FOSS India

DiscoveryBench: Data-Driven Scientific Discovery with LLM Agents



Harshit Surana

Allen Institute for Artificial Intelligence & OpenLocus









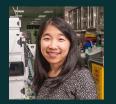




























Seattle Hub









Structure & Some Requests

- 1. Scientific Discovery with LLMs
- DiscoveryBench/ Data Driven Discovery
- 3. Best Practices in Building Al at Scale

- Important to have interactivity
- Feel free to interrupt for quick clarifications.
- Will stop at the end of each section for QnA
- If audio issues arise, please drop messages in the chat.



Scientific Discovery with LLMs

"Superintelligent tools could massively accelerate scientific discovery and innovation well beyond what we are capable of doing on our own, and in turn massively increase abundance and prosperity."

- Sam Altman, OpenAl



NOBELPRISET I KEMI 2024 THE NOBEL PRIZE IN CHEMISTRY 2024





David Baker
University of Washington
USA

"för datorbaserad proteindesign"

"for computational protein design"



Demis Hassabis Google DeepMind United Kingdom



John M. Jumper Google DeepMind United Kingdom

"för proteinstrukturprediktion"

"for protein structure prediction"



Broad Areas of Scientific Discovery where LLMs are Helping!

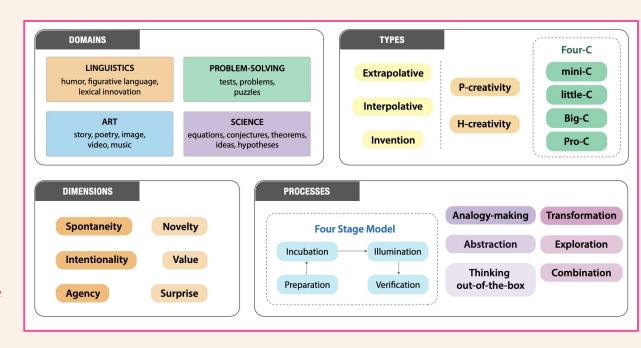
Laboratory Tools

Exp. Chemistry

Equation Discovery

Theorem Proving

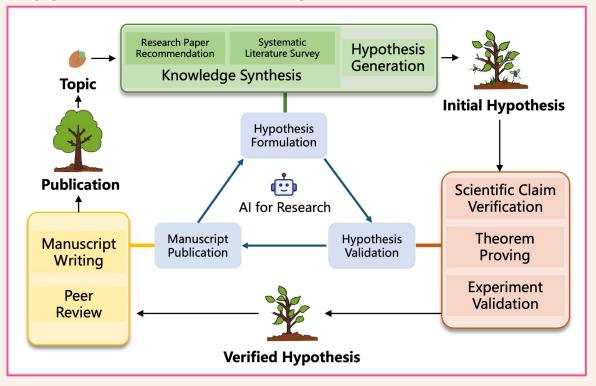
Observational Data*



What are the steps in any Research Process?

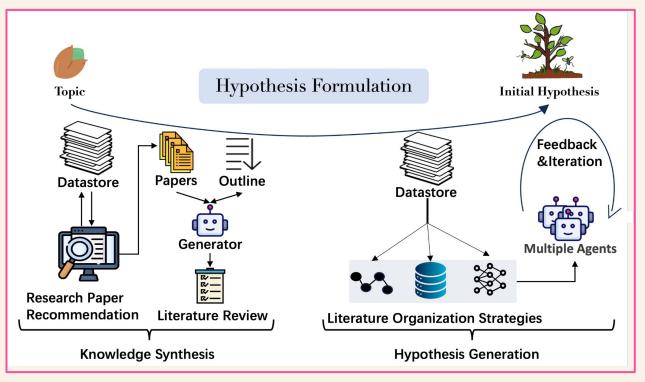
What are different methods of verifying a Hypothesis?

Topic → Hypothesis → Paper



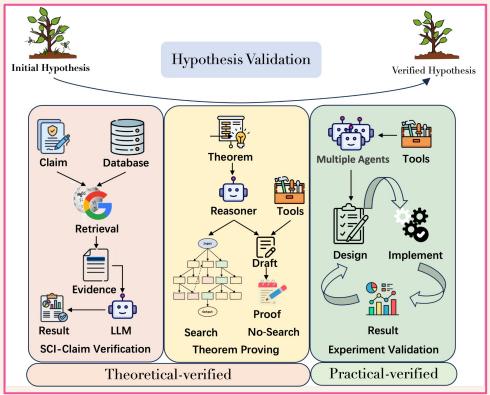
From Hypothesis to Publication: A Comprehensive Survey of Al-Driven Research Support Systems, Arxiv 2025

Hypothesis Formulation



From Hypothesis to Publication: A Comprehensive Survey of Al-Driven Research Support Systems, Arxiv 2025

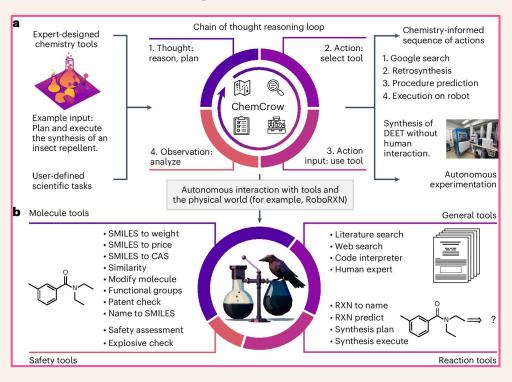
Hypothesis Validation



From Hypothesis to Publication: A Comprehensive Survey of Al-Driven Research Support Systems, Arxiv 2025

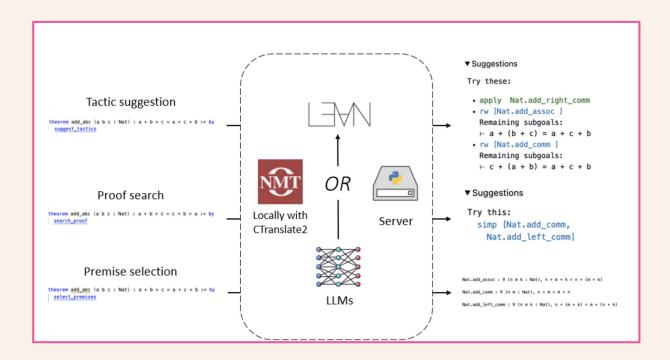


LLMs, Lab Tools & Beyond



ChemCrow: Augmenting large language models with chemistry tools, Nature 2024

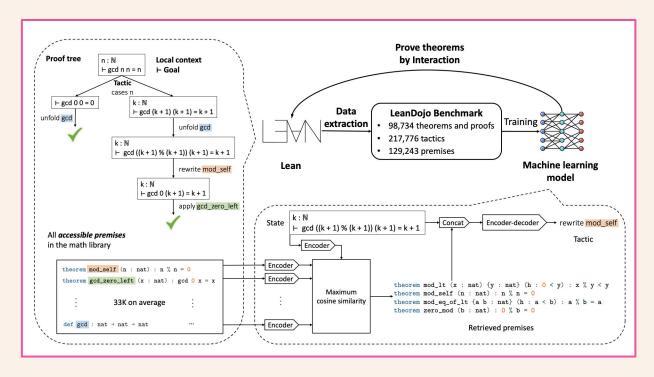
Theorem Proving Copilot



Towards Large Language Models as Copilots for Theorem Proving in Lean, NuerIPS 2023



