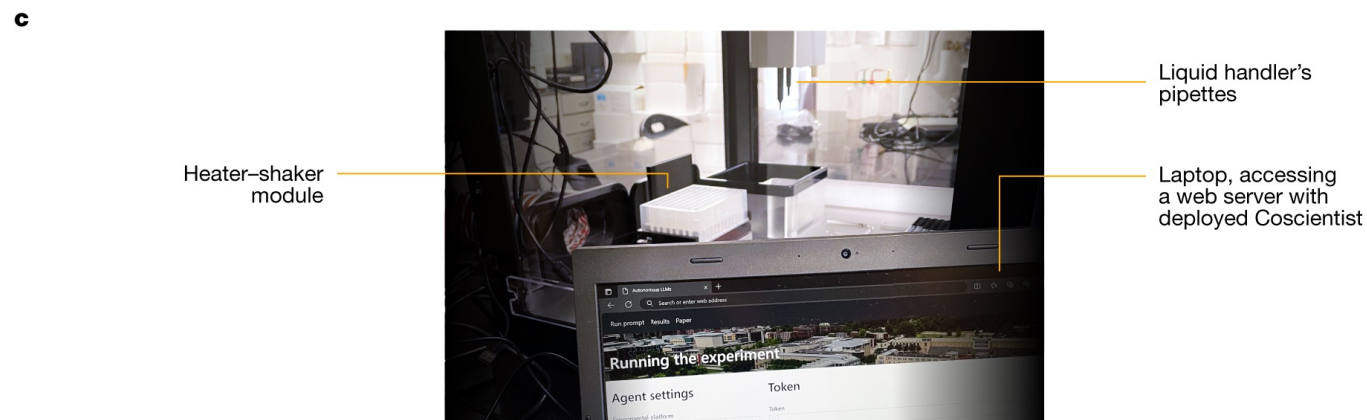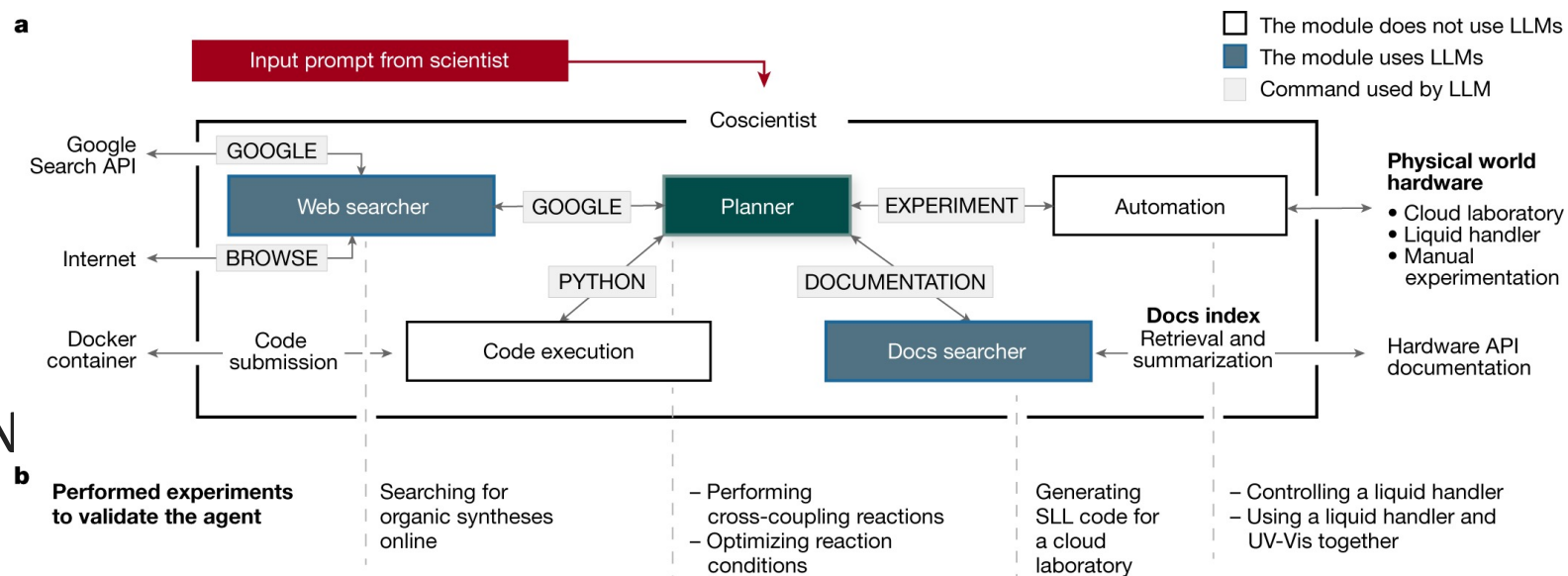Palladium-catalyzed cross-coupling reactions
Nobel prize 2010: Richard Fred Heck

