

Agents for Scientific Discovery

Boris Bolliet



UNIVERSITY OF
CAMBRIDGE



ACCELERATE
PROGRAMME
FOR SCIENTIFIC DISCOVERY



Join our growing community of over 20,000 AI agent builders



THE OPEN-SOURCE AGENTOS



Build production-ready AI agents in minutes, not months.

Enable AI-Native Organizations.



Chi Wang
Google DeepMind

[Request Demo](#)

[GitHub](#)



Accelerated inference with neural networks and agents for Cosmic Microwave Background and Large Scale Structure analyses

Boris Bolliet

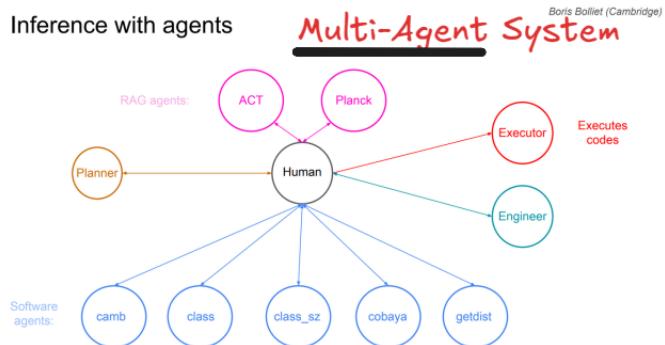
Cavendish Astrophysics and Kavli Institute for Cosmology, Cambridge

Work in collaboration with:

(Neural nets:) Kristen Surrao, Frank Qu, Hidde Jense, Colin Hill, Julien Lesgourgues, Alessio Spurio Mancini, Blake Sherwin
(Agents:) Andrew Laverick, Inigo Zubeldia, Miles Cranmer, Julien Lesgourgues, Antony Lewis, Blake Sherwin



Inference with agents



Expediting Astronomical Discovery with Large Language Models: Progress, Challenges, and Future Directions

Invited Speaker: Yuan Sen-Ting (Australian National University and Ohio State University)

The vast and interdisciplinary nature of astronomy, coupled with its open-access ethos, makes it an ideal testbed for exploring the potential of Large Language Models (LLMs) in automating and accelerating scientific discovery. In this talk, we present our recent progress in applying LLMs to tackle real-life astronomy problems. We demonstrate the ability of LLM agents to perform end-to-end research tasks, from data fitting and analysis to iterative strategy improvement and outlier detection, mimicking human intuition and deep literature understanding. However, the cost-effectiveness of closed-source solutions remains a challenge for large-scale applications involving billions of sources. To address this issue, we introduce our ongoing work at AstroMLab on training lightweight, open-source specialized models and our effort to benchmark these models with carefully curated astronomy benchmark datasets. We will also discuss our effort to construct the first LLM-based knowledge graph in astronomy, leveraging citation-reference relations. The open-source specialized LLMs and knowledge graph are expected to guide more efficient strategy searches in autonomous research pipelines. While many challenges lie ahead, we explore the immense potential of scaling up automated inference in astronomy, revolutionizing the way astronomical research is conducted, ultimately accelerating scientific breakthroughs and deepening our understanding of the Universe.



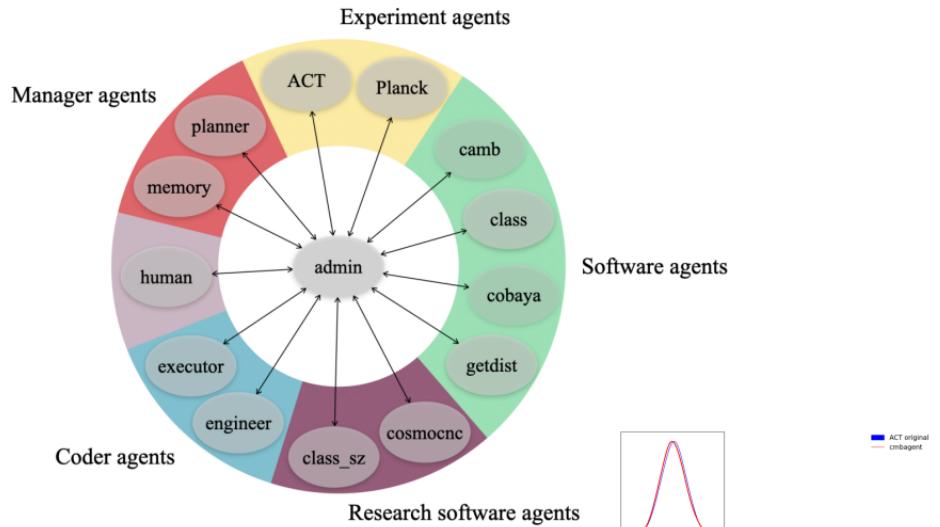
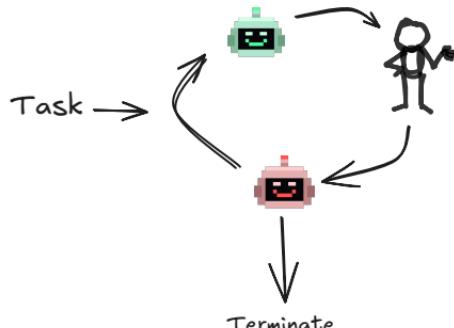
Multi-Agent System for Cosmological Parameter Analysis

Andrew Laverick¹ Kristen Surrao² Iñigo Zubeldia³ Boris Bolliet³
 Miles Cranmer³ Antony Lewis⁴ Blake Sherwin³ Julien Lesgourgues⁵

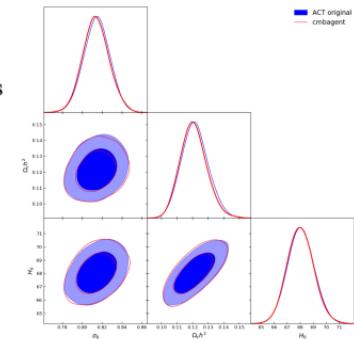
¹University of Manchester ²Columbia University
³University of Cambridge ⁴University of Sussex ⁵RWTH Aachen University

Abstract

Multi-agent systems (MAS) utilizing multiple Large Language Model (LLM) agents with Retrieval Augmented Generation and that can execute code locally may become beneficial in cosmological data analysis. Here, we illustrate a first small step towards AI-assisted analyses and a glimpse of the potential of MAS to automate and optimize scientific workflows in Cosmology. The system architecture of our example package, that builds upon the autogen/ag2¹ framework, can be applied to MAS in any area of quantitative scientific research. The particular task we apply our methods to is the cosmological parameter analysis of the Atacama Cosmology Telescope lensing power spectrum likelihood using Monte Carlo Markov Chains. Our work-in-progress code is open source and available at <https://github.com/CMBAgents/cmbagent>.



Human in the loop?... Bad.



No human-in-the-loop!

Motivated by discussions w. David Kaiser (MIT) and Bruce Bassett (SAO)'s 2024 article:
"Integrals and Integrity: Generative AI Tries to Learn Cosmology" MIT SERC journal



Task:

Download the file: https://supernova.lbl.gov/Union/figures/SCPUunion2.1_mu_vs_z.txt

Its description is:

<description>

An ASCII table with tab-separated columns: Supernova Name, Redshift, Distance Modulus, and Distance Modulus Error. For Union2.1, there is an additional column for the probability that the supernova was hosted by a low-mass galaxy.

</description>

Fit this data within flat LCDM model with two free parameters: H_0 and Ω_m . Write a simple MCMC (but optimized/fast) code to fit for H_0 and Ω_m using the SN1a data.

Make a contour plot and show the 1d posteriors, and quote the mean and 1-sigma on each parameter.

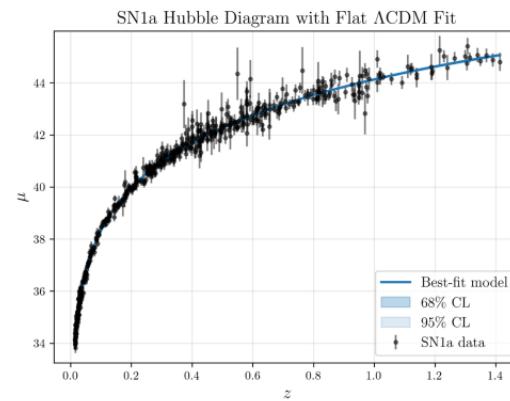
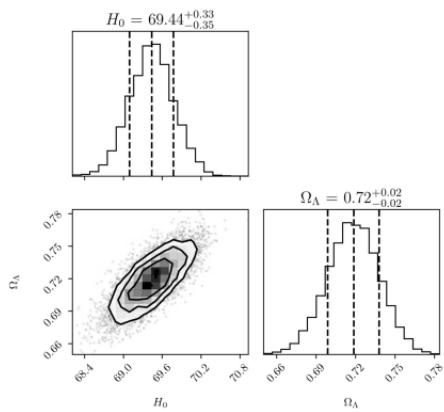
Finally show the data along with the best fit model and 68%/95% CL regions.

Comment on the results.

Constraints:

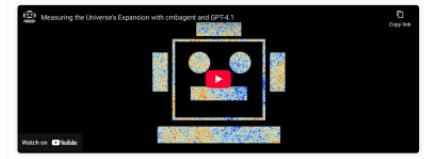
We are running this analysis on a Macbook Pro with 8 available threads. Ensure you use the resources optimally so the MCMC can run fast, i.e., within a few minutes until convergence.

Have the engineer agent do a preliminary MCMC timing step in a separate step.



Measuring the Universe's Expansion with cmbagent and GPT-4.1

Can AI, without any human-in-the-loop, reproduce the statistical data analysis that originally revealed the accelerating expansion of the universe? Yes! Note: this is the result that was awarded the 2021 Nobel Prize in Physics!



