
DISCOVERYBENCH: Towards Data-Driven Discovery with Large Language Models

Bodhisattwa Prasad Majumder*¹ Harshit Surana*^{1,2} Dhruv Agarwal*³
 Bhavana Dalvi Mishra*¹ Abhijeetsingh Meena² Aryan Prakhar² Tirth Vora²
 Tushar Khot¹ Ashish Sabharwal¹ Peter Clark¹
¹Allen Institute for AI ²OpenLocus ³University of Massachusetts Amherst

Website: <https://github.com/allenai/discoverybench>
 📄 <https://huggingface.co/datasets/allenai/discoverybench>
 *equal contributions

Abstract

Can the rapid advances in code generation, function calling, and data analysis using large language models (LLMs) help automate the search and verification of hypotheses purely from a set of provided datasets? To evaluate this question, we present DISCOVERYBENCH, the first comprehensive benchmark that formalizes the multi-step process of data-driven discovery. The benchmark is designed to systematically assess current model capabilities in discovery tasks and provide a useful resource for improving them. Our benchmark contains 264 tasks collected across 6 diverse domains, such as sociology and engineering, by manually deriving discovery workflows from published papers to approximate the real-world challenges faced by researchers, where each task is defined by a dataset, its metadata, and a discovery goal in natural language. We additionally provide 903 synthetic tasks to conduct controlled evaluations across task complexity. Furthermore, our structured formalism of data-driven discovery enables a facet-based evaluation that provides useful insights into different failure modes. We evaluate several popular LLM-based reasoning frameworks using both open and closed LLMs as baselines on DISCOVERYBENCH and find that even the best system scores only 25%. Our benchmark, thus, illustrates the challenges in autonomous data-driven discovery and serves as a valuable resource for the community to make progress.

1 Introduction

Knowledge discovery via the scientific process has been a catalyst for human progress for centuries but has, thus far, been a predominantly manual pursuit [16]. Recent breakthroughs in capabilities of large language models (LLMs) to reason and interface with the world using code [9, 40], external tools [41], and interactive agents [51, 32], however, now suggest the possibility of realizing a discovery system that is fully autonomous. Indeed, recent works [33] provide initial evidence for this paradigm within the setting of *data-driven discovery*, where both search and verification of hypotheses may be carried out using a dataset alone (i.e., after physical experiments and data collection¹), but the extent of this ability remains unclear. We, therefore, aim to systematically evaluate the following question:

How good are current state-of-the-art LLMs at automated data-driven discovery?

¹In practice, experiments and analysis are interleaved, not sequential. Our concern in this work, however, is systematically studying the data analysis part of the (interleaved) pipeline.

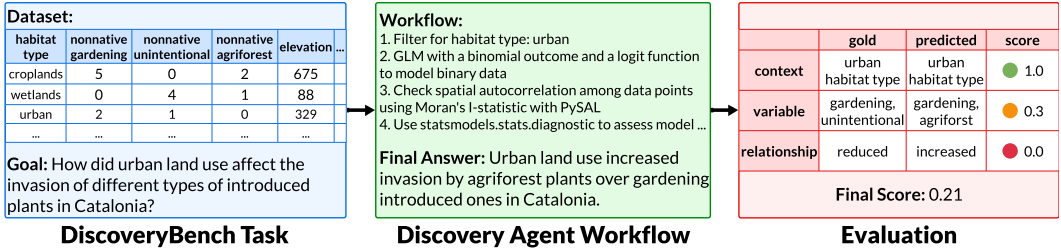


Figure 1: Each DISCOVERYBENCH task consists of a goal and dataset(s) (left). Solving the task requires both statistical analysis and scientific semantic reasoning, e.g., deciding which analysis is appropriate for the domain, and mapping goal terms to column names (center). A faceted evaluation allows open-ended final answers to be rigorously evaluated (right).

Answering this question is hard, as data-driven discovery in the wild (real-world) is diverse across domains and subject areas, which in turn makes it difficult to build a robust evaluation framework to measure progress. We address this using a pragmatic formalization of data-driven discovery, namely the search for a *relationship* that may hold between *variables* in a *context*, where (importantly) the description of those facets may not be in the language of the dataset. A data-driven discovery task then has one of these components missing, e.g., “How did urban land use affect the invasion of introduced plants in Catalonia?”. Importantly, this formalization allows for systematic, reproducible evaluation over a wide variety of real-world problems, by leveraging these facets (Fig 1, right).

Unlike prior datasets for statistical analysis [28] or AutoML [55, 15], DISCOVERYBENCH tasks also require scientific semantic reasoning, for instance, deciding which of the many possible analysis techniques are appropriate for the domain (e.g., spatial autocorrelation for plant invasion, Fig 1 center), how to clean and/or normalize the data, and how to map goal terms to dataset terms (e.g., “land use” to “habitat type”). Task solutions typically requires a multistep workflow (Fig 1, center). In this way, DISCOVERYBENCH is the first large-scale dataset to address the broader data-driven discovery pipeline, not just the statistical analysis component, and explore LLMs’ capacity for this.

Given this framework, we created DISCOVERYBENCH by manually extracting 264 discovery tasks, i.e., goal + dataset(s), from over 20 published papers, as well as creating real-world discovery workflows that solve each task. We additionally provide 903 synthetic tasks across 48 domains generated using LLMs to mimic the real-world discovery process. The synthetic benchmark allows us to conduct controlled model evaluations by varying task difficulty. Our contributions are thus:

- DISCOVERYBENCH, the first comprehensive benchmark to formalize the multi-step process of data-driven hypothesis search and verification, covering many real-world discovery tasks plus additional synthetic tasks.
- A pragmatic formalism for data-driven discovery, flexible enough to characterize many real-world tasks while constrained enough to allow for rigorous, reproducible evaluation.
- A comprehensive evaluation across state-of-the-art LLM-based reasoning methods (“discovery agents”). We find performance peaks at 25%, demonstrating the challenging nature of our task.