Theorem Proving Env with LLMs + Retrieval



LeanDojo: Theorem Proving with Retrieval-Augmented Language Models, NuerIPS 2023



Methods of Science

Methods of Scientific Inquiry

Theoretical Science

Develop models or theories to explain phenomena

Experimental Science

Conduct experiments to test pre-defined hypotheses

Observational Science

Observe & collect data, build methods to explain it

Methods of Scientific Inquiry

Theoretical Science

Develop models or theories to explain phenomena

Experimental Science

Conduct experiments to test pre-defined hypotheses

Observational Science

Observe & collect data, build methods to explain it

A lot of important science has come out of looking at **observational data**.

National Longitudinal Survey of Youth | 1979



500,000 results in S2 from 1979



37000+ papers published from 1976

Methods of Scientific Inquiry

Theoretical Science

Develop models or theories to explain phenomena

Experimental Science

Conduct experiments to test pre-defined hypotheses

Observational Science

Observe & collect data, build methods to explain it

A lot of important science has come out of looking at **observational data**.

Can we **autonomously** discover

- insights from datasets to reduce turnaround time?
- undiscovered knowledge without performing additional data collection?

National Longitudinal Survey of Youth | 1979



500,000 results in S2 from 1979



37000+ papers published from 1976

Data Driven Scientific Discovery with LLMs

Data-driven Discovery

- Comprehensive data-understanding
- Ex-ante hypothesis search/generation
- Planning & orchestrating research pathways
- Execute & verify candidate hypotheses
- Accommodating human feedback
- Reproducible and robust results

Data-driven Discovery: Following Newell & Simon (1976), we define a <u>heuristic</u> <u>search problem</u> that aims to describe a given set of observations by uncovering the laws that govern its data-generating process.

E.g., "under context c, variables v have

relationship r"

Newell, A. and Simon, H. A. Computer science as empirical inquiry: symbols and search. Commun. ACM, 1976



Dataset: National Longitudinal Surveys

Query: Study the relation between BMI and Time Preference.



2. Per

Planner

Data Understanding

based on the results

User

Formulating Initial Hypotheses
Multi-step Planning

User probes more

More interdiscplinary insights

4 Data Expert interpets

The **correlation** coefficient: -0.031, very weak negative linear relationship between dissaving and BMI.

The interaction term coefficient: 0.5259 statistically significant (p < 0.0000) ...



Hypothesis Verification and Analysis, Reproducible Results



Please connect BMI with graduation, family & demographic data, run more sophisticated model.



8 Planner replans

SES: Compare association between subject variables based on SES
 SAMPLE_SEX
 College Scores, Class Percentile

4. SAMPLE_RACE

(#) Planner

Data Understanding, Accomodating Human Feedback

111 Data Expert interprets

have a higher BMI than females.

"GENDER MALE" has a significant positive

association with BMI, indicating that males

The GLM confirms the findings from the OLS

model regarding the interactions between time preference and demographic factors.

10 Programmer executes

add_interactions(), run_glm()



Hypothesis Verification and Analysis, Data Transformation, Reproducible Results 2 Planner plans

Time preference could be 'DISSAVED'

and 'SAMESAVE' variables.

1. Initial Hypotheses:

a. Hypothesis 1: DISSAVED and BMI

are related...

2. Perform OLS & Correlation analysis...



Programmer

Hypothesis Verification and Analysis, Reproducible Results

3 Programmer executes

run_correlation(),

run_ols()

6 Data Expert proposes

Economics and Health Economics: Job status and income levels can affect health Psychology and Behavioral Economics: Stress, self-control influence saving habits and BMI... Sociology and Cultural Studies: Cultural norms and societal expectations can affect BMI....



Interdisciplinary Knowledge Integration

9 Data Expert directs

Programmer, please transform the data by adding interaction variables

Measure effects using Generalized Linear Model on 'SES', 'SAMPLE_SEX', 'SAMPLE_RACE', 'AVSAB Scores' and 'Class Percentile'



Hypothesis Verification and Analysis, Data Transformation, Reproducible Results

12 User poses question

How to mitigate the effect of testing multiple hypotheses?

User

Data Expert

Hypothesis Verification and Analysis, Data Transformation, Reproducible Results

Data-driven Discovery as a Predictive Task

Given a dataset D and a Discovery Goal G, derive the most specific hypothesis H addressing G and supported by D.

Alternatively,

A data-driven hypothesis H is a <u>declarative sentence</u> about the state of the world whose truth value may be inferred from a <u>given dataset</u> \underline{D} using a verification procedure $V: H \rightarrow \{\text{supported}, \text{unsupported}\}$, for instance, via *statistical modeling*.

Inspired by Thompson and Skau (2023), we introduce a structured formalism that breaks a hypothesis down into three hypothesis dimensions:

Context: Boundary conditions that limit the scope of a hypothesis. E.g., "for men over the age of 30"

Variables: Known set of concepts that interact in a meaningful way under a given context to produce the hypothesis. E.g., gender, age, or income

Relationship: Interactions between a given set of variables under a given context that produces the hypothesis. E.g., "quadratic relationship", "inversely proportional", or piecewise conditionals

Dataset:							
habitat type	nonnative gardening	nonnative unintentional		elevation			
croplands	5	0	2	675			
wetlands	0	4	1	88			
urban	2	1	0	329			
			•••				

Goal: How did urban land use affect the invasion of different types of introduced plants in Catalonia?

	gold	predicted	score			
context	urban habitat type	urban habitat type	1.0			
variable gardening, unintentional		gardening, agriforst	0.3			
relationship	reduced	increased	0.0			
Final Score: 0.21						

W. H. Thompson and S. Skau. On the scope of scientific hypotheses. Royal Society Open Science, 2023



Urban land use reduced invasion by gardening plants over unintentionally introduced ones.

DiscoveryBench

264 Tasks, 20+ papers, 6 domains

We replicate the **scientific process** undertaken by researchers to search for and validate a hypothesis from datasets

Data-first: Filter papers + workflows based on public datasets: National Longitudinal Surveys, Global Biodiversity Info Facility, World Bank Open Data; 2) replicate in Python.