Scan me

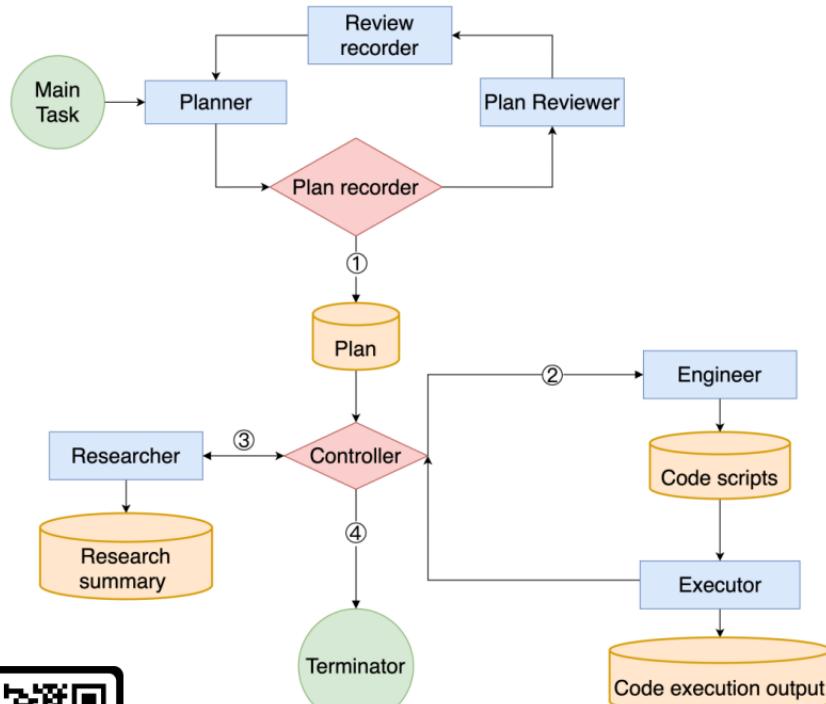
Open Source Planning & Control System with Language Agents for Autonomous Scientific Discovery

Licong Xu ^{1,2} Milind Sarkar ³ Anto I. Lonappan ⁴ Íñigo Zubeldia ^{1,2} Pablo Villanueva-Domingo ⁵
 Santiago Casas ⁶ Christian Fidler ⁶ Chetana Amancharya ⁷ Ujjwal Tiwari ⁷ Adrian Bayer ^{8,9}
 Chadi Ait Ekioui ^{10,11} Miles Cranmer ^{1,2,12} Adrian Dimitrov ¹⁰ James Ferguson ¹² Kahaan Gandhi ^{10,13,14}
 Sven Krippendorff ^{12,10} Andrew Laverick ¹⁰ Julien Lesgourgues ⁶ Antony Lewis ¹⁵ Thomas Meier ¹⁶
 Blake Sherwin ^{2,12} Kristen Surrao ¹⁷ Francisco Villascsa-Navarro ^{8,9} Chi Wang ¹⁸ Xueqing Xu ¹⁰
 Boris Bolliet ^{2,10}

Abstract

We present a multi-agent system for automation of scientific research tasks, `cmbagent`. The system is formed by about 30 Large Language Model (LLM) agents and implements a *Planning & Control* strategy to orchestrate the agentic workflow, with no *human-in-the-loop* at any point. Each agent specializes in a different task (performing retrieval on scientific papers and codebases, writing code, interpreting results, critiquing the output of other agents) and the system is able to execute

ArXiv:2507.07257



SCAN ME

Goal: Superhuman Research

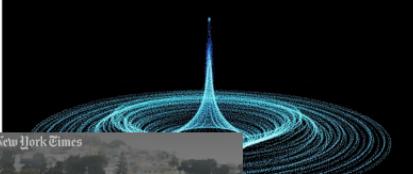
Get the open-source code here

September 18, 2025 Science

Discovering new solutions to century-old problems in fluid dynamics

Yongji Wang, Sam Blackwell

Share 



The New York Times

Top A.I. Researchers Leave OpenAI, Google and Meta for New Start-up

Founded by a co-creator of ChatGPT, Petzold Lab aims to build artificial intelligence that can accelerate discoveries in physics, chemistry and other fields.

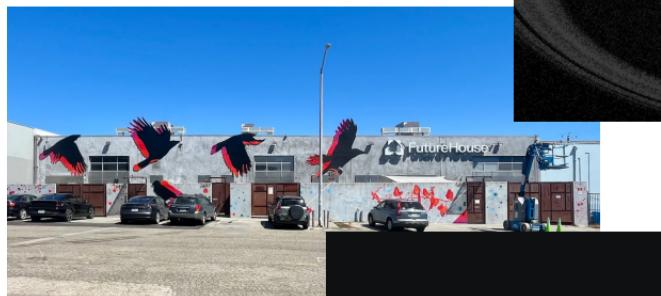
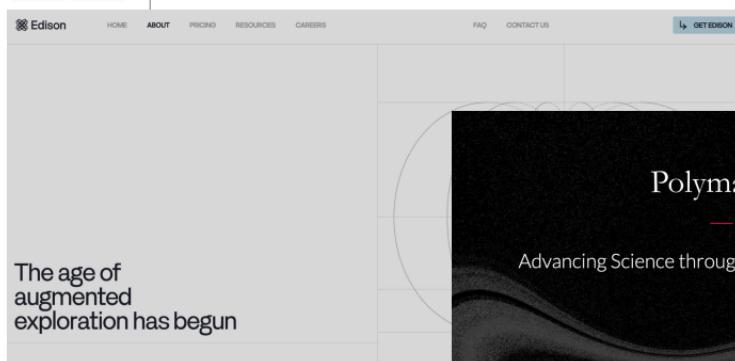


Credit, image and quote: New York Times, Sept 30, 2025

Edison

HOME ABOUT PRIMO RESOURCES CAREERS FAQ CONTACT US GET EDISON

The age of augmented exploration has begun



Polymathic

Advancing Science through Multi-Disciplinary AI

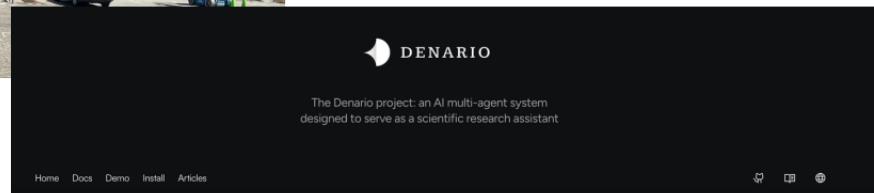
L I L A

Building Scientific Superintelligence

DENARIO

The Denario project: an AI multi-agent system designed to serve as a scientific research assistant

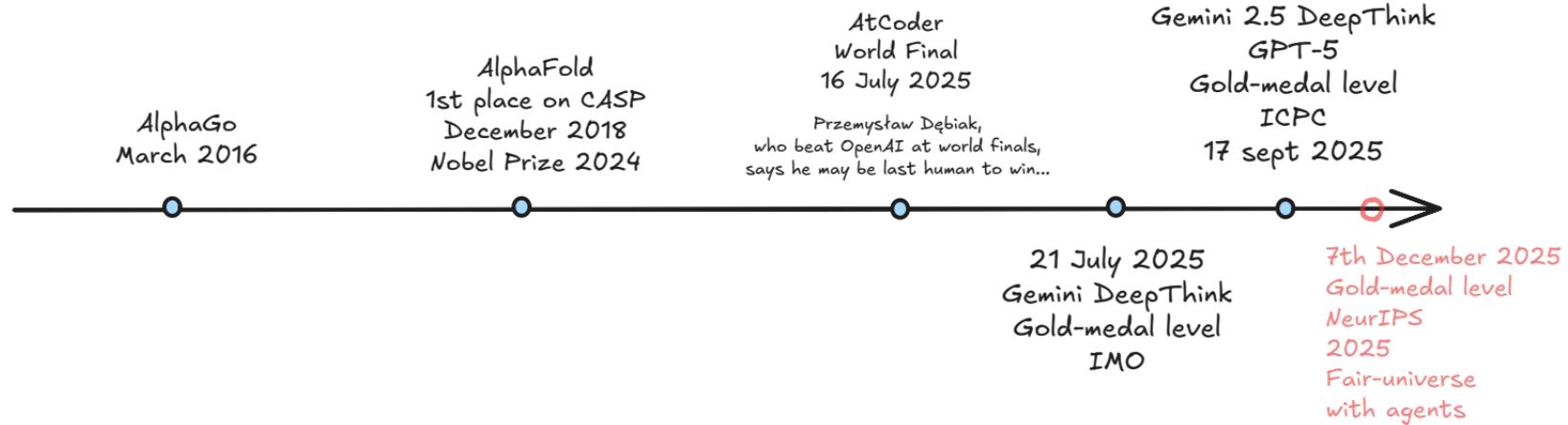
Home Docs Demo Install Articles



"The main objective of A.I. is not to automate white-collar work," said Liam Fedus, one of the start-up's founders. "The main objective is to accelerate science."

Agents for Scientific Discovery

Can research be automated?