These suggest that DISCOVERYBENCH may be a valuable resource for helping make progress on autonomous, data-driven discovery.

2 Related Work

Automated data-driven discovery has been a long-standing dream of AI [33, 21]. Although there have been a range of **data-driven discovery systems**, from early ones that fit equations to idealized data, e.g., Bacon [23], to more modern ones handling complex real-world problems, e.g., AlphaFold [19], their associated datasets are task-specific and customized to a pre-built pipeline. In contrast, DISCOVERYBENCH aims to be a general test over multiple tasks, including testing whether systems can design appropriate pipelines themselves.

A number of datasets and tools are available for **AutoML**, a related technology aimed at automating workflows for building optimal machine learning models [18, 55, 24]. AutoML tools include packages like Scikit [13], and embedded in cloud platforms such as Google Cloud Platform, Microsoft Azure,

and Amazon Web Services. However, associated datasets for AutoML are primarily used for training models, rather than for open-ended discovery tasks.

Similarly, there are several datasets that test **statistical analysis** in various fields, e.g., [42, 25, 50]. Software packages like Tableau, SAS, and R also support users in that task. However, these datasets and tools are designed specifically for data analysis, while DISCOVERYBENCH aims to automate the broader pipeline including ideation, semantic reasoning, and pipeline design, where statistical analysis is just one component.

One recent dataset similar in spirit to ours is QRData [28]. QRData also explores LLM capabilities but targets statistical/causal analysis for well-defined (mainly) textbook questions that have unique, (mainly) numeric gold answers. In contrast, DISCOVERYBENCH has no prescribed boundaries on statistical techniques to apply, uses open-ended questions and answers, and complex tasks drawn from state-of-the-art published work.

3 Formalization

We begin by formalizing what we mean by a data-driven hypothesis and how the structure of a complex hypothesis may be viewed as a hypothesis semantic tree.

A **data-driven hypothesis** h in \mathcal{H} (the space of such hypotheses) is a declarative sentence about the state of the world whose truth value may be inferred from a given dataset D using a verification procedure $\mathcal{V}_D : \mathcal{H} \rightarrow \{\text{supported, unsupported}\}$, for instance, via statistical modeling.

Each hypothesis may further be expressed using a propositional formula ϕ over a set of **sub-hypotheses** $h_i \in \mathcal{H}$ using logical connectives, e.g., disjunctions and conjunctions, such that $h := \phi(h_1, \dots, h_n)$ and $\mathcal{V}_D(h) = \phi(\mathcal{V}_D(h_1), \dots, \mathcal{V}_D(h_n))$. For instance, suppose h is the hypothesis “for men younger than 20, popularity of product A varies proportional to their age (h_1), while there exists an inverse relationship for those older than 40 (h_2)”, then h can be expressed as the conjunction $h_1 \wedge h_2$.

Inspired by recent work of Thompson and Skau [47], we additionally introduce a structured formalism that breaks a hypothesis down into **three hypothesis dimensions**:

- **Contexts** (c): Boundary conditions that limit the scope of a hypothesis. E.g., “for men over the age of 30” or “in Asia and Europe” or unbounded/full dataset when not specified.
- **Variables** (v): Known set of concepts that interact in a meaningful way under a given context to produce the hypothesis. E.g., gender, age, or income. Note that each hypothesis is associated with a target variable and a set of independent variables.
- **Relationships** (r): Interactions between a given set of variables under a given context that produces the hypothesis. E.g., “quadratic relationship”, “inversely proportional”, or piecewise conditionals.

With slight abuse of notation, we can now equivalently define hypothesis $h := \psi(c, v, r)$, where $\psi(\cdot, \cdot, \cdot)$ returns the declarative sentence “under context c , variables v have relationship r .” For instance, for sub-hypothesis h_1 in our example above, $c_1 :=$ “men younger than 20”, $v_1 := \{\text{gender, consumer_age, product_popularity}\}$, and $r_1 :=$ “popularity is proportional to age”.

Hypothesis Semantic Tree. Observe that each independent variable in a hypothesis may itself be a target variable for a prior hypothesis. To emphasize this hierarchical nature, we introduce the concept of a *hypothesis semantic tree* whose nodes are variables (independent or derived) and whose sub-trees represent hypotheses, as follows. Consider a hypothesis h . A semantic hypothesis tree \mathcal{T}_h with h as the *primary* hypothesis is a Markov tree whose root node is the target variable of h , each of whose leaf nodes is an independent variable that is not derived further, and each of whose internal nodes is the target variable of an *intermediate* hypothesis. In other words, each sub-tree $\mathcal{T}_{h'}$ rooted at an internal node v of \mathcal{T}_h is itself a hypothesis semantic tree for a hypothesis h' with v as the target variable. In particular, a sub-tree rooted

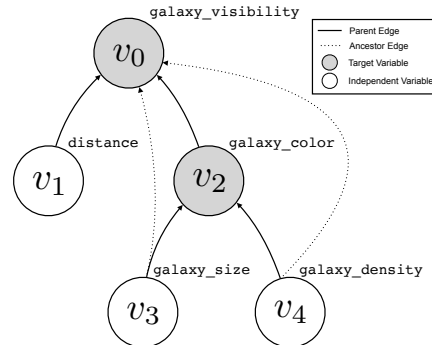


Figure 2: Hypothesis Semantic Tree

at v and all its immediate children C_v implicitly encodes h' as $\psi(c, \{v\} \cup C_v, r)$ where r is the relationship between v and C_v under context c as specified in h' . More generally, \mathcal{T}_h can encode many different hypotheses by choosing one node as the target variable and considering nodes at arbitrary descendant levels as independent variables.