Replication took up to 90 person-hours per dataset, often (30%) not resulting in success.

Code-first: Checked 785 repos + datasets, 85% had missing or non-adaptable code to Python, or closed datasets. Only few passed the check.

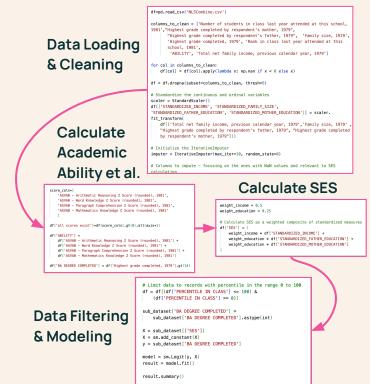
Papers from Nature, AER, etc.

Task Dataset: Dataset contains information from National Longitudinal Survey of Youth (NLSY79). It includes information about the Demographics, Family Background, Education ...

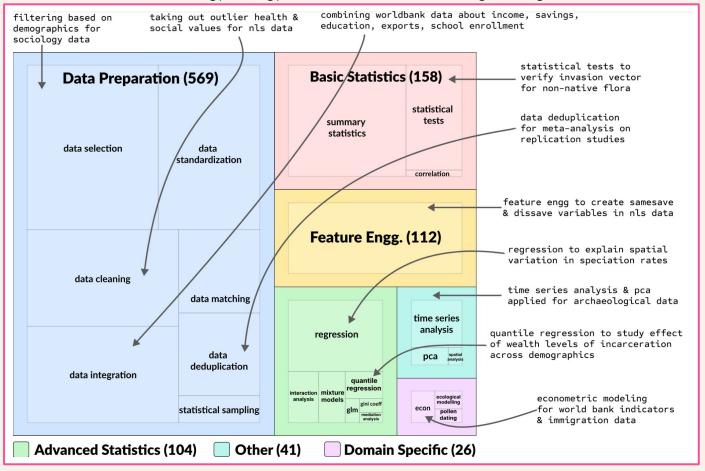
Discovery Goal: How does socioeconomic status affect the likelihood of completing a BA degree?

Target Hypothesis:

Socioeconomic status has a positive relationship with college degree completion with a coefficient of 0.4729 with statistical significance.



DB-Real (6 domains: sociology, biology, humanities, economics, engineering, & meta-science)



Discovery Agents

All discovery agents have access to a python environment, capable of generating and executing programs on the datasets

CodeGen

generates the entire code at one go to solve the task, with help of a demonstration example in the context.

After code execution and based on the result, it generates the NL hypothesis and summarizes the workflow

ReAct

solves the task by generating thought and subsequent codes in a multi-turn fashion.

A traditional sequential-decision maker.

DataVoyager*

is a multi-component data-driven discovery agent.

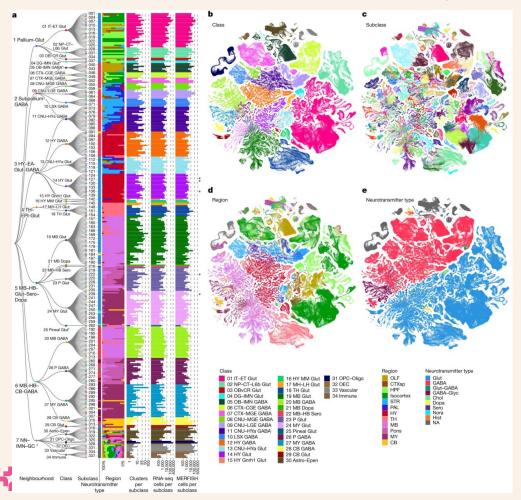
It has four components: planner, code generator, data analysis, and critic, that orchestrate the discovery process.

Reflexion (Oracle)

is an extension of CodeGen agent, where at the end of one trial, we provide an "oracle" feedback about task completion, and it generates a reflection to improve in the next trial till it solves the task, or maximum trials (3) are reached.

A NIMHANS (neuro) Example?

Culmination of decades of work summarized in an image



Landmark Mouse Brain Paper (Ai1)

DataVoyager Prompt:

- 1. Find the unique classes present in the data.
- 2. Select the "20 MB GABA" class and plot the distribution of subclasses and neurotransmitters.
- 3. Select the "23 P Glut" class and plot the distribution of subclasses and neurotransmitters.
- 4. Plot the supertypes and neurotransmitters distribution for the same filtered class.
- 5. Interpret the unique supertypes for this filtered class.

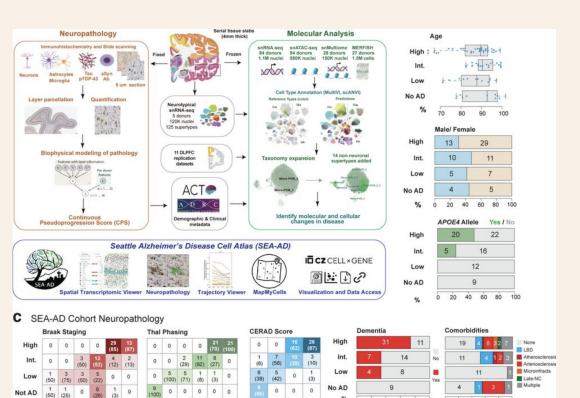
A high-resolution transcriptomic and spatial atlas of cell types in the whole mouse brain

https://www.nature.com/articles/s41586-023-06812-z

DataVoyager Traces for WMB Data







Integrated multimodal cell atlas of Alzheimer's disease

Data Input: Not sure of exact data used in paper or preprocessing. We currently used the following:

Donor Data

66 columns, and info includes metadata for each donor such as 'Primary Study Name', 'Age at Death', 'Sex', 'Race', 'CERAD score', 'Overall CAA Score', 'Highest Lewy Body Disease', 'Total Microinfarcts', 'Atherosclerosis', 'Arteriolosclerosis', 'LATE', 'RIN', and whether the donor is 'Severely Affected'.

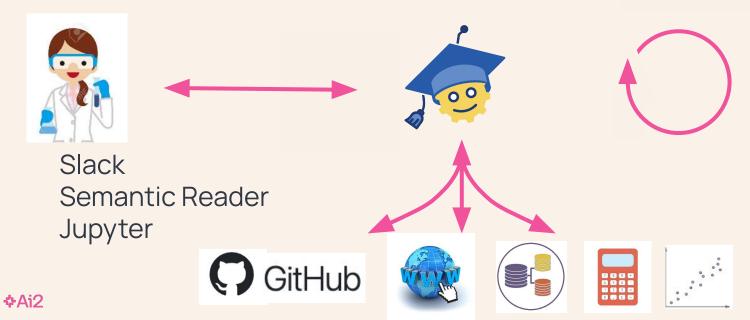
MTG Data

394 columns, measurements related to AT8 positive areas in different layers of Grey matter, pTDP43 positive areas, and other neuropathological quantifications for each donor.