In data-driven fields, AI outperforms humans at most tasks

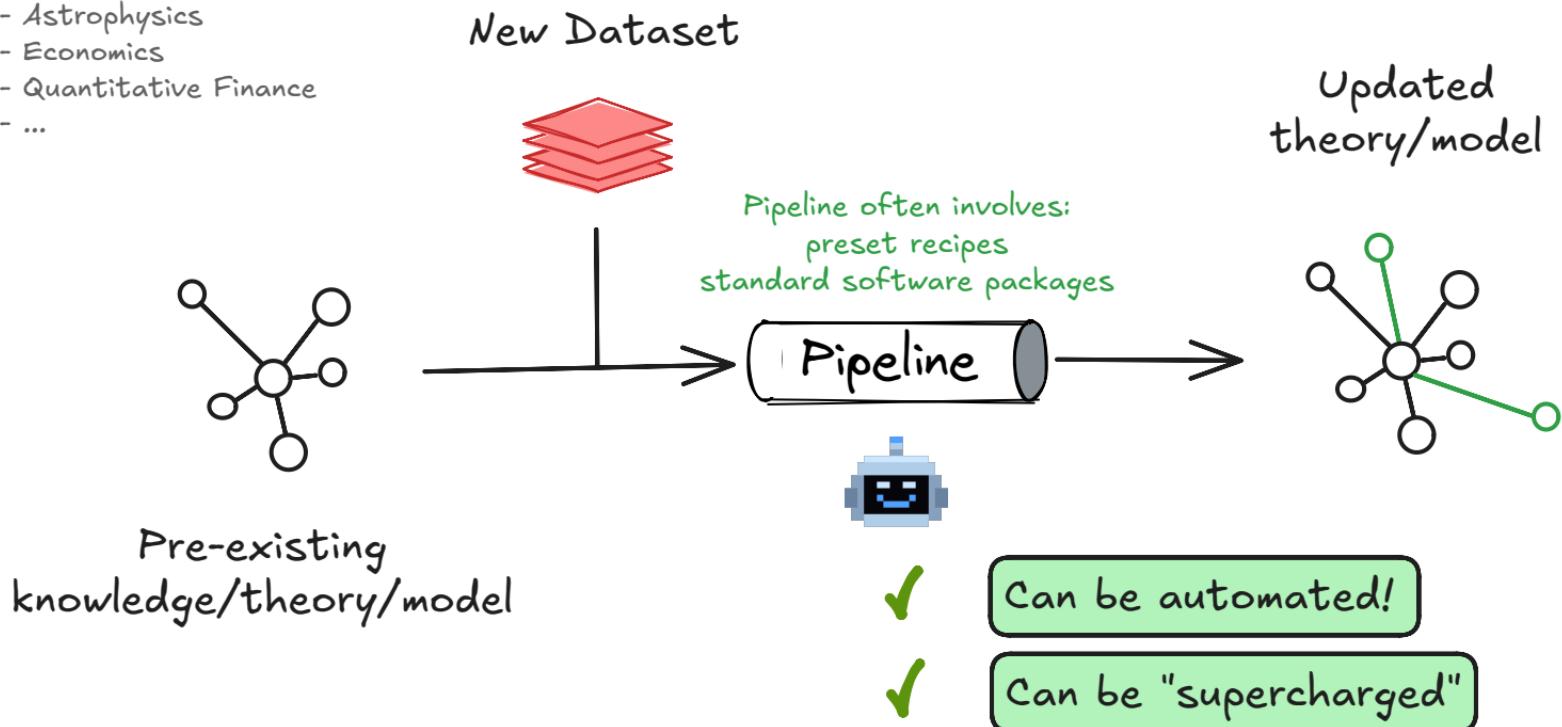
Massive opportunity for acceleration and discovery

"Superhuman" research?

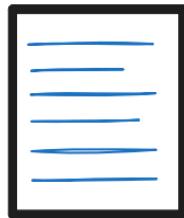
Rank	Participant	Date	ID	Method Name	Results	
					Part	Short Answer
1	○	2025-07-15	404234	tomorrow_collapse	11.7021	0
2	○	2025-07-15	404235	tomorrow_collapse	11.7021	0
3	○	2025-07-14	402350	andreas_collapse	11.7021	0
4	○	2025-07-15	404236	tomorrow_collapse	11.7026	0
5	○	2025-07-14	402351	tomorrow_collapse	11.7026	0
6	○	2025-07-15	404237	tomorrow_collapse	11.7026	0
7	○	2025-07-14	402352	tomorrow_collapse	11.7026	0
8	○	2025-07-15	404238	tomorrow_collapse	11.7026	0
9	○	2025-07-14	402353	tomorrow_collapse	11.7026	0
10	○	2025-07-15	404239	tomorrow_collapse	11.7026	0
11	○	2025-07-12	401916	tomorrow_collapse	11.6020	0
12	○	2025-07-15	404240	tomorrow_collapse	11.6020	0
13	○	2025-07-15	404241	tomorrow_collapse	11.6020	0
14	○	2025-07-15	404242	tomorrow_collapse	11.6020	0
15	○	2025-07-15	404243	tomorrow_collapse	11.6020	0
16	○	2025-07-15	404244	tomorrow_collapse	11.6020	0
17	○	2025-07-15	404245	tomorrow_collapse	11.6020	0
18	○	2025-07-15	404246	tomorrow_collapse	11.6020	0
19	○	2025-07-15	404247	tomorrow_collapse	11.6020	0
20	○	2025-07-15	404248	tomorrow_collapse	11.6020	0
21	○	2025-07-15	404249	tomorrow_collapse	11.6020	0
22	○	2025-07-15	404250	tomorrow_collapse	11.6020	0
23	○	2025-07-15	404251	tomorrow_collapse	11.6020	0
24	○	2025-07-15	404252	tomorrow_collapse	11.6020	0
25	○	2025-07-15	404253	tomorrow_collapse	11.6020	0
26	○	2025-07-15	404254	tomorrow_collapse	11.6020	0
27	○	2025-07-15	404255	tomorrow_collapse	11.6020	0
28	○	2025-07-15	404256	tomorrow_collapse	11.6020	0
29	○	2025-07-15	404257	tomorrow_collapse	11.6020	0
30	○	2025-07-15	404258	tomorrow_collapse	11.6020	0
31	○	2025-07-15	404259	tomorrow_collapse	11.6020	0
32	○	2025-07-15	404260	tomorrow_collapse	11.6020	0
33	○	2025-07-15	404261	tomorrow_collapse	11.6020	0
34	○	2025-07-15	404262	tomorrow_collapse	11.6020	0
35	○	2025-07-15	404263	tomorrow_collapse	11.6020	0
36	○	2025-07-15	404264	tomorrow_collapse	11.6020	0
37	○	2025-07-15	404265	tomorrow_collapse	11.6020	0
38	○	2025-07-15	404266	tomorrow_collapse	11.6020	0
39	○	2025-07-15	404267	tomorrow_collapse	11.6020	0
40	○	2025-07-15	404268	tomorrow_collapse	11.6020	0
41	○	2025-07-15	404269	tomorrow_collapse	11.6020	0
42	○	2025-07-15	404270	tomorrow_collapse	11.6020	0
43	○	2025-07-15	404271	tomorrow_collapse	11.6020	0
44	○	2025-07-15	404272	tomorrow_collapse	11.6020	0
45	○	2025-07-15	404273	tomorrow_collapse	11.6020	0
46	○	2025-07-15	404274	tomorrow_collapse	11.6020	0
47	○	2025-07-15	404275	tomorrow_collapse	11.6020	0
48	○	2025-07-15	404276	tomorrow_collapse	11.6020	0
49	○	2025-07-15	404277	tomorrow_collapse	11.6020	0
50	○	2025-07-15	404278	tomorrow_collapse	11.6020	0
51	○	2025-07-15	404279	tomorrow_collapse	11.6020	0
52	○	2025-07-15	404280	tomorrow_collapse	11.6020	0
53	○	2025-07-15	404281	tomorrow_collapse	11.6020	0
54	○	2025-07-15	404282	tomorrow_collapse	11.6020	0
55	○	2025-07-15	404283	tomorrow_collapse	11.6020	0
56	○	2025-07-15	404284	tomorrow_collapse	11.6020	0
57	○	2025-07-15	404285	tomorrow_collapse	11.6020	0
58	○	2025-07-15	404286	tomorrow_collapse	11.6020	0
59	○	2025-07-15	404287	tomorrow_collapse	11.6020	0
60	○	2025-07-15	404288	tomorrow_collapse	11.6020	0
61	○	2025-07-15	404289	tomorrow_collapse	11.6020	0
62	○	2025-07-15	404290	tomorrow_collapse	11.6020	0
63	○	2025-07-15	404291	tomorrow_collapse	11.6020	0
64	○	2025-07-15	404292	tomorrow_collapse	11.6020	0
65	○	2025-07-15	404293	tomorrow_collapse	11.6020	0
66	○	2025-07-15	404294	tomorrow_collapse	11.6020	0
67	○	2025-07-15	404295	tomorrow_collapse	11.6020	0
68	○	2025-07-15	404296	tomorrow_collapse	11.6020	0
69	○	2025-07-15	404297	tomorrow_collapse	11.6020	0
70	○	2025-07-15	404298	tomorrow_collapse	11.6020	0
71	○	2025-07-15	404299	tomorrow_collapse	11.6020	0
72	○	2025-07-15	404300	tomorrow_collapse	11.6020	0
73	○	2025-07-15	404301	tomorrow_collapse	11.6020	0
74	○	2025-07-15	404302	tomorrow_collapse	11.6020	0
75	○	2025-07-15	404303	tomorrow_collapse	11.6020	0
76	○	2025-07-15	404304	tomorrow_collapse	11.6020	0
77	○	2025-07-15	404305	tomorrow_collapse	11.6020	0
78	○	2025-07-15	404306	tomorrow_collapse	11.6020	0
79	○	2025-07-15	404307	tomorrow_collapse	11.6020	0
80	○	2025-07-15	404308	tomorrow_collapse	11.6020	0
81	○	2025-07-15	404309	tomorrow_collapse	11.6020	0
82	○	2025-07-15	404310	tomorrow_collapse	11.6020	0
83	○	2025-07-15	404311	tomorrow_collapse	11.6020	0
84	○	2025-07-15	404312	tomorrow_collapse	11.6020	0
85	○	2025-07-15	404313	tomorrow_collapse	11.6020	0
86	○	2025-07-15	404314	tomorrow_collapse	11.6020	0
87	○	2025-07-15	404315	tomorrow_collapse	11.6020	0
88	○	2025-07-15	404316	tomorrow_collapse	11.6020	0
89	○	2025-07-15	404317	tomorrow_collapse	11.6020	0
90	○	2025-07-15	404318	tomorrow_collapse	11.6020	0
91	○	2025-07-15	404319	tomorrow_collapse	11.6020	0
92	○	2025-07-15	404320	tomorrow_collapse	11.6020	0
93	○	2025-07-15	404321	tomorrow_collapse	11.6020	0
94	○	2025-07-15	404322	tomorrow_collapse	11.6020	0
95	○	2025-07-15	404323	tomorrow_collapse	11.6020	0
96	○	2025-07-15	404324	tomorrow_collapse	11.6020	0
97	○	2025-07-15	404325	tomorrow_collapse	11.6020	0
98	○	2025-07-15	404326	tomorrow_collapse	11.6020	0
99	○	2025-07-15	404327	tomorrow_collapse	11.6020	0
100	○	2025-07-15	404328	tomorrow_collapse	11.6020	0
101	○	2025-07-15	404329	tomorrow_collapse	11.6020	0
102	○	2025-07-15	404330	tomorrow_collapse	11.6020	0
103	○	2025-07-15	404331	tomorrow_collapse	11.6020	0
104	○	2025-07-15	404332	tomorrow_collapse	11.6020	0
105	○	2025-07-15	404333	tomorrow_collapse	11.6020	0
106	○	2025-07-15	404334	tomorrow_collapse	11.6020	0
107	○	2025-07-15	404335	tomorrow_collapse	11.6020	0
108	○	2025-07-15	404336	tomorrow_collapse	11.6020	0
109	○	2025-07-15	404337	tomorrow_collapse	11.6020	0
110	○	2025-07-15	404338	tomorrow_collapse	11.6020	0
111	○	2025-07-15	404339	tomorrow_collapse	11.6020	0
112	○	2025-07-15	404340	tomorrow_collapse	11.6020	0
113	○	2025-07-15	404341	tomorrow_collapse	11.6020	0
114	○	2025-07-15	404342	tomorrow_collapse	11.6020	0
115	○	2025-07-15	404343	tomorrow_collapse	11.6020	0
116	○	2025-07-15	404344	tomorrow_collapse	11.6020	0
117	○	2025-07-15	404345	tomorrow_collapse	11.6020	0
118	○	2025-07-15	404346	tomorrow_collapse	11.6020	0
119	○	2025-07-15	404347	tomorrow_collapse	11.6020	0
120	○	2025-07-15	404348	tomorrow_collapse	11.6020	0
121	○	2025-07-15	404349	tomorrow_collapse	11.6020	0
122	○	2025-07-15	404350	tomorrow_collapse	11.6020	0
123	○	2025-07-15	404351	tomorrow_collapse	11.6020	0
124	○	2025-07-15	404352	tomorrow_collapse	11.6020	0
125	○	2025-07-15	404353	tomorrow_collapse	11.6020	0
126	○	2025-07-15	404354	tomorrow_collapse	11.6020	0
127	○	2025-07-15	404355	tomorrow_collapse	11.6020	0
128	○	2025-07-15	404356	tomorrow_collapse	11.6020	0
129	○	2025-07-15	404357	tomorrow_collapse	11.6020	0
130	○	2025-07-15	404358	tomorrow_collapse	11.6020	0
131	○	2025-07-15	404359	tomorrow_collapse	11.6020	0
132	○	2025-07-15	404360	tomorrow_collapse	11.6020	0
133	○	2025-07-15	404361	tomorrow_collapse	11.6020	0
134	○	2025-07-15	404362	tomorrow_collapse	11.6020	0
135	○	2025-07-15	404363	tomorrow_collapse	11.6020	0
136	○	2025-07-15	404364	tomorrow_collapse	11.6020	0
137	○	2025-07-15	404365	tomorrow_collapse	11.6020	0
138	○	2025-07-15	404366	tomorrow_collapse	11.6020	0
139	○	2025-07-15	404367	tomorrow_collapse	11.6020	0
140	○	2025-07-15	404368	tomorrow_collapse	11.6020	0
141	○	2025-07-15	404369	tomorrow_collapse	11.6020	0
142	○	2025-07-15	404370	tomorrow_collapse	11.6020	0
143	○	2025-07-15	404371	tomorrow_collapse	11.6020	0
144	○	2025-07-15	404372	tomorrow_collapse	11.6020	0
145	○	2025-07-15	404373	tomorrow_collapse	11.6020	0
146	○	2025-07-15	404374	tomorrow_collapse	11.6020	0
147	○	2025-07-15	404375	tomorrow_collapse	11.6020	0
148	○	2025-07-15	404376	tomorrow_collapse	11.6020	0
149	○	2025-07-15	404377	tomorrow_collapse	11.6020	0
150	○	2025-07-15	404378	tomorrow_collapse	11.6020	0
151	○	2025-07-15	404379	tomorrow_collapse	11.6020	0
152	○	2025-07-15	404380	tomorrow_collapse	11.6020	0
153	○	2025-07-15	404381	tomorrow_collapse	11.6020	0
154	○	2025-07-15	404382	tomorrow_collapse	11.6020	0
155	○	2025-07-15	404383	tomorrow_collapse	11.6020	0
156	○	2025-07-15	404384	tomorrow_collapse	11.6020	0
157	○	2025-07-15	404385	tomorrow_collapse	11.6020	0
158	○	2025-07-15	404386	tomorrow_collapse	11.6020	0
159	○	2025-07-15	404387	tomorrow_collapse	11.6020	0
160	○	2025-07-15	404388	tomorrow_collapse	11.6020	0
161	○	2025-07-15	404389	tomorrow_collapse	11.6020	0
162	○	2025-07-15	404390	tomorrow_collapse	11.6020	0