For instance, in Fig 2, we show a semantic tree \mathcal{T}_h with the following primary hypothesis (h): “*The visibility of a galaxy reduces when the blue spectrum dominates and the distance of the galaxy from Earth increases*”, where `galaxy_visibility` (v_0) is the target variable with independent variables `distance` (v_1) and `galaxy_color` (v_2). Consider now the sub-tree rooted at v_2 , which encodes the following intermediate hypothesis: “*Visibility of blue light from galaxies increases with an increase in galaxy size and decrease in star density*”, where `galaxy_color` (v_2) is the target variable with independent variables `galaxy_size` (v_3) and `galaxy_density` (v_4). Note further that due to there existing ancestor edges from v_3 and v_4 to v_0 , \mathcal{T}_h also encodes the hypothesis: “*The visibility of a galaxy reduces with distance from Earth combined with an increase in galaxy size and decrease in star density*”, where the target variable is v_0 and the independent variables are v_3 and v_4 .

(Task) Dataset. We formally define a dataset D on which hypothesis search and verification is performed as a collection of tuples $\{\mathbf{x}_i\}_{i=1}^m$ that supports multiple hypothesis semantic trees resulting in a *semantic forest* $\mathcal{F} := \cup_i \mathcal{T}_{h(i)}$, where each \mathbf{x}_i is a row in the dataset and $x \in \mathbf{x}_i$ is an observation for a particular column. Further, \mathbf{x}_i may span only a subset of nodes in \mathcal{F} , i.e., not all nodes in \mathcal{F} may be observed. Specifically, while roots and leaves are always observed, internal nodes (target variables for intermediate hypotheses) may be latent. Therefore, multiple versions of D may be collected for \mathcal{F} with different degrees of observability of internal nodes, altering the difficulty of the discovery task.

4 DISCOVERYBENCH

We now introduce a novel benchmark, DISCOVERYBENCH, for discovering data-driven hypotheses. In this benchmark, a *data-driven discovery task* is defined as follows: Given one or more task dataset(s) D and a discovery goal G , derive a hypothesis $h = \psi(c, v, r)$ addressing G with the highest specificity for the context c , variables v , and relationship r supported by D . Optionally, a workflow of deriving such a hypothesis can be outputted to augment information already present in the hypothesis. DISCOVERYBENCH has two components: **DB-REAL** encompassing data-driven hypotheses and workflows derived from published scientific papers and **DB-SYNTH** capturing systemic variations in data-driven hypotheses and workflows obtained from synthetically generated datasets. We release our dataset under the ODC-BY license: <https://github.com/allenai/discoverybench>.

4.1 DB-REAL: Collecting data-driven hypotheses in the wild

Our goal is to replicate the scientific process undertaken by researchers to search for and validate a hypothesis from one or more datasets. We focus on six scientific domains where data-driven research is the cornerstone of scientific progress: sociology, biology, humanities, economics, engineering, and meta-science. Our data collection follows either a **data-first** or **code-first** approach.

For the **data-first approach**: 1) we filter papers based on open public datasets (D) such as National Longitudinal Surveys (NLS), Global Biodiversity Information Facility (GBIF), and World Bank Open Data (WBOD) that have workflow details; 2) we then try to replicate these workflows in Python. For this data-first approach, replication took up to 90 person-hours per dataset, often (30%) not resulting in success. This highlights building data-driven discovery benchmarks from real studies is not only challenging and time-consuming, but automating discovery can also be key for scientific progress and reproducibility.

The data-first approach by design is limited to well-known aforementioned public datasets. To improve diversity in domains, datasets (D), and workflows, we also adopted a **code-first approach** to look beyond popular public datasets. In this approach, we 1) search for code repositories based on scientific papers with available datasets and 2) attempt to replicate them in Python with existing code or from scratch with interpretation of the associated paper. We looked at 785 data points in Zenodo, EU’s Open Research Repository, with a filter for computational notebooks. Over 85% of the repositories either had missing code, code that could not be easily translated to Python, or a

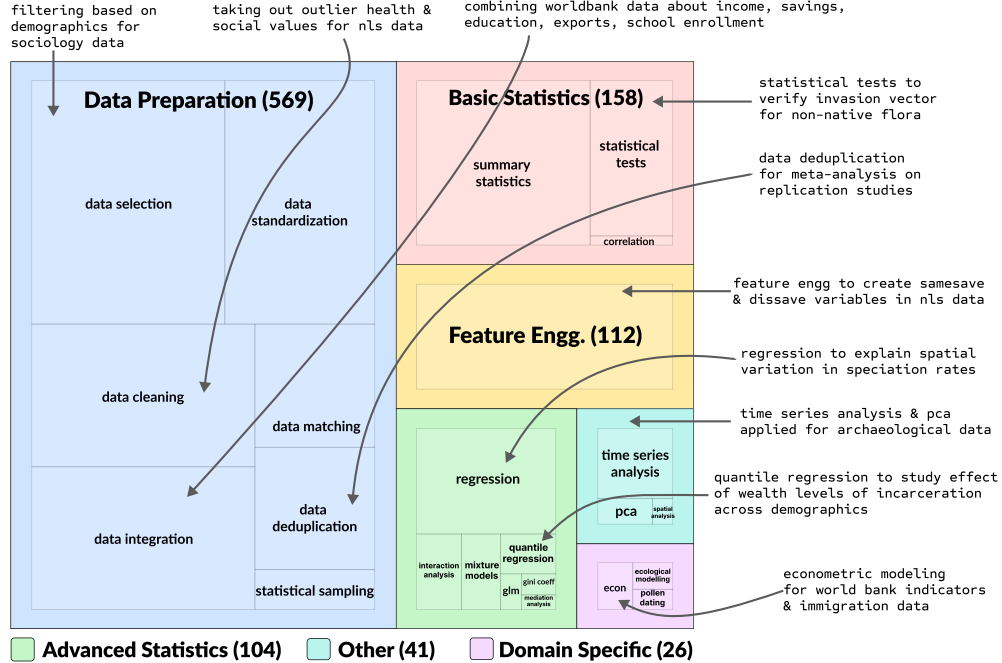


Figure 3: Workflow categories in DB-REAL with representative examples.

proprietary/non-open dataset. We finalized a candidate list of 14 repositories, but in the end, 3 of them passed all our checks for their hypotheses to be included in the benchmark².