DataVoyager is a part of Asta - a Science Copilot

- 1) Collaborates naturally with human scientists
- 2) Uses tools on humans' behalf

3) Learns & improves over time.



... with Infra. for Personalization & Optimization

Research Functions



Surfaces

Slack Web Semantic Reader Jupyter

Interaction Store

Evaluation Harness

Some Accolades

New Directions in the Area Influenced by our Work



BLADE: Benchmarking Language Model Agents for Data-Driven Science





POPPER: Automated Hypothesis Validation with Agentic Sequential Falsifications





ScienceAgentBench: Toward Rigorous Assessment of Language Agents for Data-Driven Scientific Discovery





QR Data: Are LLMs Capable of Data-based Statistical and Causal Reasoning? Benchmarking Advanced Quantitative Reasoning with Data

Recognized with other major works including the aforementioned Nobel Prize in 2024

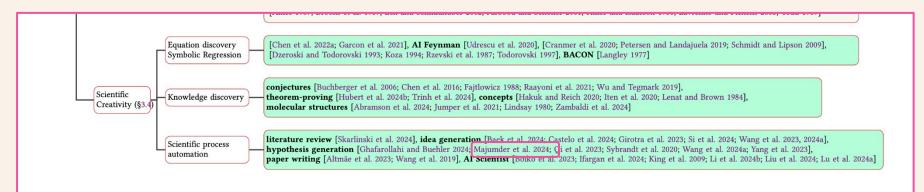


Fig. 2. Taxonomy of creativity in AI covering areas of linguistic creativity, creative problem-solving, artistic and scientific creativity. Note that this taxonomy is not exhaustive, but rather a representative view of the key works.

Creativity in AI: Progresses and Challenges, EPFL Zurich. Mete Ismayilzada, Debjit Paul, Antoine Bosselut, Lonneke van der Plas

Learnings

Learnings for in LLMs + Open Source

Instrument the **Data Pipeline**, Not Just the Model

Building **Evaluations** that Actually Matter

Adopting & Extending Frameworks (LangChain, AutoGen, etc.)

Treat **Prompts** and Pipelines as **Code Assets**

Exploit Small Models First, Big Models Last

Exploit Big Models First, Small Models Last

A/B (or **Ablations**) Everything Behind a Feature Flag/ Config



How to contribute to fundamental Al from India?

What happens with petabytes of data?

Or millions of lines of code?

Appendix

DataVoyager

All of the subagents and the controller is based on with a specialized prompt.

They see the full history, but the agent-specific prompt elicit specialized behavior.

Programmer has the function-calling ability which allows it to "execute" the generated code in an interactive python shell.

AutoGen User Proxy

system_message="An admin that takes input from the user."

termination_criteria code_execution_config max_consecutive_auto_reply llm_config human_input_mode

initiate_chat="Starting message
for that experiment"

Code Execution Env

timeout cache_seed config_list

functions_for_python_cell

functions_for_shell

Group Agent Chat



Planner

Prompt: You interpret scientific queries, devise hypotheses, segment tasks into sequential subtasks with a focus on statistical methods, assign roles to team members, and ensure coordinated progress.



Data Expert

Prompt: You specialize in analyzing statistical data and user queries, offering detailed inferences, formulating and testing hypotheses, and collaborating with programmers for sophisticated data modeling and inferential insights.



Programmer

Prompt: You develop and code tasks assigned by the Planner, utilizing specific function calls and adhering to guidelines to produce outputs in JSON format, with a focus on robust coding and comprehensive logging.

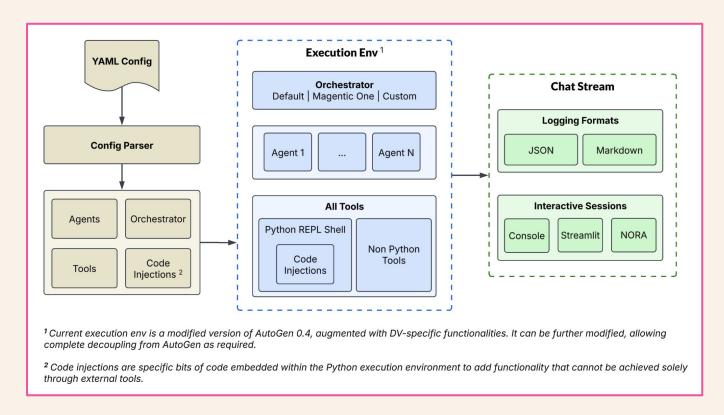


Stats Functions



Prompt: As the Critic, you evaluate and assure the quality of research processes and outcomes, scrutinize research methodologies, assess data quality, and provide constructive feedback on analytical methods and findings.

DataVoyager Core Architecture



LLM Benchmarks

- DataGLUE
- NORA End to End Science
- Bio Neuro/ Immunology

Collabs

Auto-exploration with UMass/ Andrew McCollum group

Better memory for coding/agents with Zurich U

Better Agents/ Models for

- Data Driven Discovery
- End to End Science
- Neuro/Immunology/ Genetics

Better Cost Efficiency

Pareto optimal model usage/ model routing for DDD

100X+ difference b/w cheaper & advanced models

Related Projects List

Augmenting code editing with Al features

- Fork open-source Al code editors like melty, zed or pearai or aider
- Add new features like Al diagraming for architecture etc. in them
- https://aider.chat/docs/leaderboards/

Enterprise heavy requirements

- Add guardrails to prevent certain data leakage
- Serve data while making sure it is not hallucinated
- Can be RAG or database query driven

Other Ideas

- AutoRAG
- RAG + agents on domain specific areas like medicine
- RAG + agents for indian docs
 - Analyze data
 - Make it available in Indian languages

DSBench

https://github.com/LigiangJing/DSBench

- Has files such as .txt, .xlxs
- Has metadata/background about the data
- MCQ
- Requires computation
- Autogen-based GPT-4 achieves 87.84% task success

DSBench benchmark consists of 466 data analysis tasks and 74 data modeling tasks. We evaluate several state-of-the-art LLMs, LVLMs, and agents, and find that our benchmark is challenging for the existing models.



An election has been heldin Excelstan.Excelstan is a small countryincluding 9 Districts. There are 1000 voters. Each voter is assigned a District Code based on wherehey live. The Dist rict Code is a number between 105 194 and determines what District the voter votes for.