8

# Coscientist system architecture

# Coscientist system architecture



**a**

Input prompt from scientist

Coscientist

| | The module does not use LLMs |
| | The module uses LLMs |
| | Command used by LLM |

Google Search API ← GOOGLE

Web searcher ← GOOGLE → Planner → EXPERIMENT → Automation

Internet ← BROWSE

**Physical world hardware**
- Cloud laboratory
- Liquid handler
- Manual experimentation

PYTHON

DOCUMENTATION

Docker container ← Code submission → Code execution

Docs searcher

**Docs index** Retrieval and summarization

Hardware API documentation

**b** **Performed experiments to validate the agent**

Searching for organic syntheses online

– Performing cross-coupling reactions
– Optimizing reaction conditions

Generating SLL code for a cloud laboratory

– Controlling a liquid handler
– Using a liquid handler and UV-Vis together

Processing robot platform

Synthesis robot platform

Characterization tools

Mobile robot for Sample transferring and loading

Liquid handler's pipettes

Heater–shaker module

Laptop, accessing a web server with deployed Coscientist

# Command

- ➢ GOOGLE
- ➢ PYTHON
- ➢ DOCUMENTATION
- ➢ EXPERIMENT



**a**

Input prompt from scientist

☐ The module does not use LLMs
▣ The module uses LLMs
☐ Command used by LLM

Coscientist

Google Search API — GOOGLE — Web searcher — GOOGLE — Planner — EXPERIMENT — Automation

Internet — BROWSE

**Physical world hardware**
- Cloud laboratory
- Liquid handler
- Manual experimentation

Docker container — Code submission — Code execution — PYTHON

DOCUMENTATION — Docs searcher

**Docs index** Retrieval and summarization

Hardware API documentation

**b**

| **Performed experiments to validate the agent** | Searching for organic syntheses online | – Performing cross-coupling reactions<br>– Optimizing reaction conditions | Generating SLL code for a cloud laboratory | – Controlling a liquid handler<br>– Using a liquid handler and UV-Vis together |

**c**

Liquid handler's pipettes

Heater–shaker module

Laptop, accessing a web server with deployed Coscientist

11

# Web search module

# Ranking

➢5 for a **very detailed** and **chemically accurate** procedure description

➢4 for a **detailed and chemically accurate** description but without reagent quantities

➢3 for a correct chemistry description that **does not include step-by-step procedure**

➢2 for **extremely vague or unfeasible** descriptions

➢1 for **incorrect responses or failure** to follow instructions
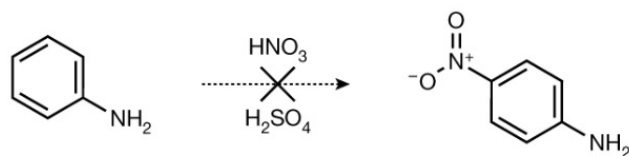
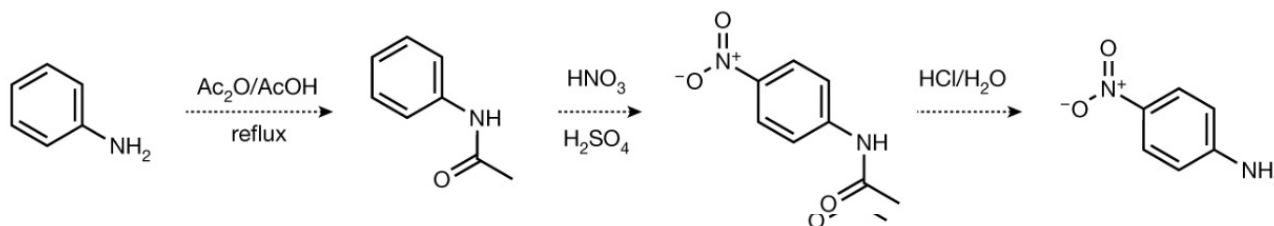All scores below 3 indicate task failure