Discovery flow in data-driven fields

- Chemistry
- Biology
- Material Science
- Astrophysics
- Economics
- Quantitative Finance
- ...



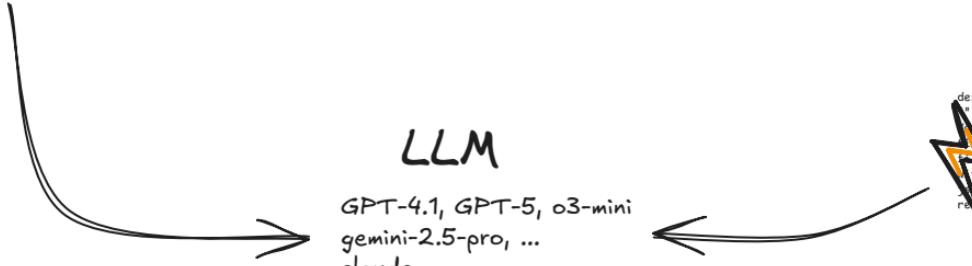
What is an Agent?



System Message

You are an idea maker agent.

You must provide a high quality set of ideas and update your ideas based on recommendations. Ideas should be based on the data/problem of interest, and feasibility given the data available



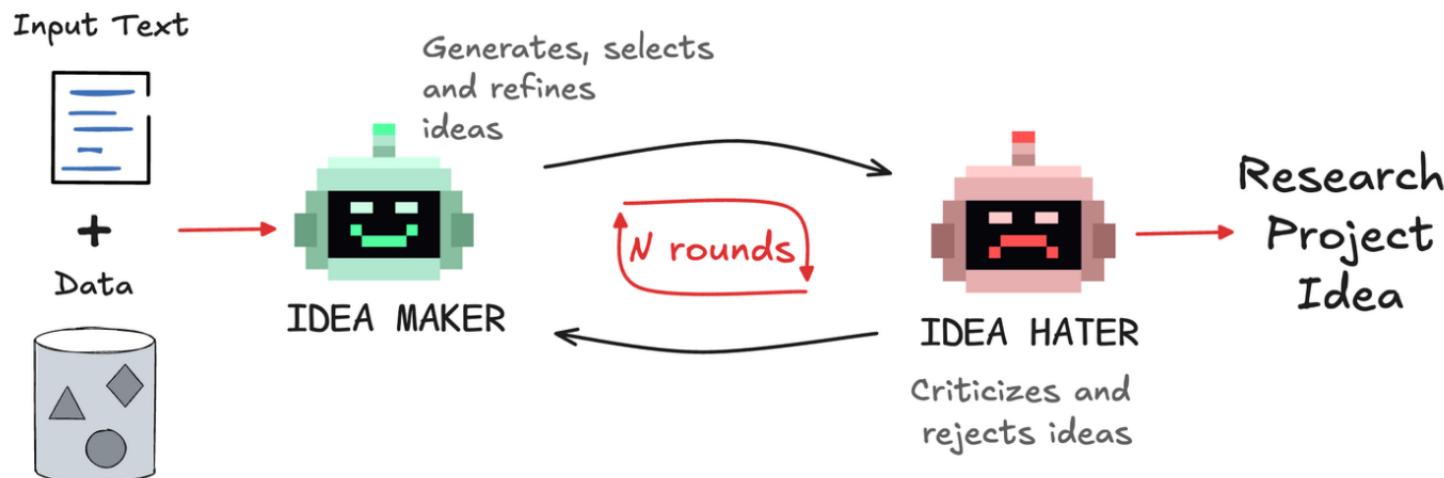
Tools

```
def record_ideas(ideas: list):
    """ Record ideas. You must record the entire list of ideas and their descriptions.
    You must not alter the list. """
    timestamp = datetime.datetime.now().strftime("%Y%b%d_%H%M%S")
    filepath = os.path.join(cmbagent_instance.work_dir, f'ideas_{timestamp}.json')
    with open(filepath, 'w') as f:
        json.dump(ideas, f)
    return f"\nideas saved in {filepath}\n"
```