1	105 - 114	Alpha	100,739	62	166
_		-	102,052	19	147
2	115 – 124	Beta	102,128	51	160
3	125 - 134	Gamma	116,072	36	135
4	135 - 144	Delta	119,995	19	142
5	145 – 154	Epsibn	122,875	55	124
6	155 – 164	Zeta	128,927	24	139
_	100		146,040	24	151
7	165 - 174	Eta	148,402	69	137
8	175 - 184	Theta	157,647	67	171
9	185 - 194	lota	162,858	50	163

District Code District

Voter ID 212,231 122 222,632 237,420 251,708 256,580 319,000 323,903 437,795



Question

Inputs

Table 1: district code in Excelstan. Excefile: Voter information and their voted district (not show) How many voters are there in the Delta District? A. 113 B. 114 C. 115 D. 116 E. 117



Given the problem description and data files, a DS agent generates executable codes to solve the problem.











data_excel = pd.read_excel('Election_Voting.xlsx', sheet_name="Data" delta_voting_count = 0 for dis code in data excel['District Code']: if int(dis_code) >= 135 and int(dis_code) <= 144: delta voting count += 1 print('Number of voters in Delta district:', delta_voting_count) Number of voters in Delta district: 115



After executing the code, the DS agent finalize the answer based on execution results: There are 115 voters in Delta district, the answer is C.





DSEval

https://github.com/MetaCopilot/dseval

- Has files such as .txt, .xlxs
- Standard datasets like titanic, twitter
- QA with no options
- Requires computation
- GPT-4-based agents reach ~70% pass rate
- Has other evaluations/validators

Query pop2010 pop2023 pop205 2973190 1.24E+09 1.43E+09 1.67E+09 Calculate the population density China 9424703 1.35E+09 1.43E+09 1.31E+09 for each country in 2023 and 2050. 9147420 3.11E+08 3.4E+08 3.75E+08 Result should be a new frame with "Country" as the index and "2023 Density" and "2050 Density" as the columns. Correct Code Interpreter pd.DataFrame({ 'Country': pop['country'], '2023 Density': pop['pop2023'] / pop['landAreaKm'], '2050 Density': pop['pop2050'] / pop['landAreaKm'] }).set_index('Country') **Intact Violation + Wrong Output** CoML pop['2023 Density'] = pop['pop2023'] / pop['landAreaKm'] growth = (pop['pop2023'] / pop['pop2010']) ** (1 / (2023-2010)) - 1 pop['2050 Population'] = pop['pop2023'] * (1+growth)**(2060-2023) pop['2050 Density'] = pop['2050 Population'] / pop['landAreaKm'] pop[['country', '2023 Density', '2050 Density']].set index('country') **Intact Violation + Presentation Error** pop = pop.set index('country') pop['2023 Density'] = pop['pop2023'] / pop['landAreaKm'] pop['2050 Density'] = pop['pop2050'] / pop['landAreaKm'] pop[['2023 Density', '2050 Density']] Crash Jupyter Al dens 2023 = pop.div(pop['landAreaKm'], axis=0)

dens 2050 = pop.div(pop['landAreaKm'], axis=0)*(1+growth)**(2050-2023)

'2023 Density': density_2023,
'2050 Density': density 2050})

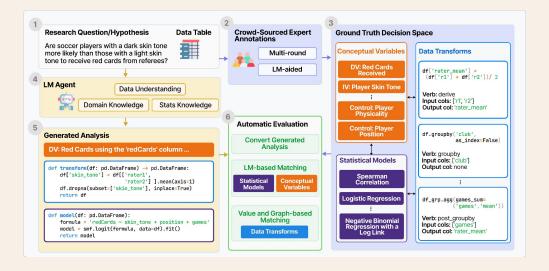
pd.DataFrame({'Country': pop['country'],



BLADE

https://github.com/behavioral-data/BLAD E/tree/main

- Has files, .csvs
- Has metadata about the datasets
- Research Q/Hypothesis as a goal (most in the format: Is this true?)
- MCQs
- Requires computation

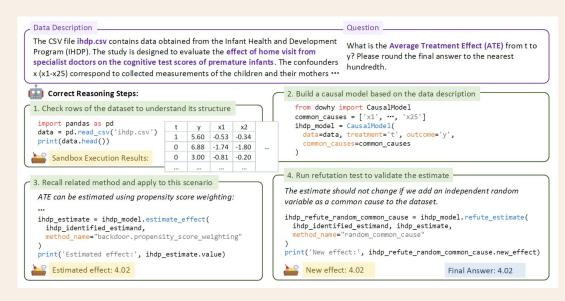




QRData

https://xxxiaol.github.io/QRData/

- Subset of benchmark has files, .csvs
- MCQ
- Requires computation
- GPT-4 achieves ~60% accuracy





Memory with **DataVoyager**

Agent Structure

Group Agent Chat

AutoGen User Proxy

system_message="An admin that takes input from the user."

termination_criteria code_execution_config max_consecutive_auto_reply 11m_config human input mode

initiate_chat="Starting message for that experiment"

Code Execution Env

timeout cache seed config_list

functions_for_python_cell

functions_for_shell

Project Memory



a focus on statistical Library functions, methods, assign roles to User preferences, coordinated progress. Domain knowledge,

Private to project (can be proprietary info)

Explored hypotheses



Planner

Prompt: You interpret scientific queries, devise hypotheses, segment tasks into sequential subtasks with team members, and ensure



Data Expert

Prompt: You specialize in analyzing statistical data and user queries, offering detailed inferences, formulating and testing hypotheses, and collaborating with programmers for sophisticated data modeling and inferential insights.



Programmer

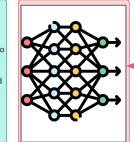
Prompt: You develop and code tasks assigned by the Planner, utilizing specific function calls and adhering to guidelines to produce outputs in JSON format, with a focus on robust coding and comprehensive logging.



Stats Functions



Critic*



Finetuned critic

> Learned from a big corpus/ Voluntarily shared logs



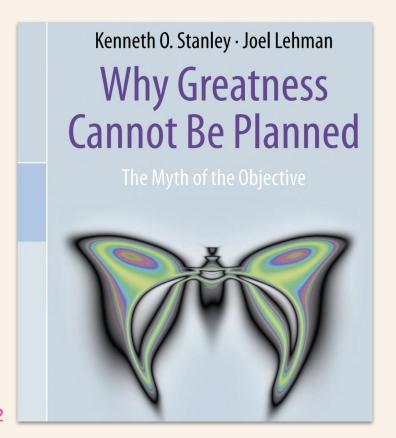
But, what if the goal is unknown?

Can <u>open-ended</u> exploration lead to <u>useful</u> discoveries?