# Web search module

# Evaluation

**b**

Incorrect synthesis steps but makes chemical sense ②
(GPT-3.5, no search)

Correct synthesis, including detailed experimental procedure ⑤
(GPT-4 with search)



**c**

Incorrect synthesis steps, does not make chemical sense (GPT-4, no search) ①



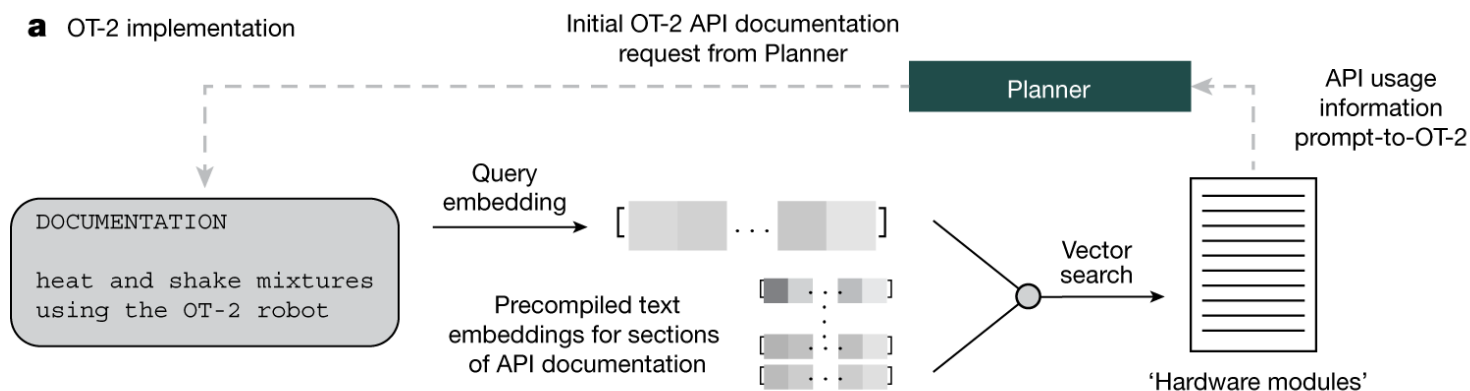Correct synthesis logic but no reagents and experimental procedure ③