An agent is an LLM instructed to play a role and to select tools

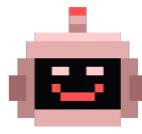


Example: Scientific Research Project Idea Generation



with Adrian Bayer
Flatiron/Princeton

A multi-agent system is a group of agents that collaborate



Idea Agent



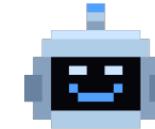
Literature Agent



Simulation agent



Paper writing Agent

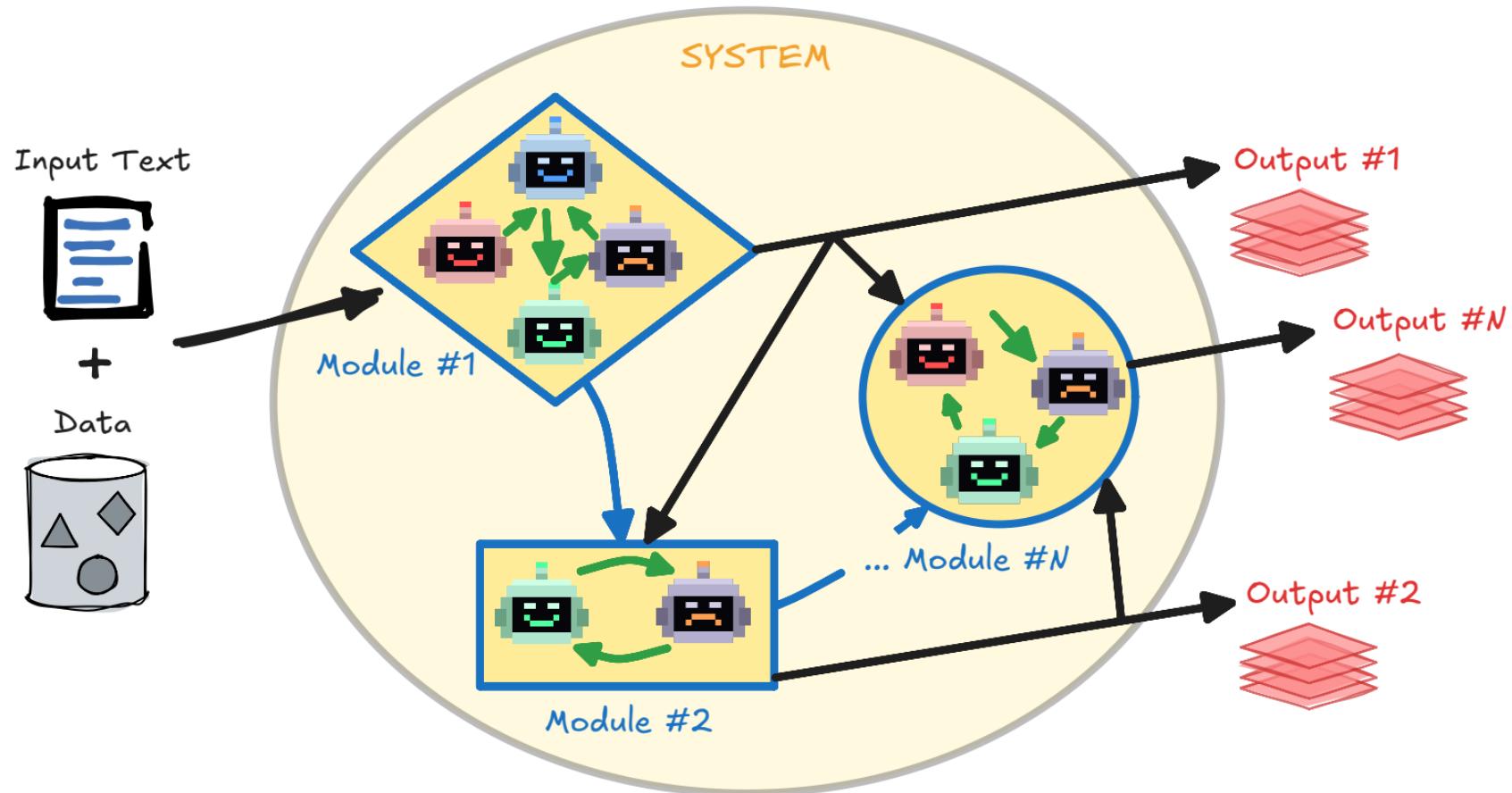


Referee Agent



...

The Lab of tomorrow



The DENARIO project: Modular Automation of Scientific Research with Multi-Agent Systems

Boris Bolliet*, Francisco Villaescusa-Navarro*, Pablo Villanueva-Domingo*,
Adrian E. Bayer, Aidan Acquah, Chetana Amancharla, Almog Barzilay Siegal, Pablo Bermejo,
Camille Bilodeau, Pablo Cárdenas Ramírez, Miles Cranmer, Urbano L. França, ChangHoon Hahn,
Yan-Fei Jiang, Raul Jimenez, Jun-Young Lee, Antonio Lerario, Osman Mamun, Thomas Meier,
Anupam Anand Ojha, Pavlos Protopapas, Shimanto Roy, Pedro Tarancón-Álvarez, Ujjwal Tiwari, Matteo Viel,
Digvijay Wadekar, Chi Wang, Bonny Y. Wang, Licong Xu, Yossi Yovel, Shuwen Yue, Wenhan Zhou, Qiyao Zhu,
Jiajun Zou, Íñigo Zubeldia



Francisco Villaescusa-Navarro



Pablo Villanueva-Domingo



*Equal Contribution. Listing order of BB, PVD, FVN is random.



SCAN ME



Google DeepMind



HARVARD
UNIVERSITY



PRINCETON
UNIVERSITY



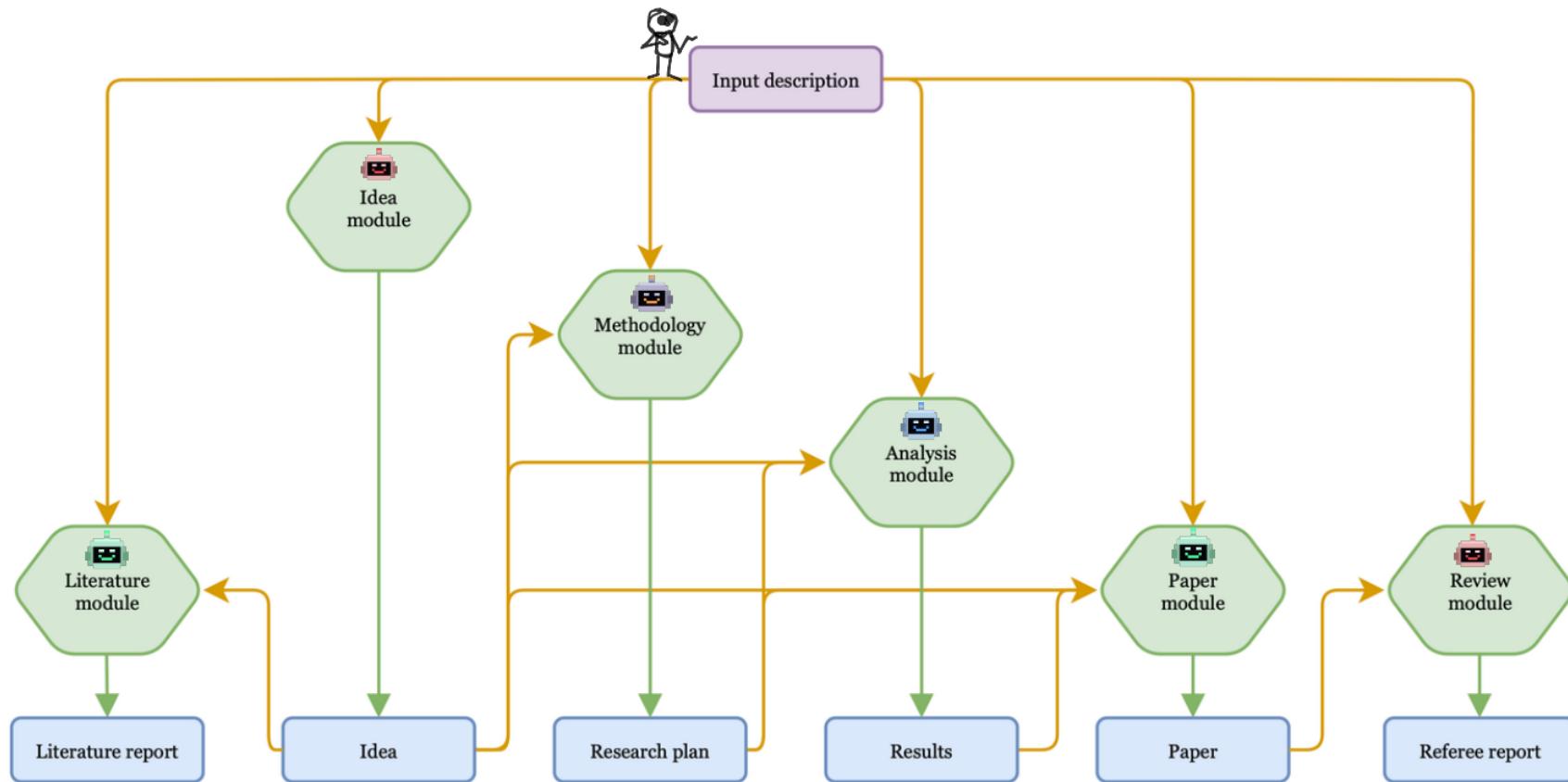
Massachusetts
Institute of
Technology



UNIVERSITY OF
OXFORD

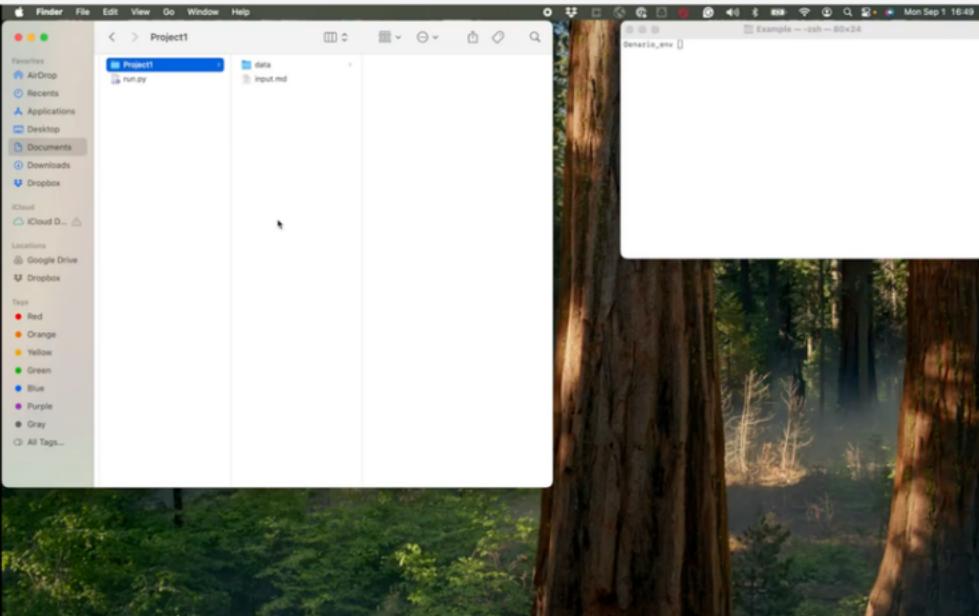


Denario: End-to-End Scientific Research



Denario - End to End Research

End-to-end research with Denario, from hypothesis generation to paper writing.