Upon replication of the result or implementation of the full procedure as described in the paper, we include the (dataset D , hypothesis h , implementation workflow) tuple to the benchmark.

During the process, the implementation workflow may lead to other hypotheses that are not directly reported in the paper but can be supported by the data. We included them in DISCOVERYBENCH, which leads to a good mix of already reported science-worthy hypotheses as well as novel hypotheses grounded in datasets. This is particularly useful as our goal is to evaluate LLMs’ ability to solve a discovery task that is realistic but never reported before.

Finally, the task datasets are supplemented with a dataset description, natural language descriptions of the columns, and additional background knowledge related to the domain or the datasets. Some of our tasks, for instance, archaeology, require domain knowledge to derive a particular hypothesis.

Inferring task difficulty. Using the hypothesis semantic tree defined in Section 3, we say that the difficulty of a discovery task is proportional to the *path length* from an observed node to the target hypothesis node in the tree. However, knowing the tree structure from a task dataset alone is impractical due to incomplete *a priori* information about unobserved intermediate nodes and edges between observed nodes. To infer task difficulty, we, therefore, approximate the path length between the target and leaf nodes using the length of the implementation workflow required to derive a target hypothesis. Specifically, for each step in the workflow, we add 1 to the discovery path length. In some cases, we derive two tasks: easy and hard from the same hypothesis, where for easy, we provide the derived variables as observed variables in the dataset (e.g., BMI), and for hard, it would require deriving intermediate variables (BMI from height and weight) to reach the target. Additionally, given the view of a task dataset as encoding the union of multiple semantic trees rooted at different hypotheses, i.e., a semantic forest \mathcal{F} , we further posit that task difficulty increases as the number of trees in the forest ($|\mathcal{F}|$) increases. Intuitively, discovery becomes harder as the hypothesis search space increases. In practice, this setting is observed when a task requires access to multiple datasets.

Forming discovery goals. By definition, each hypothesis can be fully specified by the declarative sentence as $h := \psi(c, v, r)$. To systematically construct the discovery goals for the task, we first mask one of each dimension, context c , variable v , relationship r , and generate a discovery goal to

²Some repositories include hypotheses from multiple papers as their background.

identify the masked information given the rest of the hypothesis and the task dataset(s). For instance, for a target hypothesis, “*The effect of socioeconomic status on college degree completion is higher for females (0.4995) than males (0.4467)*”, we form a discovery goal as “*How does socioeconomic status impact on college degree completion for females compared to males?*” seeking the relationship r to be discovered from the dataset(s) given the relevant variables v and context c . Additionally, we ensure each discovery goal leads to only one answer, i.e., the target hypothesis.

4.1.1 Features of DB-REAL benchmark

DISCOVERYBENCH incorporates a broad landscape of data-driven discovery. With over 500 instances of data preparation activities such as cleaning, deduplication, and integration, captures the complexity of real-world data pre-processing for discovery. Tasks also demand a spectrum of statistical methods, from statistical tests to mixture models, and include domain-specific approaches in econometric and ecological modeling, as reflected in the Fig 3³.

	Train	Test
# tasks	25	239
# unique hypotheses	14	144
# tasks need > 1 dataset	4	110
# domains	3	6

Table 1: Statistics for DB-REAL

Table 1 shows the diversity of tasks both in train and test split for DB-REAL. Most importantly, the benchmark incorporates 114 (4 + 110) tasks that require more than one related datasets to be analyzed, with a maximum of 6 datasets for a task. Each workflow within the dataset can be viewed as a composition of unit actions—such as code generation for statistical tests—that LLMs excel at, showing how our tasks require the chaining of such atomic actions to address complex scenarios for data-driven discovery. We measure the complexity of these workflows by quantifying the number of unit actions involved, referring to this as the *workflow length*, whose distribution can be seen in Fig 5.

4.2 DB-SYNTH: Generating data-driven hypotheses using LLMs

To scale data collection, we next introduce a supplementary benchmark, which is synthetically constructed to enable controlled model evaluations. Our goal is to reverse-engineer the process of hypothesis discovery to synthesize datasets and discovery tasks of varying difficulty that require analysis workflows similar to those in the real-world benchmark. Our approach leverages the broad pre-trained knowledge of LLMs in four stages:

Domain sampling: First, we prompt the model to generate a list of diverse topics or *domains* along with their natural language descriptions. E.g., “Ancient architecture” → “Related to historic buildings, architectural marvels, and ancient construction techniques”.

Semantic tree construction: For each domain, we then build a semantic tree \mathcal{T}_h , recursively deriving nodes starting from a primary hypothesis h . Specifically, we prompt the model with the domain and a sampled real-world workflow (e.g., “within-cluster analysis”) to generate a hypothesis and its target variable. Setting the target variable as root, we then derive child nodes by generating the independent variables required to verify h using $\mathcal{V}(\cdot)$. We operationalize this by generating a column name and description for each child node (along with a data type and range) and a pandas expression⁴ [49] over only independent variables in \mathcal{T}_h such that its execution results in the target variable. We repeat this with each leaf in \mathcal{T}_h as the root of a new semantic sub-tree, generating intermediate hypotheses and a new set of variables until the desired height of \mathcal{T} is reached.⁵ We also generate a set of distractor columns disjoint from nodes in \mathcal{T} , thus resulting in a synthetic semantic forest \mathcal{F} .

Data generation: We then construct a task dataset $D := \{\mathbf{x}_i\}_{i=1}^m$ by generating synthetic data in a bottom-up manner (i.e., from leaves to root) for each node in \mathcal{F} . Starting with various sampling strategies for leaf nodes (see more in Sec D), for each subsequent level in \mathcal{F} , we create new columns for nodes by simply executing their pandas expressions. Finally, to mimic real-world challenges in data collection, we probabilistically perturb each instance $x \in \mathbf{x}_i$ by adding noise or dropping values to create missing data⁶. Note that at this stage, D contains a column for each node in \mathcal{F} .