- No explicit research **question** or discovery **goal**
- Only given a dataset
- Optionally, a set of known hypotheses

- Should align with human notions of "interestingness"
 - Discoveries made may answer human-posed research queries in the future
- Serves as a building block for future discoveries

Not just a philosophical point



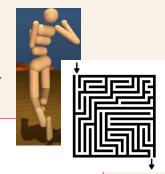
Could address **practical limitations** in **goal-driven discovery**:

- Low dataset coverage: Performs shallow analyses; Struggles to sufficiently cover the dataset; misses promising experiment designs.
- Low diversity: Repeated execution does not produce diverse outputs. <u>Inference-time scaling laws do not</u> <u>seem to hold in our discovery setting.</u>

Open-endedness in data-driven discovery

Benefits:

- **Immediate feedback:** Closed loop of proposing and executing experiments on a given dataset allows for a continual, autonomous process to collect verifiable hypotheses.
- **Expanding search space:** As hypotheses are collected, new opportunities arise to further explore known hypotheses, combine them, or continue sampling new ones.
 - o Solutions quickly plateau in static search spaces. E.g., maze solving, bipedal walking.

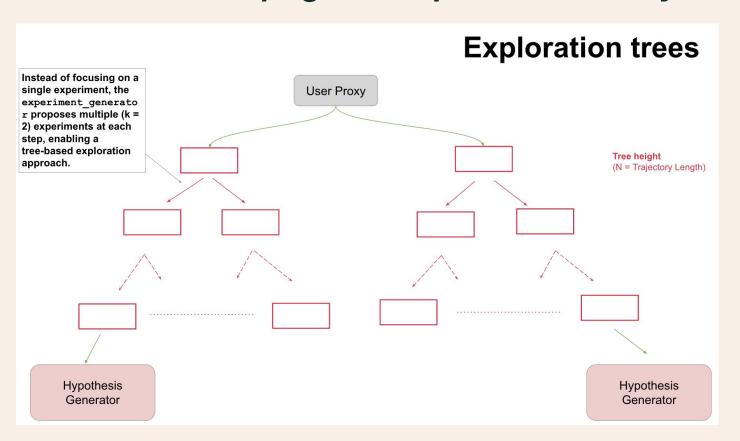


3 main challenges:

- 1. Repeated sampling is insufficient. **Need a method to sample diverse trajectories.**
- 2. Is diversity/novelty a sufficient metric? What reward function should be used to guide exploration?
 - Interestingness? Utility? Hard to define!
- 3. Even with useful reward functions, **how should models be steered to continually generate interesting hypotheses**, especially given the dynamical nature of "interestingness"?
 - Is prompting enough? Can we do more?

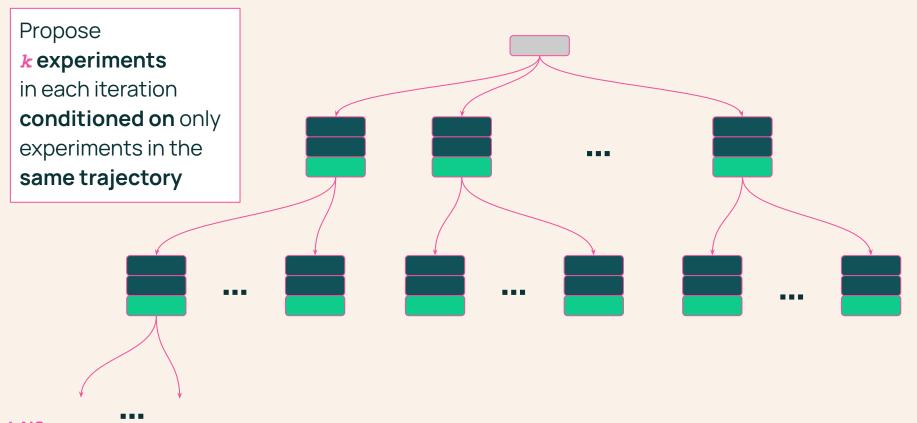


AutoDataVoyager - Exploration Trajectories





AutoDV: Experiment Tree



Experiment Tree: Samples

E: Analyze the influence of family size on educational achievement, specifically focusing on the completed BA degree status.

R: Significant difference in family size between those who completed a BA degree and those who did not, with smaller family sizes associated with degree completion.

H: Individuals from smaller families are more likely to complete a BA degree compared to those from larger families, as evidenced by the significant difference in mean family sizes between the two groups.

E: Investigate potential nonlinear relationships between socioeconomic status (SES) and academic performance indicators such as ASVAB scores and class percentile.

R: Slight nonlinear relationship with SES, but linear relationship remains dominant pattern.

H: There is a slight nonlinear relationship between socioeconomic status (SES) and academic performance indicators, with a quadratic term indicating diminishing returns at higher SES levels. However, the linear relationship remains the dominant pattern, suggesting that SES consistently predicts academic performance across its range.

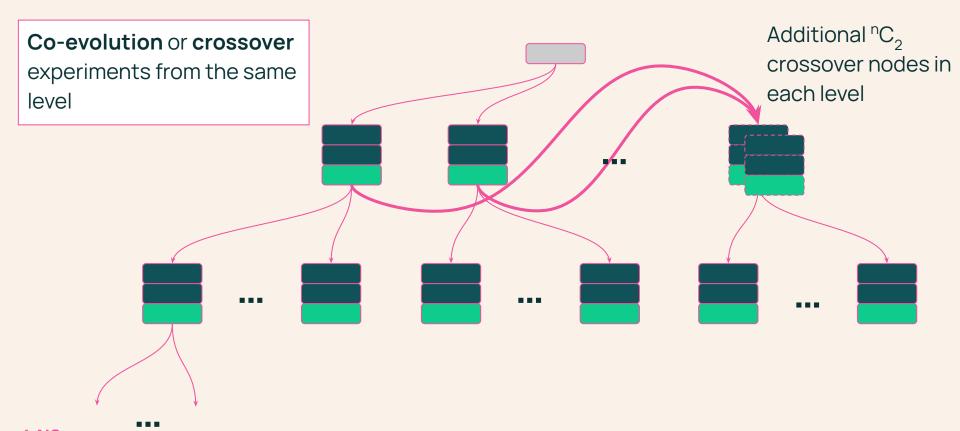
E: Conduct a clustering analysis to identify distinct profiles of students based on their ASVAB scores, class percentile, SES, race, and gender.

R: Three distinct clusters identified, showing varying academic performance and SES levels.

H: Distinct profiles of students can be identified based on ASVAB scores, class percentile, and SES, with higher SES associated with better academic performance. These profiles may also reflect racial and gender disparities, with certain groups more likely to belong to clusters with lower academic performance.

Higher diversity and higher complexity in experiment design.