SCAN ME

Open Conference of AI Agents for Science 2025

The 1st open conference where AI serves as both primary authors and reviewers of research papers

Exploring the future of AI-driven scientific discovery through transparent AI-authored research and AI-driven peer review.



247 submissions
48 accept
Only 5 with >95% AI

<https://agents4science.stanford.edu>

All Papers	Accepted	Rejected	Total: 247 papers	Accepted: 48	Rejected: 199				
Paper Title	Status	Primary Topic	Secondary Topic	Human Review	AI Reviewer	Autonomy Scores			
				1	2	3	Hypothesis Development	Experimental Design	Data Analysis
PsySpace: Simulating Emergent Psychological Dynamics in Long-Duration Space Missions using Multi-Agent LLMs	Accepted	Computer & Data Sciences	Human-Computer Interaction	5	3	6	4	D	D
Green by Design: Energy-Guided Reranking of LLM-Generated Programs	Accepted	Computer & Data Sciences	Artificial Intelligence & Machine Learning	4	3	6	3	D	D
QITT-Enhanced Multi-Scale Substructure Analysis with Learned Topological Embeddings for Cosmological Parameter Estimation	Accepted	Computer & Data Sciences	Artificial Intelligence & Machine Learning	4	3	6	3	D	D

Deep Research: Cmbagent

Get the open-source code here:



Planning

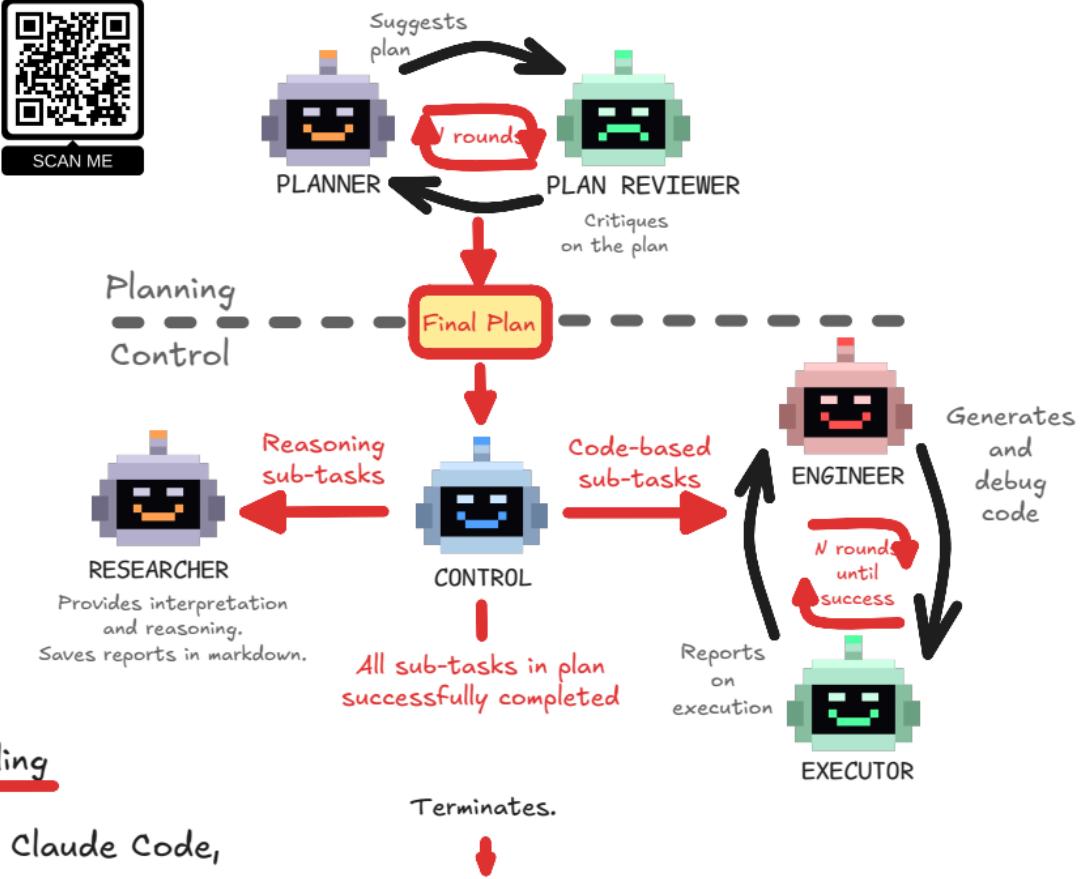
decompose main task into sub-tasks
propose-critique loop

Control/Execution

Solve each sub-task
generate-evaluate loop

Long workflows enabled by context handling

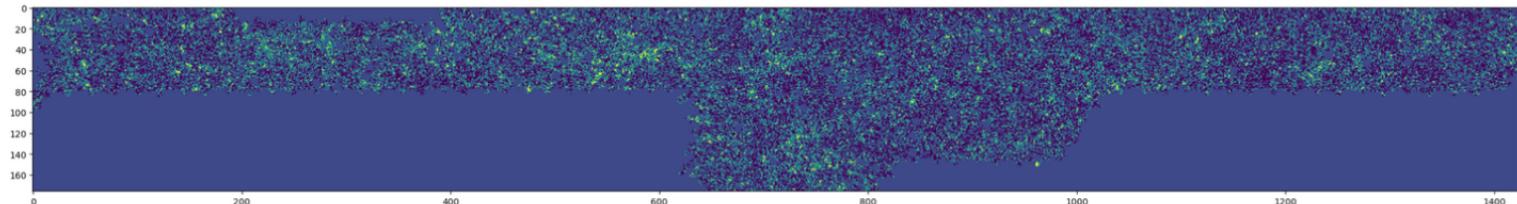
Other Example: ChatGPT DeepResearch, Claude Code,
OpenAI Codex, gemini cli, cursor agents,...
Ours is for scientific research



e.g., Xu et al (2025)
<https://github.com/CMBAgents/cmbagent>

FAIR-Universe NeurIPS Competition 2025

```
# noisy training convergence map
Visualization.plot_noisy_training_convergence_map(kappa=data_obj.kappa,
                                                    mask=data_obj.mask,
                                                    pixelsize_arcmin=data_obj.pixelsize_arcmin,
                                                    ng=data_obj.ng)
```

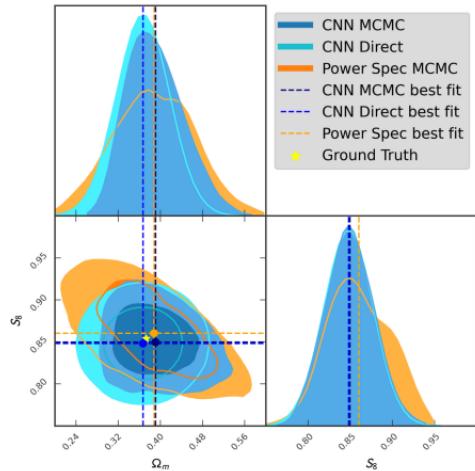
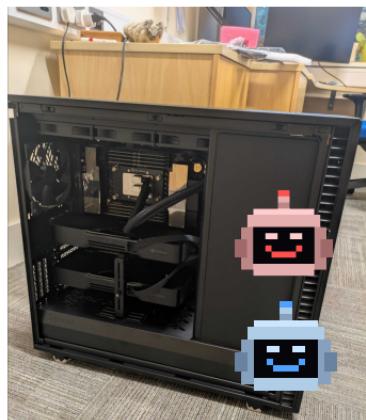


Input: Weak lensing maps
Output: Parameter value prediction

Strategy:

Human research followed by
Autonomous Program Synthesis
with Human-on-the-loop

Data volume ~1TB
Run time ~4h



Main-task Prompt

<TASK>

Find and train a neural network that maximises the score. Do better than current. Best model will achieve above 11.

The Simple_CNN model provided in the example yields a score of around 8.2-8.5. You don't need to re-do this, explore better alternatives from the start.

Note: `fair_universe` is a package from which methods can be imported. You are not operating from within the package, so you must not use relative imports.

Computational resources: We are running on an NVIDIA RTX PRO 6000 Blackwell Workstation Edition with 96GB RAM. Ensure this is used well.

Here, the task is not about refining the scoring script/parts, but about finding the model that yields the highest score.

</TASK>

<PREVIOUS RUN INSIGHTS>

<FIRST ITERATION>
Final Analysis and Actionable Recommendations for Achieving Target Score

1. Interpretation of Ensemble Model Performance

The previous iteration successfully implemented and evaluated an ensemble of four 'ResNet18' models, each trained on a unique noise realization of the dataset. This approach, combined with refined hyperparameters and the introduction of data augmentation, yielded a significant performance improvement over the initial single-model experiment.

- **Baseline 'Simple_CNN' Score:** ~8.2-8.5
- **Single 'ResNet18' (Model ID 1) Score:** 8.91
- **Ensemble 'ResNet18' Score:** **9.59**



Get the open-source code here

Planning Prompt

Use engineer for preprocessing, training and scoring. Use researcher for insights/interpretation and suggestions.

Step 1: Preprocess data and save preprocessed data for first model training. Use engineer.

Step 2: A first model is trained with engineer.

Step 3: The first model is scored with engineer.

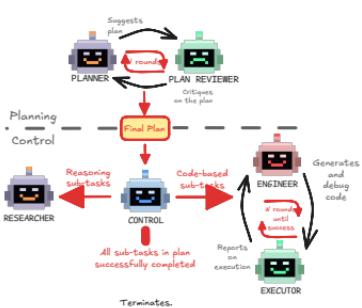
Step 4: Interpret results with researcher and suggest improved model based on these results.

Step 5: Preprocess data according to recommendations and findings. Use engineer.

Step 6: Train improved model with engineer.

Step 7: Score improved model with engineer.

Step 8: Review results and provide insights on what to try next (assuming the same workflow and constraints).