³A task may require multiple data preparation and analytical activities

⁴The pandas expression encodes the structured hypothesis $\psi(c, v, r)$.

⁵In practice, with probability 0.6, we choose whether a node is derived further or marked as a leaf.

⁶Each value is noised independently; therefore, each row has sufficient true data useful for discovery.

Task generation: For each internal node h in \mathcal{F} , we now create multiple task datasets $D_h^{(l)}$ from D , varying the difficulty of the discovery task based on the path length l between h and the observed independent variables in \mathcal{F} . Finally, we follow the same strategy for goal formulation as DB-REAL. We generate 903 tasks over 48 diverse domains and assign them to train, dev, and test sets using a 60/20/20 split, where each task is additionally tagged with a difficulty level from 1-4. While we evaluate our agents on the test, the training set can serve as supervised data for improving models.

4.3 Evaluation

We evaluate task performance by measuring the alignment of the predicted and gold hypotheses in natural language.⁷ We take inspiration from recent works in LLM benchmarking [43, 54, 52, 14, 26, 27] and design a model-based evaluation strategy using `gpt-4-preview-0125` as the *evaluator*, conditioned on our structured formalism of data-driven hypotheses.

Recall the propositional form $h := \phi(h_1, \dots, h_n)$ of a hypothesis h that decomposes it into sub-hypotheses. We first use our GPT-4 based evaluator to independently decompose the gold (h^g) and predicted (h^p) hypotheses into their respective sub-hypotheses $\{h_i^g\}_{i=1}^n$ and $\{h_j^p\}_{j=1}^m$, asking it to also identify, for each sub-hypothesis h_k , its context, variables, and relationship dimensions (prompt in Listing 1). Given this structured representation of the gold and predicted hypotheses, we then compute a **hypothesis match score (HMS)**, which measures the degree to which two hypotheses align on each dimension, as follows.

To compute HMS, we match each predicted sub-hypothesis h_j^p with a gold sub-hypothesis h_i^g when their contexts are judged as equivalent by our GPT-4 based evaluator (prompt in Listing 2).⁸ Let M denote this set of context-matched pairs of predicted and gold sub-hypotheses. At this point, treating each sub-hypothesis context as a single unit, we can compute an F1 score, ctx_{F1} , capturing how aligned the n contexts of sub-hypothesis of h^g with the m contexts of sub-hypotheses of h^p . Then, for each matched pair of sub-hypotheses, we measure how well the variables and relations align, using an F1 score for the variables (var_{F1}) and an accuracy score for the relation (rel_{acc}). Specifically, for each sub-hypothesis pair in M , we extract the set of interacting variables in the gold and predicted sub-hypotheses using the GPT-4 based evaluator (prompt in Listing 3). We compute the alignment between these two sets of variables as an F1 score, var_{F1} , similar to how ctx_{F1} was computed. For relationships, we compute relationship accuracy with reference to the relationship between the gold variables (rel_{acc}) based on evaluator judgments using the following scoring heuristic: 100 if there is an exact match of the relation, 50 when the predicted relationship is broader than the gold relationship but encompasses it, and 0 otherwise (prompt in Listing 4). Finally, we compute $\text{HMS} \in [0, 100]$ as the average alignment of the variable and relationship dimensions over context-matched sub-hypotheses, weighted by the overall context alignment:

$$\text{HMS}(h^p, h^g) = \text{ctx}_{\text{F1}}(h^p, h^g) \times \frac{1}{|M|} \sum_{i=1}^{|M|} \left(\text{var}_{\text{F1}}(h_i^p, h_i^g) \times \text{rel}_{\text{acc}}(h_i^p, h_i^g) \right)$$

5 Experiments

5.1 Discovery Agents

We benchmark state-of-the-art LLM-based few-shot reasoning methods as discovery agents with two closed models, GPT-4o and GPT-4-0125-preview (GPT-4p), and one open, Llama-3-70B, model powering the reasoning methods. A discovery agent takes the task description, paths to the task dataset(s) D , metadata about the datasets (description, column descriptions), and the goal, G , to produce a natural language (NL) hypothesis specified by context, variables, and relationship.

- **CodeGen** generates the entire code at one go to solve the task, where we provide a demonstration of a solution code in the context. After code execution and based on the result, it generates the NL hypothesis and summarizes the workflow.
- **ReAct** [51] solves the task by generating thought and subsequent codes in a multi-turn fashion.

⁷We deliberately take an outcome-based approach as > 1 discovery path may lead to the same hypothesis.

⁸Note that at most one gold sub-hypothesis is matched with a predicted sub-hypothesis.

	GPT-4o	GPT-4p	Llama-3
DB-REAL			
NoDataGuess	0.0	4.7	11.5
CodeGen	15.5	16.3	12.1
React	15.4	15.6	13.5
DataVoyager	15.4	13.9	11.5
Reflexion (Oracle)	24.5	19.5	22.5
DB-SYNTH			
CodeGen	14.1	8.7	10.9
React	11.6	7.4	12.0
DataVoyager	5.7	6.9	11.7
Reflexion (Oracle)	15.7	12.9	23.2

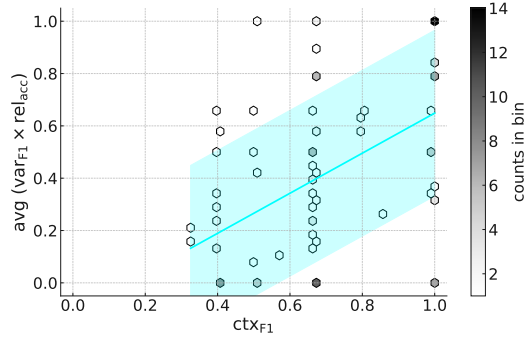


Figure 4: (Left) Hypothesis Matching Scores (HMS) across agent-LLM pairs in DB-REAL and DB-SYNTH. (Right) Scatter plot for ctx_{F1} and average $var_{F1} \times rel_{acc}$, showing accurate contexts increases the probability of predicting variables and relations accurately. Scores are for the best model on DB-REAL and only include data points (44.2%) where both scores are non-zero.