AutoDV: Crossovers



Ai2

Walkthrough of the DataVoyager system

https://x.com/surana_h/status/178609 7912147239157

https://x.com/mbodhisattwa/status/17 61061506127655244

LLMs have revolutionized NLP.

But they often fall short in key areas.

LLM Limitations (for Advanced Tasks)

- Hallucinations
- Limited context
- Loses memory over longer context
- Non-trivial to verify tasks

- Opaque
- Poor adaptability for out of distribution domain & tasks
- Cannot integrate into specific user workflows

LLM Limitations (for Advanced Tasks)

- Hallucinations
- Limited context
- Loses memory over longer context
- Non-trivial to verify tasks

- Opaque
- Poor adaptability for out of distribution domain & tasks
- Cannot integrate into specific user workflows

Cool demo, but how do I *use* it?

An Example for NLP Code

Non-Agentic Workflow

Write a Python script that uses a Hugging Face model for sentiment analysis on a given bio dataset. Please type out the code in one go without waiting or even using the backspace.



Agentic Workflow

- Begin by refining the problem scope into a detailed specification
- Iterate with domain experts: doctors, clinical data scientists
- Design a modular preprocessing pipeline that can adapt to new domain rules
- Start with a baseline model like BioBERT and iterate
- Define domain-specific metrics that matter including accuracy & F1-score
- Evaluate results on the defined metrics & iterate the code



What are the solutions?

Agentic Design Patterns can tackle the limitations to make more productive use of the LLMs.

- 1. Reflection
- 2. Tool use
- 3. Planning
- 4. Multi-agent collab!

Reflection

Verify & reflect the LLM output by external feedback (i.e. unit tests) & LLMs. Use the reflection to iterate the results.

- Self Refine
- Reflexion

Self-Refine: Iterative Refinement with Self-Feedback. Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh. Peter Clark. NIPS 24.

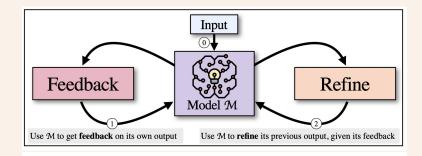


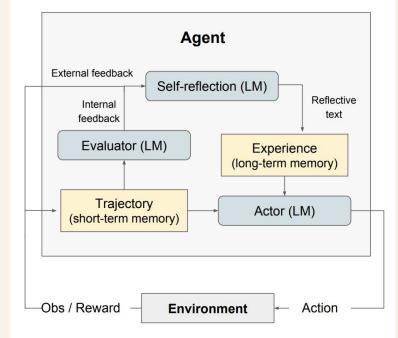










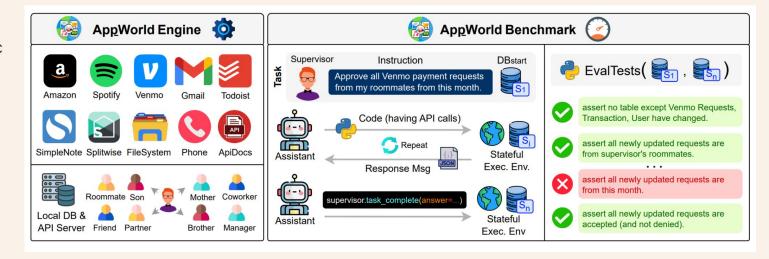


Tool Use

Connect with various tools & functions like:

Browser, APIs, code exec, **custom code** (domain specific code hard for LLMs to generate), search engine etc.

- AppWorld
- Gorilla LLM Exec
- Function calling

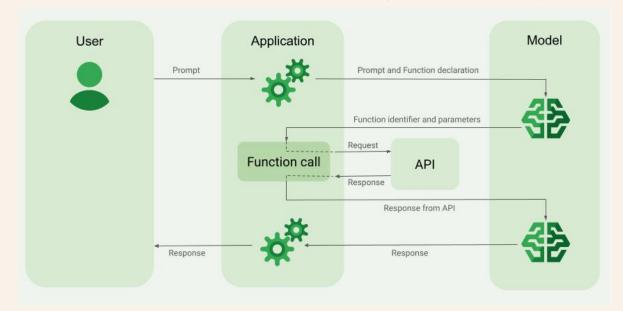


Tool Use

Connect with various tools & functions like:

Browser, APIs, code exec, **custom code** (domain specific code hard for LLMs to generate), search engine etc.

- AppWorld
- Gorilla LLM Exec
- Function calling





^{*} showcased toolcalling is from Google Gemini - but a similar function call works across platforms

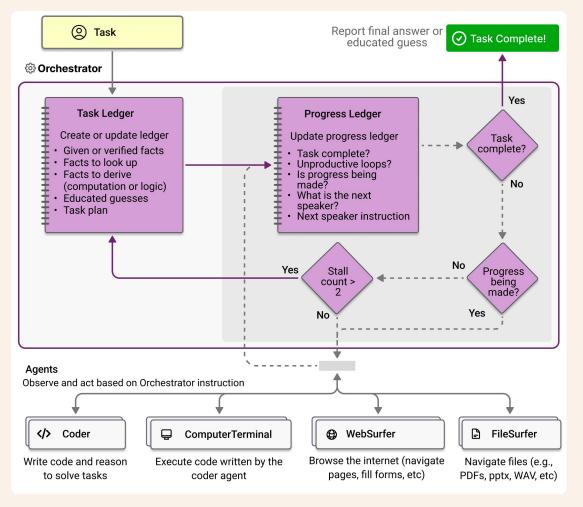
Planning

Plan tasks, track them while performing a task and reason over them if needed.

- Autogen Magento
- MetaGPT
- OpenHands
- HuggingGPT

Magentic-One: A Generalist Multi-Agent System for Solving Complex Tasks. Fourney et. al Microsoft Tech Report 2024.





Multi-agent Collab

Agents work together to solve a complex task.

- Autogen
- CrewAl
- ChatDev
- DataVoyager

ChatDev: Communicative Agents for Software Development. Qian et. al. ACL 2024

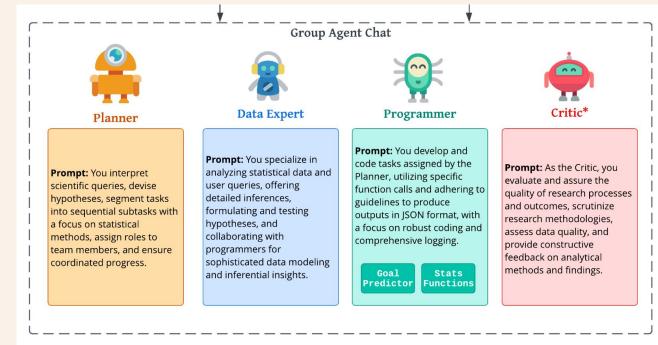




Multi-agent Collab

Agents work together to solve a complex task.