Xu et al (2025)
<https://github.com/CMBagents/cmagent>

445	eff	2025-10-21 21:07	2099810	EffL_mse0.3_MNN	10.9162
446		2025-11-15 16:57	421419	666	10.5339
447	mishraani	2025-11-15 16:57	421439	464444	10.9199
448	finayak	2025-10-22 09:22	2099801	uploaded	10.9126
449	ashwagell	2025-11-04 04:53	410889	a	10.9034
450	ashwagell	2025-11-05 04:53	410899	c	10.9034
451	cmagent	2025-10-27 18:37	3988004	cmagent	10.901
452	ashwagell	2025-11-04 22:45	422580	fb11111	10.9009
453	yanling007	2025-10-31 21:58	407864	25-10-31-17-52	10.9003
454	datame	2025-11-01 01:11	423603	BOT	10.4942
455	zhukuli	2025-11-05 22:45	417136	CNN 05	10.4993
456	remay	2025-11-06 01:00	415600	per_cnn_sqe_mse	10.4883

FAIR-Universe NeurIPS Competition 2025

202

PARTICIPANTS

1760

SUBMISSIONS

1. Final leaderboard evaluated solely on (i):

RANK	PARTICIPANT	FINAL SCORE	MEAN MSE (STANDARDIZED)	MEAN COVERAGE
1st	cmbagent	11.7029	0.1033	0.7000
2nd	eiffi	11.6535	0.1038	0.7087
3rd	Shubhosit	11.5987	0.1032	0.6583

We will award the prizes to **cmbagent**, **eiffi**, and **Shubhosit** for extraordinary performance on the original cosmologies.

- **cmbagent**: team members Erwan Ally, Boris Bolliet, Tom Borret, Celia Lecat, Andy Nilipour, Sébastien Pierre, Licong Xu
- **eiffi - Transatlantic Dream Team**: team members Noe Dia, Sacha Guérini, Wassim Kablan, François Lanusse, Julia Linhart, Laurence Perreault-Levasseur, Benjamin Remy, Sammy Sharieff, Andreas Tersenov, Justine Zeghal
- **shubhosit** - Shubhosit Naskar

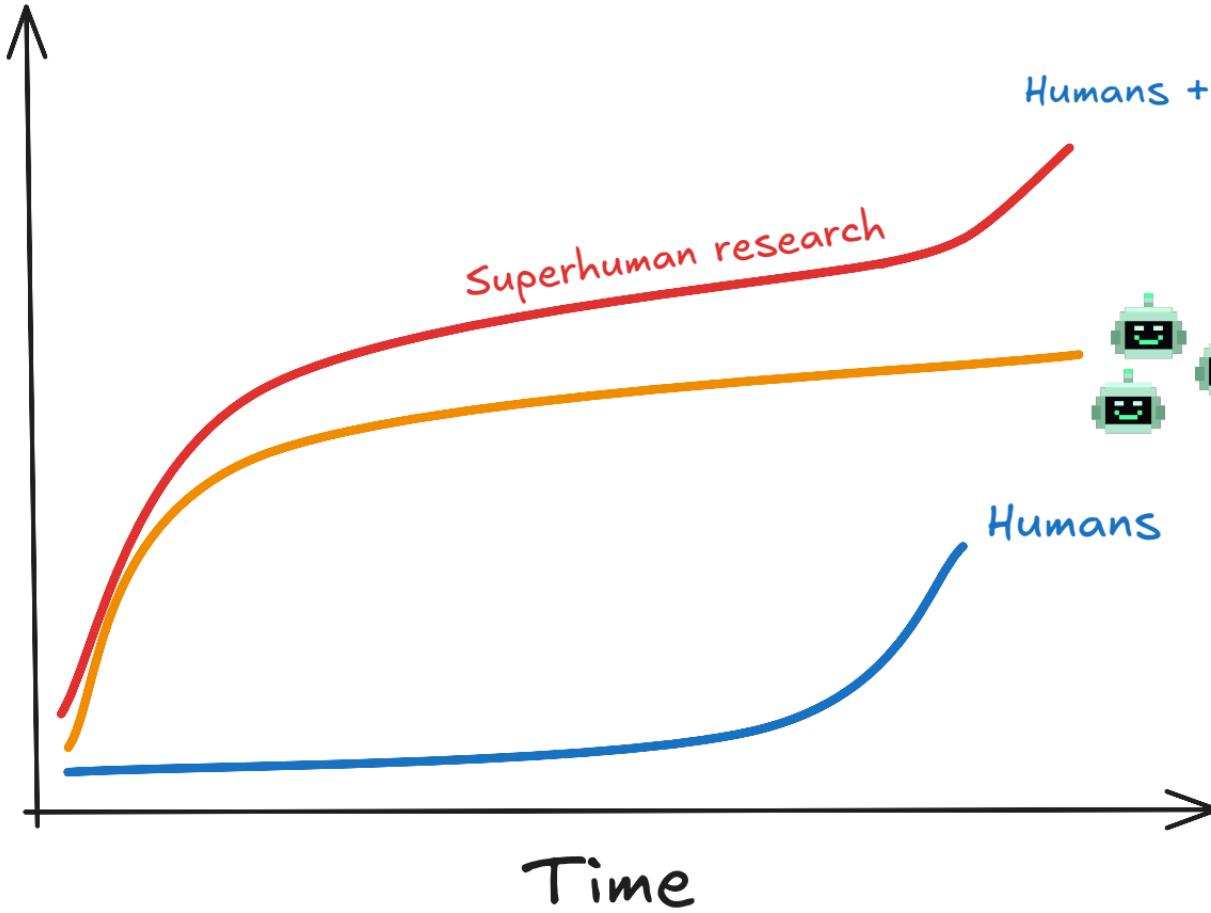
Final leaderboard (i)

Open-phase leaderboard

RANK	USERNAME	SUBMISSION ID	SCORE	MSE	R ²	COVERAGE	METHOD NAME
1	cmbagent (team)	424324	11.7321	0.1022	0.8958	0.7056	tomborrett_cmbagent
2	eiffi (team)	424310	11.6612	0.1028	0.8952	0.7095	b-remy_e36uxobt
3	shubhosit	417744	11.6192	0.1025	0.8956	0.6583	CNNv48_
4	THUML (team)	424064	11.5209	0.1064	0.8916	0.6907	F1
5	adscft	424300	11.5142	0.1056	0.8924	0.7279	cnnv11152
6	piyush555	394590	11.4681	0.1077	0.8902	0.7103	dns_l
7	jhu_suicee	423555	11.4590	0.1077	0.8903	0.6512	STILI
8	azhang81	424251	11.4192	0.1098	0.8882	0.7017	m6
9	mmayr	424022	11.2759	0.1133	0.8846	0.7080	NL128_LD512_NB2_NSA6_NH16_LR2e4_BS32_NP200_VPRC
10	jagongcalves	418273	11.2437	0.1160	0.8818	0.6835	20251008_173243
11	DOT (team)	422646	11.2203	0.1160	0.8818	0.7016	CNN

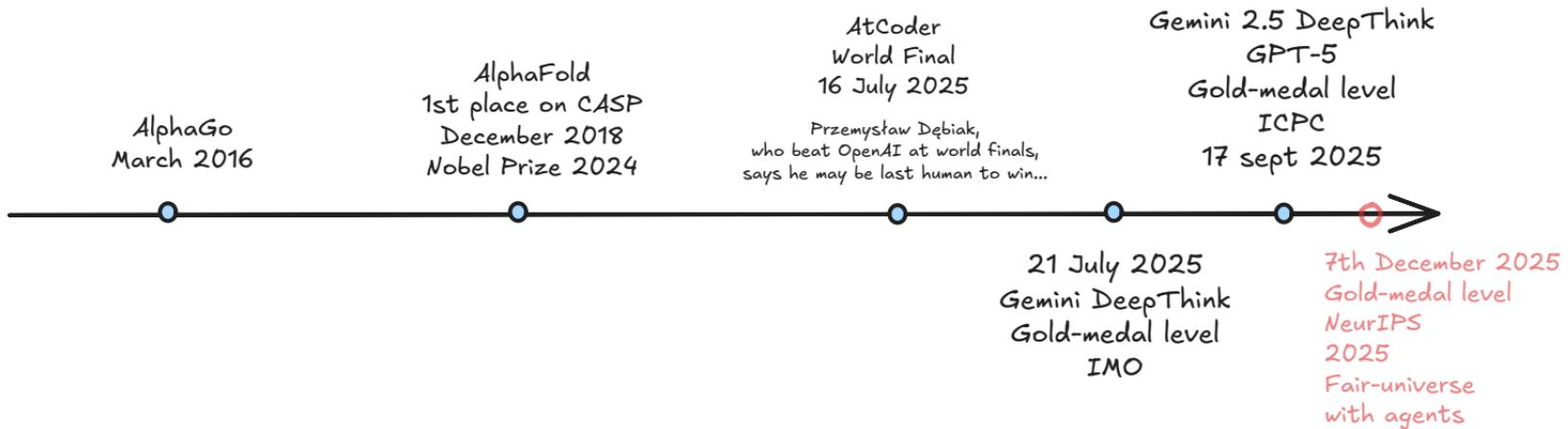
#	Participant	Date	ID	Method Name	Score
1	cmbagent	2025-11-15 23:48	424324	tomborrett_cmbagent	11.7321
2	cmbagent	2025-11-15 23:44	423570	andynilipour_cmbagent	11.7257
3	cmbagent	2025-11-15 23:11	424297	tomborrett_cmbagent	11.7226
4	cmbagent	2025-11-14 23:36	423521	tomborrett_cmbagent	11.7203
5	cmbagent	2025-11-15 18:18	424175	licongxiu_cmbagent	11.7188
6	cmbagent	2025-11-14 11:19	422877	licongxiu_cmbagent	11.7166
7	cmbagent	2025-11-15 18:06	424162	andynilipour_cmbagent	11.695
8	cmbagent	2025-11-14 18:58	423244	licongxiu_cmbagent	11.6879
9	cmbagent	2025-11-13 13:03	422151	licongxiu_cmbagent	11.6686
10	cmbagent	2025-11-12 22:50	421614	licongxiu_cmbagent	11.6628
11	eiffi	2025-11-15 23:31	424310	b-remy_e36uxobt	11.6612
12	cmbagent	2025-11-13 23:21	422546	licongxiu_cmbagent	11.6602

Understanding & Performance



It's not about automating menial tasks... Not about replacing students/researchers...