- **DataVoyager** is a multi-component data-driven discovery agent from [33]. It has four components, planner, code generator, data analysis, and critic, that orchestrate the discovery process.
- **Reflexion (Oracle)** [44] is an extension of CodeGen agent, where at the end of one trial, we provide the “oracle” HMS score as an evaluation signal, and it generates a reflection to improve (when $HMS < 1$) in the next trial till it solves the task, or maximum trials (3) are reached.
- **NoDataGuess** guesses the hypothesis (in DB-REAL) just from the dataset description and the goal without accessing the datasets where we measure LLM’s memorization of already published works.

5.2 Main Results

Fig 4(left) shows that overall performance for all framework-LLM pairs is low for both DB-REAL and DB-SYNTH, highlighting the challenging nature of the task and the benchmark. Most importantly, effective reasoning prompts such as React and planning with a self-critic (DataVoyager) do not help improve the simple CodeGen agent. But with oracle feedback, Reflexion (Oracle) significantly improves over CodeGen (base) performance. Analysis reveals that almost all non-reflexion agents solve the *easiest* (in terms of workflow category and length) instances from the benchmark. GPT-4o refuses to hallucinate in the NoDataGuess baseline, whereas surprisingly Llama-3 performs similarly in both data and no-data modes. We additionally observe that the models’ performance in DB-REAL and DB-SYNTH are similar, indicating our synthetic benchmark captures complexities of the real workflow but provides a systematic way to analyze the models’ performance.

5.3 Analysis

Context is important. Fig 4(right) shows the trends of the ctx_{F1} and combined $var_{F1} \times rel_{acc}$. A positive trend signifies that to predict variables and relationships accurately, precise and accurate context prediction is necessary. However, correct identification of context is an important first step, although it does not guarantee success.

Workflow complexity barrier. Almost all agents struggle more with tasks involving complex statistical techniques, complex data preparation methods, or domain-specific models. The top three workflow categories where the best non-oracle model was highly performant are correlation analysis (55%), data selection (18%), and summary statistics (18%), whereas the lowest three workflow categories are spatial analysis (0%), pollen dating (0%), and ecological modeling (0%).

Domain knowledge dependency. To check if additional domain helps agents perform better, we collect targeted domain knowledge for the archaeology-related tasks that needed significant domain knowledge during data collection. When added as additional hints, we find that DataVoyager’s (GPT-4p) performance jumps from 9.9% (w/ domain knowledge) to 17.5% (w/o domain knowledge).

Performance across domains and goal types. Fig 5(a) depicts that biology (0%) and engineering (7%) perform the worst due to their higher dependence on advanced statistical methods, while economics (25%) and sociology (23%) perform better. Additionally, Fig 5(b) shows goals related to discovering a relationship given context and variables are more easily solved than the other two types

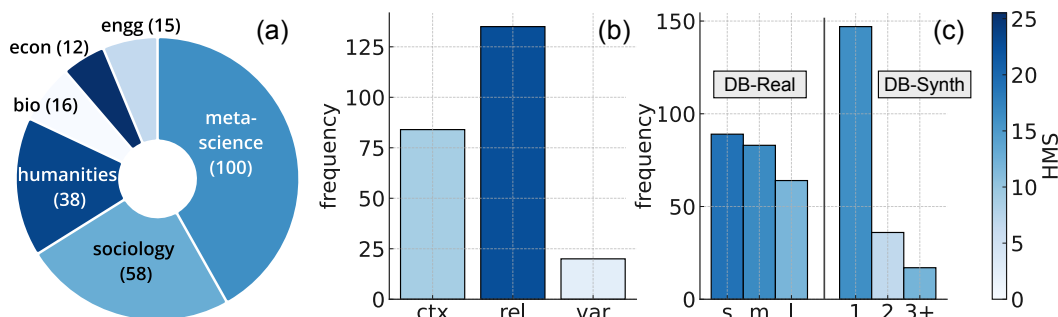


Figure 5: Best non-oracle agent’s performance (HMS) (a) across domains, (b) for goal types (dimension to be discovered), and (c) for different workflow lengths. In (c) workflow length categories for DB-REAL are: s: < 10, m: > 10, < 20, l: > 20. For DB-SYNTH, it is the semantic tree height.

of goals, finding context and variables. This is explained by the complexity of the hypothesis search, which is broader for finding the right context or a set of variables given a fixed relationship, whereas finding the relationship given context and variables is easier.

Impact of workflow length. Inherently, the difficulty of the tasks is measured by the gold workflow length (DB-REAL) or the height of the semantic tree (DB-SYNTH). Figure 5(c) shows a decreasing trend in performance as workflow length (hence, complexity) increases. The performance drops significantly even for medium-length workflows, highlighting current agents’ limitations.

6 Conclusion

We present DISCOVERYBENCH, the first data-driven discovery benchmark consisting of 264 discovery tasks that capture real scientific workflows extracted from published works. We supplement this with 903 structurally generated synthetic tasks, tailored to evaluate discovery agents at various levels of difficulty. We benchmark state-of-the-art reasoning frameworks with the most advanced LLMs, but the best agent’s performance only peaks at 25% underscoring the challenging nature of the task and the benchmark. We hope our timely contribution can increase interest and efforts in making progress on reliable and reproducible autonomous scientific discovery using large generative models.