# Documentation search module

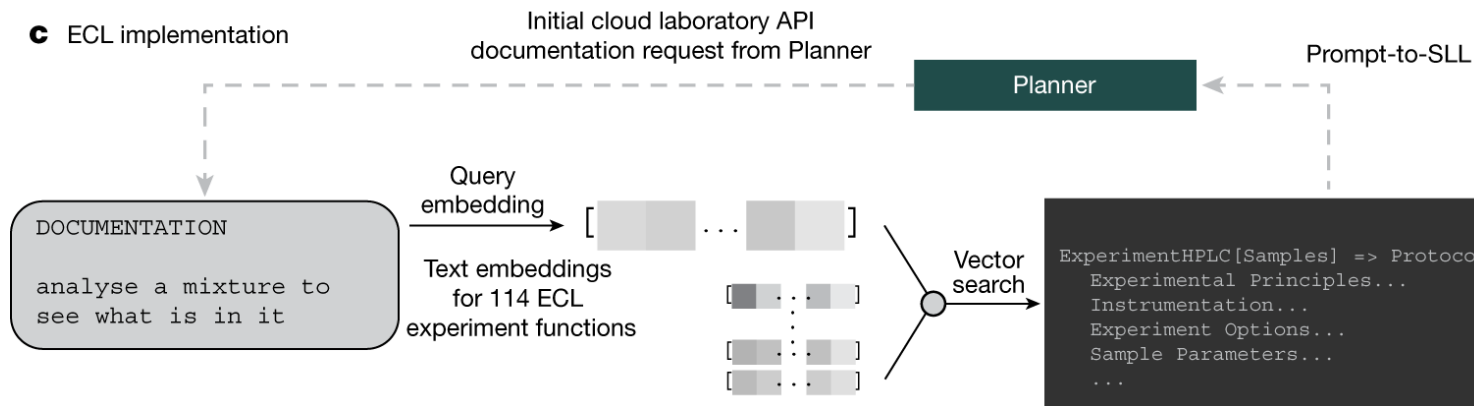# Documentation search module



**a** OT-2 implementation

Initial OT-2 API documentation request from Planner

Planner

API usage information prompt-to-OT-2

DOCUMENTATION

heat and shake mixtures using the OT-2 robot

Query embedding

Precompiled text embeddings for sections of API documentation

Vector search

'Hardware modules'

**b** Valid OT-2 API code

```
# Heat and shake the reaction
hs_mod.set_target_temperature(75)
hs_mod.wait_for_temperature()
hs_mod.set_and_wait_for_shake_speed(500)

# Deactivate heater and shaker
hs_mod.deactivate_heater()
hs_mod.deactivate_shaker()
hs_mod.open_labware_latch()
```

Proper usage of heater–shaker module

**c** ECL implementation

Initial cloud laboratory API documentation request from Planner

Planner

Prompt-to-SLL

DOCUMENTATION

analyse a mixture to see what is in it

Query embedding

Text embeddings for 114 ECL experiment functions

Vector search

```
ExperimentHPLC[Samples] => Protocol
    Experimental Principles...
    Instrumentation...
    Experiment Options...
    Sample Parameters...
    ...
```

**d** Valid ECL SLL code

```
# Generated HPLC Experiment SLL Function Call
ExperimentHPLC[
    Object[Sample, ...],
    Instrument -> Model[Instrument, ...]
]
```

Targeted experiment options are set by the Planner

Why is it important?

17
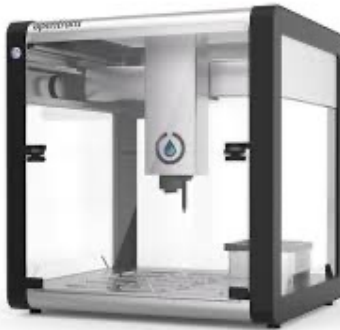
# Documentation search module



**a** OT-2 implementation

Initial OT-2 API documentation request from Planner

Planner

API usage information prompt-to-OT-2

DOCUMENTATION

heat and shake mixtures using the OT-2 robot

Query embedding

Precompiled text embeddings for sections of API documentation

Vector search

'Hardware modules'

**b** Valid OT-2 API code

```
# Heat and shake the reaction
hs_mod.set_target_temperature(75)
hs_mod.wait_for_temperature()
hs_mod.set_and_wait_for_shake_speed(500)

# Deactivate heater and shaker
hs_mod.deactivate_heater()
hs_mod.deactivate_shaker()
hs_mod.open_labware_latch()
```

Proper usage of heater–shaker module

**c** ECL implementation

Initial cloud laboratory API documentation request from Planner

Planner

Prompt-to-SLL

DOCUMENTATION

analyse a mixture to see what is in it

Query embedding

Text embeddings for 114 ECL experiment functions

Vector search

```
ExperimentHPLC[Samples] => Protocol
    Experimental Principles...
    Instrumentation...
    Experiment Options...
    Sample Parameters...
    ...
```

**d** Valid ECL SLL code

```
# Generated HPLC Experiment SLL Function Call
ExperimentHPLC[
    Object[Sample, ...],
    Instrument -> Model[Instrument, ...]
]
```

Targeted experiment options are set by the Planner

Why is it important?