It's about building systems that do research with us, better than us alone.



Research is a competition.

Critique System

skeptical
The AI reviewer

 Paper loaded: 2106.03846.pdf
Drag a new file here or [change paper](#)

Select review type (at least one)

- Full review
- Logic
- Evidence
- Interpretation, clarity and reproducibility
- General consistency
- Figures
- Statement verification

Review thoroughness

- Standard
- High
- Very high

Your email address (or addresses)

The review should reach you in a few minutes.

Type an email and press space

Start review

Paper uploaded. Ready to review.

SKEPTICAL = SKEPTIC + ETHICAL

Skeptical review: The superclustering of hot gas: cosmological sensitivity in the Websky simulations

Summary

The manuscript presents an oriented tSZ stacking analysis using cosmological simulations, with particular emphasis on multipole radial profiles, integrated multipole power, environmental constraints, and cosmology-gas degeneracies. The introductory and methodological sections appear well developed. The cosmological parameters and key foundational and recent works are appropriately cited there. The work is potentially impactful, especially regarding the use of higher-order moments and environmental selections to probe cosmology and gas physics, and could be a valuable contribution if its evidential basis is strengthened.

However, there are extensive gaps in referencing and external validation throughout the core scientific sections, particularly in "2.1 Varying Λ CDM parameters", "2.2 Fixing σ_8 at $z = 0.5$ ", "2.3 The Cosmo2 variations", "2.4 The Cosmo2 variations", "2.5 A quick look at the simulations", "3 Stacking methods", "4.1 Multipole radial profiles", "4.2 Integrated multipole power", "4.3 Source of higher-order moment information", "4.4 Environmental constraints", "5.4 Degeneracy with cosmology", and "6 Conclusions". Many statements about cosmological effects, simulation behavior, stacking methodology, performance, and the relative roles of gas physics and cosmology are not supported by citations, benchmarking, or explicit qualification of novel empirical findings. In addition, several specific statements are only partially supported or incorrectly referenced, including claims about the intergalactic medium and halo-based gas pasting, pressure-profile prescriptions and their literature basis, cosmological parameter constraints, and the interpretation of anisotropic structures and AGN feedback effects.

The lack of external validation and shortcomings substantially limit the reliability and generalizability of the paper's main conclusions, particularly regarding cosmological sensitivity, degeneracy-breaking via higher-order moments, and the impact of environmental constraints and gas modeling. A major revision is needed to (i) add appropriate references, (ii) correct or qualify partially supported statements, and (iii) clearly distinguish between well-established results, prior literature, and genuinely new empirical findings of this work. Once these issues are addressed, the manuscript's scientific contribution and impact would be significantly strengthened.

Major issues

1. Widespread lack of citations for key cosmological evolution statements (Section 2.1-2.4): In "2.1 Varying Λ CDM parameters", the description of how changing Ω_m affects the transitions between radiation-, matter-, and dark-energy-dominated eras and the matter power spectrum turnover is presented without explicit literature support. In "2.2 Fixing σ_8 at $z = 0.5$ ", while the dependence of the tSZ power spectrum on σ_8 is referenced, the subsequent claim that the anisotropic tSZ signal specifically traces differences in structure

1



Inigo Zubeldia
(Cambridge)

Multi-Modal Reasoning (2511.14631)

Enhancing Agentic Autonomous Scientific Discovery with Vision-Language Model Capabilities

Kahaan Gandhi^{1,2,3} Boris Bolliet^{2,4} Iñigo Zubeldia^{4,5}

Abstract

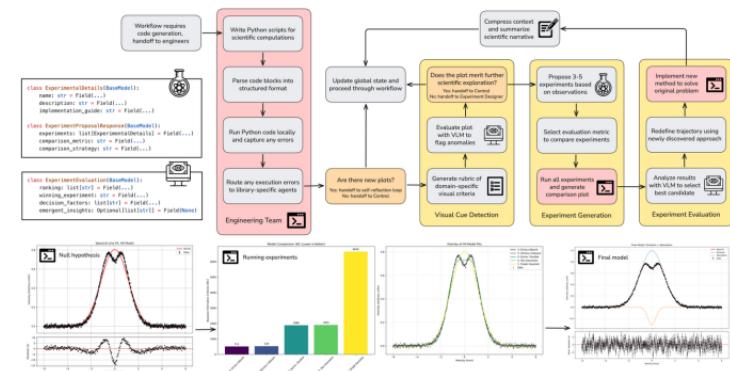
We show that multi-agent systems guided by vision-language models (VLMs) improve end-to-end autonomous scientific discovery. By treating plots as verifiable checkpoints, a VLM-as-a-judge evaluates figures against dynamically generated domain-specific rubrics, enabling agents to correct their own errors and steer exploratory data analysis in real-time. Case studies in cosmology and astrochemistry demonstrate recovery from faulty reasoning paths and adaptation to new datasets without human intervention. On a 10-task benchmark for data-driven discovery, VLM-augmented systems achieve pass@1 scores of 0.7–0.8, compared to 0.2–0.3 for code-only and 0.4–0.5 for code-and-text baselines, while also providing auditable reasoning traces that improve interpretability.

soning and communication are more subjective and require discretion. When orchestrated into multi-agent systems for end-to-end automation, these harder-to-verify tasks often emerge as failure modes.

For autonomous systems to become credible scientific collaborators, they must move beyond analysis alone and communicate findings in ways interpretable to the research community. In data-intensive fields, figures are the primary medium for both communication and interpretation. They compress large datasets into digestible representations while also guiding the research process: plots reveal anomalies, prompt the reconsideration of hypotheses, and steer subsequent steps. This feedback loop is central to human discovery workflows but remains largely absent in current systems, where frontier models fail to handle domain-specific conventions in plots (Joseph et al., 2025).

To address this gap, we extend `cmbagent`, a fully au-

Coding Agent	Judge	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8	Q9	Q10	pass@1
Gemini 2.5 Pro	None	✓	✗	✗	✗	✗	✗	✗	✓	✗	✗	0.2
GPT-4.1	None	✓	✗	✗	✗	✗	✗	✗	✗	✗	✓	0.2
Claude Opus 4.1	None	✓	✗	✗	✗	✗	✗	✗	✓	✗	✓	0.3
<code>cmbagent</code> + Gemini 2.5 Pro	None	✓	✗	✓	✓	✓	✓	✓	✓	✓	✗	0.5
<code>cmbagent</code> + GPT-4.1	None	✓	✗	✓	✓	✓	✓	✓	✓	✓	✗	0.4
<code>cmbagent</code> + GPT-4.1	GPT-4.0 (LLM)	✓	✗	✓	✓	✓	✓	✗	✗	✗	✓	0.5
<code>cmbagent</code> + GPT-4.1	Gemini 2.5 Pro (LLM)	✓	✗	✓	✓	✗	✗	✗	✓	✗	✓	0.5
<code>cmbagent</code> + GPT-4.1	GPT-4.0 (VLM)	✓	✓	✓	✓	✗	✓	✗	✓	✓	✓	0.7
<code>cmbagent</code> + GPT-4.1	Gemini 2.5 Pro (VLM)	✓	✓	✓	✓	✓	✓	✗	✓	✗	✓	0.8



Problem Statement: Given a new dataset, test whether the null hypothesis remains supported or should be rejected in favor of an alternative model.

Hypothesis (H_0): The spectral line is modeled as a single Gaussian profile on a constant continuum with independent Gaussian noise:

$$I(v; \theta) = c_0 + A \exp \left[-\frac{(v - \mu)^2}{2\sigma^2} \right].$$

Prior Context: Prior datasets did not provide sufficient evidence to reject H_0 .

New Dataset: `path/to/data.npz` with keys “`v`” (velocity), “`I`” (intensity), and “`sigma`” (per-channel noise).

Tasks: Test H_0 against the new dataset. If rejected, identify and fit an alternative line-profile model.

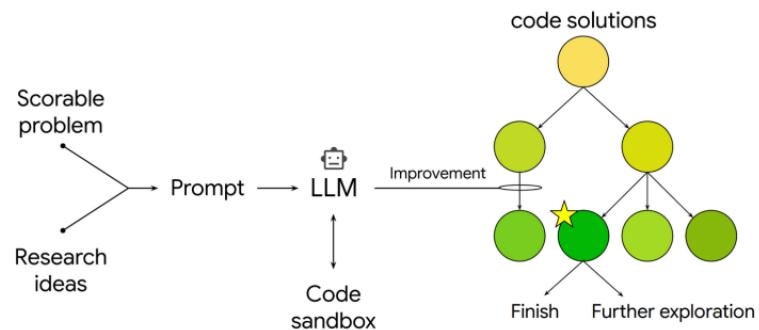


Kahaan Gandhi
(Caltech)

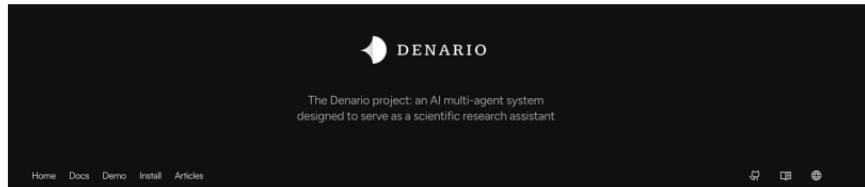
Other areas of ongoing efforts:

- Data mining + hypothesis generation
(w. U. Demirbozan, G. Farren)
- Program & Agent Synthesis
(New York x Cambridge)
- Disseminating access
(Platforms for research)
- More competitions

Stay tuned!



(Aygun,.., Brenner, Google: 2509.06503)



Agents for Scientific Discovery

Boris Bolliet



UNIVERSITY OF
CAMBRIDGE



ACCELERATE
PROGRAMME
FOR SCIENTIFIC DISCOVERY