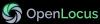
- Autogen
- CrewAl
- ChatDev
- DataVoyager
 - Agents
 - Memory
 - Functions
 - Literature



Data-driven Discovery with Large Generative Models. Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Sanchaita Hazra, Ashish Sabharwal, Peter Clark. ICML 2024.





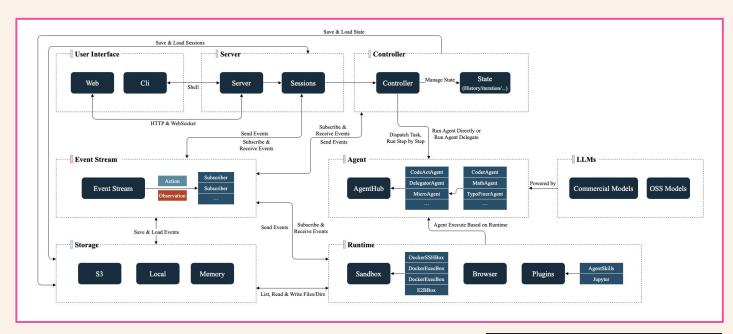


Let's go over 2 concrete examples!

Software Engg with LLMs

(not just code generation)

OpenHands by CMU & UIUC



Full coding env; **not just code gen**

Access to agents

Access to runtime

Access to tools

Pluggable LLMs

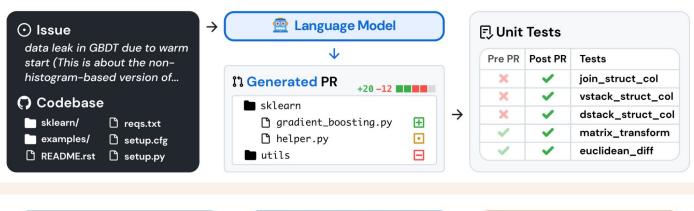
Top LLM repo on Github



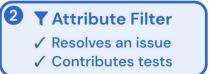




SWE Bench by Princeton



☐ 12 Scrape PRs
☐ 12 popular repositories
☐ >90% Python Code



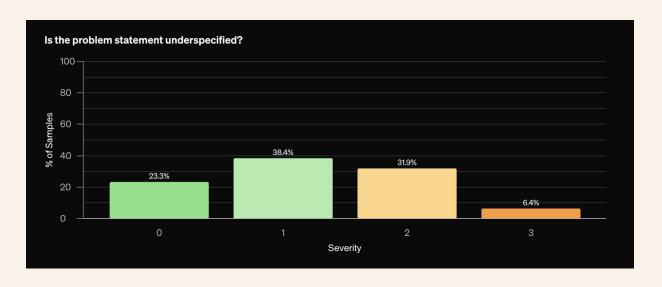


task instances from real-world Python repositories by connecting GitHub issues to merged pull request solutions that resolve related tests.

Provided with the issue text and a codebase snapshot, models generate a patch that is evaluated against real tests.

Can Language Models Resolve Real-World GitHub Issues? Carlos E. Jimenez*, John Yang*, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, Karthik Narasimhan. ICLR 2024

SWE Bench Verified by OpenAl



SWE-bench Verified manually screened by 93 SWEs for 1,699 random samples.

Whether we consider the issue description to be underspecified and hence unfair to be testing on.

Whether the FAIL_TO_PASS unit tests filter out valid solutions.

https://openai.com/index/introducing-swe-bench-verified/

Neil Chowdhury, James Aung, Chan Jun Shern, Oliver Jaffe, Dane Sherburn, Giulio Starace, Evan Mays, Rachel Dias, Marwan Aljubeh, Mia Glaese, Carlos E. Jimenez, John Yang, Kevin Liu, Aleksander Madry



OpenHands - CodeAct Agent

First agent to cross 50% in SWE-Bench Verified

Task planning by developing capabilities for bug detection, codebase management, and optimization

Made a number of fixes to make it easier for agents to **traverse directories**

Switched to use **function calling**, a method used by language models to more precisely specify the functions available to them

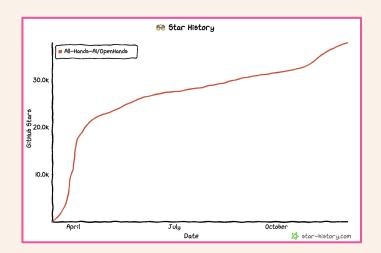
Agent + Model	Score
OpenHands + CodeAct v2.1 + Claude 3.5 Sonnet	53.00
Anthropic Tools + Claude 3.5 Sonnet	49.00
Anthropic Tools + Claude 3.5 Haiku	40.60
Composio SWEkit + Claude 3.5 Sonnet	40.60
SWE-agent + Claude 3.5 Sonnet	33.60
SWE-agent + Claude 3 Opus	18.20
RAG + Claude 3 Opus	7.00
RAG + Claude 2	4.40

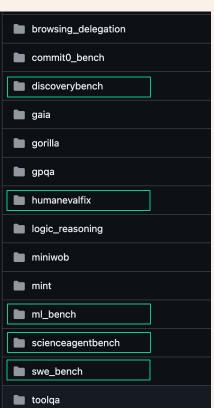
Contributing to OSS LLM Dev with OpenHands

We are contributing for better evals & data science agents - can discuss more offline.

Do try contributing to OpenHands. Their aim is to have full replication of production-grade applications with LLMs

https://github.com/All-Hands-Al/OpenHands/blob/main/CONTRIBUTING.md

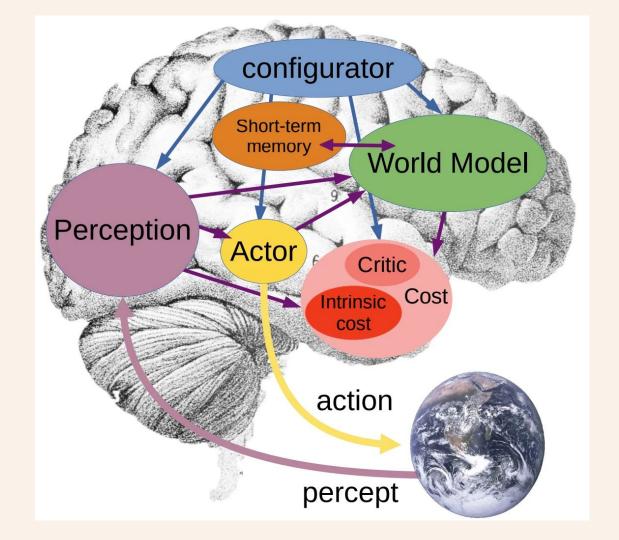






V Jepa

By Meta/ Yann Le Cunn



Thank you!