References

- [1] K. L. Alexander, C. Riordan, J. Fennessey, and A. M. Pallas. Social background, academic resources, and college graduation: Recent evidence from the national longitudinal survey. *American Journal of Education*, 90(4):315–333, 1982.
- [2] A. P. S. Alves, M. Kalinowski, G. Giray, D. Mendez, N. Lavesson, K. Azevedo, H. Villamizar, T. Escovedo, H. Lopes, S. Biffel, et al. Status quo and problems of requirements engineering for machine learning: Results from an international survey. In *International Conference on Product-Focused Software Process Improvement*, pages 159–174. Springer, 2023.
- [3] R. Apel and G. Sweeten. The impact of incarceration on employment during the transition to adulthood. *Social problems*, 57(3):448–479, 2010.
- [4] E. N. Appiah. The effect of education expenditure on per capita gdp in developing countries. *International Journal of Economics and Finance*, 9(10):136–144, 2017.
- [5] J. Brinkmann. Copper output, demand for wood and energy expenditure—evaluating economic aspects of bronze age metallurgy. *How’s life*, pages 11–34, 2019.
- [6] J. P. Brozio, J. Müller, M. Furholt, W. Kirleis, S. Dreibrodt, I. Feeser, W. Dörfler, M. Weinelt, H. Raese, and A. Bock. Monuments and economies: What drove their variability in the middle-holocene neolithic? *The Holocene*, 29(10):1558–1571, 2019.
- [7] J. P. Brozio, J. Kneisel, S. Schaefer-Di Maida, J. Laabs, I. Feeser, A. Ribeiro, and S. Schultrich. Patterns of socio-economic cultural transformations in neolithic and bronze age societies in the central northern european plain: Human-environmental interaction concerning bourdieu’s forms

- of capital. In *Perspectives on Socio-environmental Transformations in Ancient Europe*, pages 105–142. Springer, 2024.
- [8] F. O. Cerezer, C. S. Dambros, M. T. Coelho, F. A. Cassemiro, E. Barreto, J. S. Albert, R. O. Wüest, and C. H. Graham. Accelerated body size evolution in upland environments is correlated with recent speciation in south american freshwater fishes. *Nature Communications*, 14(1):6070, 2023.
- [9] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [10] M. Dal Corso, W. Kirleis, J. Kneisel, N. Taylor, M. Wieckowska-Lüth, and M. Zanon. *How’s Life? Living Conditions in the 2nd and 1st Millennium BCE*. Sidestone Press, 2019.
- [11] C. Dougherty. Numeracy, literacy and earnings: Evidence from the national longitudinal survey of youth. *Economics of education review*, 22(5):511–521, 2003.
- [12] I. Feeser, W. Dörfler, J. Kneisel, M. Hinz, and S. Dreibrodt. Human impact and population dynamics in the neolithic and bronze age: Multi-proxy evidence from north-western central europe. *The Holocene*, 29(10):1596–1606, 2019.
- [13] M. Feurer, A. Klein, K. Eggenberger, J. Springenberg, M. Blum, and F. Hutter. Efficient and robust automated machine learning. In *NeurIPS*, 2015.
- [14] J. Fu, S.-K. Ng, Z. Jiang, and P. Liu. Gptscore: Evaluate as you desire. *ArXiv*, abs/2302.04166, 2023. URL <https://api.semanticscholar.org/CorpusID:256662188>.
- [15] P. Gijsbers, M. L. P. Bueno, S. Coors, E. LeDell, S. Poirier, J. Thomas, B. Bischl, and J. Vanschoren. Amlb: an automl benchmark. *ArXiv*, abs/2207.12560, 2022. URL <https://api.semanticscholar.org/CorpusID:251066648>.
- [16] D. J. Glass and N. Hall. A brief history of the hypothesis. *Cell*, 134(3):378–381, 2008.
- [17] R. Heyard and L. Held. Meta-regression to explain shrinkage and heterogeneity in large-scale replication projects. Technical report, Center for Open Science, 2024.
- [18] H. Jin, F. Chollet, Q. Song, and X. Hu. Autokeras: An automl library for deep learning. *J. Mach. Learn. Res.*, 24:6:1–6:6, 2023.
- [19] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Židek, A. Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [20] E. Kaiser and W. Schier. *Time and Materiality: Periodization and Regional Chronologies at the Transition from Bronze to Iron Age in Eurasia (1200-600 BCE)*. Edited by Elke Kaiser and Wolfram Schier. Verlag Marie Leidorf GmbH, 2021.
- [21] H. Kitano. Artificial intelligence to win the nobel prize and beyond: Creating the engine for scientific discovery. *AI magazine*, 37(1):39–49, 2016.
- [22] J. Kneisel. Chronology and transformation. the transition from bronze to iron age in northern europe. *Time and materiality: Periodization and regional chronologies at the transition from Bronze to Iron Age in Eurasia (1200–600 BCE)*, pages 237–263, 2021.
- [23] P. Langley. Data-driven discovery of physical laws. *Cogn. Sci.*, 5:31–54, 1981. URL <https://api.semanticscholar.org/CorpusID:39694251>.
- [24] E. LeDell and S. Poirier. H2O AutoML: Scalable automatic machine learning. In *Proceedings of the AutoML Workshop at ICML*, 2020.
- [25] M. Y. Li, E. B. Fox, and N. D. Goodman. Automated statistical model discovery with language models. *ArXiv*, abs/2402.17879, 2024. URL <https://api.semanticscholar.org/CorpusID:268041863>.
- [26] X. Li, T. Zhang, Y. Dubois, R. Taori, I. Gulrajani, C. Guestrin, P. Liang, and T. B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval, 2023.
- [27] B. Y. Lin, K. Chandu, F. Brahma, Y. Deng, A. Ravichander, V. Pyatkin, R. L. Bras, and Y. Choi. Wildbench: Benchmarking language models with challenging tasks from real users in the wild, 2024. URL <https://huggingface.co/spaces/allenai/WildBench>.