✓**Technical Integration**
   ✓Capable of addressing the complexity of **software components** and their interactions
   ✓Crucial for integrating LLMs with **laboratory automation**

✓**Effective Utilization of Technical Documentation**
   ✓Enables Coscientist to **understand** and use technical documentation

✓**Learning New Languages and Systems**
   ✓Demonstrates GPT-4's ability to **learn new programming languages** (like ECL SLL)

# Controlling laboratory hardware

# Controlling laboratory hardware

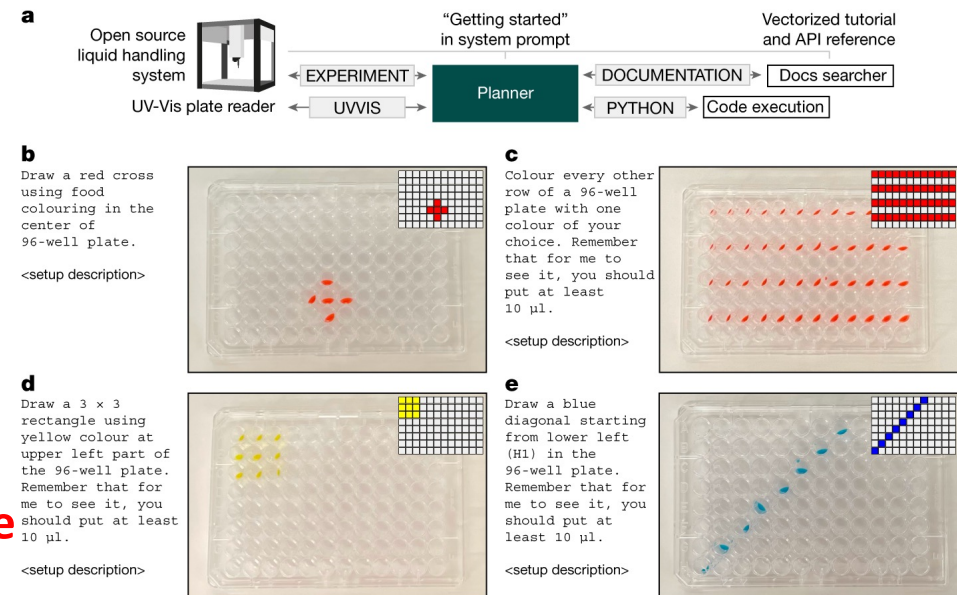Robotic liquid handler control capabilities and integration with analytical tools.
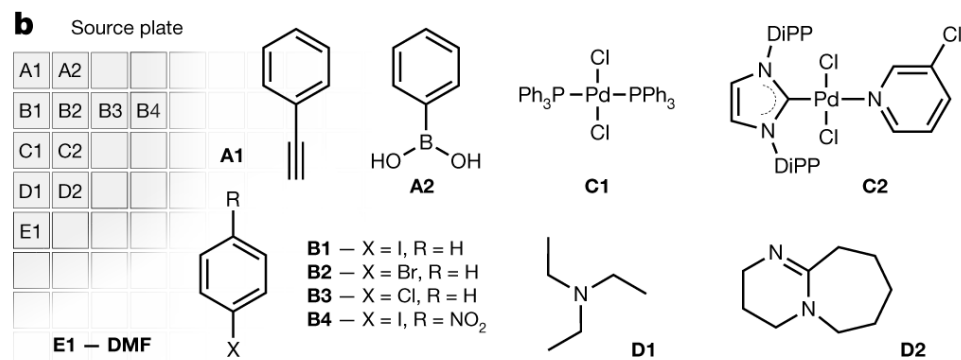
# Controlling laboratory hardware



OT-2

➤ Coscientist conducts experiments using an open-source **liquid handler OT-2** and its **API**, obtaining information from **documentation**.

➤ Coscientist is capable of integrating multiple modules, like using **UV-Vis spectrometry**, and solving complex tasks through data analysis, such as **determining sample colors and positions without prior information**.

# Integrated chemical experiment design

# Controlling laboratory hardware

Cross-coupling Suzuki and Sonogashira reaction experiments designed and performed by Coscientist



➤ No human decision

# Cross-coupling Suzukiand Sonogashira reaction

# Cross-coupling Suzuki and Sonogashira reaction



> The model is able to **provide explanations** for specific choices
> demonstrating its **capability to handle concepts such as reactivity and selectivity**.

# Discussion

# Discussion

➢ This study demonstrates an AI system capable of **autonomously designing** and **executing complex scientific experiments**.

➢ The system combines **large language models** with research tools, showcasing **advanced reasoning and experimental design capabilities**.

# Thanks for listening!