- [28] X. Liu, Z. Wu, X. Wu, P. Lu, K.-W. Chang, and Y. Feng. Are llms capable of data-based statistical and causal reasoning? benchmarking advanced quantitative reasoning with data. *arXiv preprint arXiv:2402.17644*, 2024.
- [29] L. Lorenz. *Kommunikationsstrukturen mittelneolithischer Gesellschaften im nordmitteleuropäischen Tiefland*. Verlag Dr. Rudolf Habelt GmbH, in Kommission, 2018.
- [30] A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegrefe, U. Alon, N. Dziri, S. Prabhunoye, Y. Yang, S. Gupta, B. P. Majumder, K. Hermann, S. Welleck, A. Yazdanbakhsh, and P. Clark. Self-refine: Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=S37h0erQLB>.
- [31] S.-D. Maida et al. Unter hügeln: Bronzezeitliche transformationsprozesse in schleswig-holstein am beispiel des fundplatzes von mang de bargen (bornhöved, kr. segeberg) band 1, 2023.
- [32] B. P. Majumder, B. D. Mishra, P. Jansen, O. Tafjord, N. Tandon, L. Zhang, C. Callison-Burch, and P. Clark. CLIN: A continually learning language agent for rapid task adaptation and generalization. *arXiv preprint arXiv:2310.10134*, 2023.
- [33] B. P. Majumder, H. Surana, D. Agarwal, S. Hazra, A. Sabharwal, and P. Clark. Data-driven discovery with large generative models. *ICML*, 2024.
- [34] G. I. P. Ottaviano, G. Peri, and G. C. Wright. Immigration, offshoring, and american jobs. *American Economic Review*, 103(5):1925–1959, 2013.
- [35] L. C. Pal. Impact of education on economic development. *Khazanah Pendidikan Islam*, 5(1): 10–19, 2023.
- [36] A. Palmisano, A. Bevan, A. Kabelindde, N. Roberts, and S. Shennan. Long-term demographic trends in prehistoric italy: Climate impacts and regionalised socio-ecological trajectories. *Journal of world prehistory*, 34(3):381–432, 2021.
- [37] E. W. Parkinson, T. R. McLaughlin, C. Esposito, S. Stoddart, and C. Malone. Radiocarbon dated trends and central mediterranean prehistory. *Journal of world prehistory*, 34:317–379, 2021.
- [38] N. Rambeli, D. A. A. Marikan, J. M. Podivinsky, R. Amiruddin, and I. Ismail. The dynamic impact of government expenditure in education on economic growth. *International Journal of Business and Society*, 22(3):1487–1507, 2021.
- [39] M. Riera, J. Pino, L. Sáez, P. Aymerich, and Y. Melero. Effect of introduction pathways on the invasion success of non-native plants along environmental gradients. *Biological Invasions*, pages 1–20, 2024.
- [40] B. Roziere, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, T. Remez, J. Rapin, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- [41] T. Schick, J. Dwivedi-Yu, R. Dessì, R. Raileanu, M. Lomeli, E. Hambro, L. Zettlemoyer, N. Cancedda, and T. Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.
- [42] Z. Shao, F. Wang, Y. Xu, W. Wei, C. Yu, Z. Zhang, D. Yao, G. Jin, X. Cao, G. Cong, C. S. Jensen, and X. Cheng. Exploring progress in multivariate time series forecasting: Comprehensive benchmarking and heterogeneity analysis. *ArXiv*, abs/2310.06119, 2023.
- [43] S. Shashidhar, A. Chinta, V. Sahai, Z. Wang, and H. Ji. Democratizing llms: An exploration of cost-performance trade-offs in self-refined open-source models. *ArXiv*, abs/2310.07611, 2023. URL <https://api.semanticscholar.org/CorpusID:263834891>.
- [44] N. Shinn, F. Cassano, B. Labash, A. Gopinath, K. Narasimhan, and S. Yao. Reflexion: Language agents with verbal reinforcement learning. In *NeurIPS*, 2023. URL <https://api.semanticscholar.org/CorpusID:258833055>.
- [45] P. K. Smith, B. Bogin, and D. Bishai. Are time preference and body mass index associated?: Evidence from the national longitudinal survey of youth. *Economics & Human Biology*, 3(2): 259–270, 2005.
- [46] C. Sommerfeld. *Gerätegeld Sichel: studien zur monetären Struktur bronzezeitlicher Horte im nördlichen Mitteleuropa*, volume 19. Walter de Gruyter, 2013.

- [47] W. H. Thompson and S. Skau. On the scope of scientific hypotheses. *Royal Society Open Science*, 10(8):230607, 2023.
- [48] H. Weatherly, K. Lopez, and C. Tierra. The impact of education on gdp per capita. 2022.
- [49] Wes McKinney. Data Structures for Statistical Computing in Python. In Stéfan van der Walt and Jarrod Millman, editors, *Proceedings of the 9th Python in Science Conference*, pages 56 – 61, 2010. doi: 10.25080/Majora-92bf1922-00a.
- [50] Z. Yang, X. Liu, T. Li, D. Wu, J. Wang, Y. Zhao, and H. Han. A systematic literature review of methods and datasets for anomaly-based network intrusion detection. *Comput. Secur.*, 116: 102675, 2022.
- [51] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. R. Narasimhan, and Y. Cao. ReAct: Synergizing reasoning and acting in language models. In *ICLR*, 2023. URL https://openreview.net/forum?id=WE_vluYUL-X.
- [52] Z. Yuan, J. Liu, Q. Zi, M. Liu, X. Peng, and Y. Lou. Evaluating instruction-tuned large language models on code comprehension and generation. *ArXiv*, abs/2308.01240, 2023. URL <https://api.semanticscholar.org/CorpusID:260379087>.
- [53] K. Zaw, D. Hamilton, and W. Darity. Race, wealth and incarceration: Results from the national longitudinal survey of youth. *Race and Social Problems*, 8:103–115, 2016.
- [54] Z. Zeng, J. Yu, T. Gao, Y. Meng, T. Goyal, and D. Chen. Evaluating large language models at evaluating instruction following. *ArXiv*, abs/2310.07641, 2023. URL <https://api.semanticscholar.org/CorpusID:263834884>.
- [55] S. Zhang, C. Gong, L. Wu, X. Liu, and M. Zhou. AutoML-GPT: Automatic machine learning with GPT. *ArXiv*, abs/2305.02499, 2023.

A FAQs

1. **Dataset or Benchmark:** Is this a dataset or a benchmark? **A benchmark**
2. **Benchmark:** For benchmarks, the supplementary materials must ensure that all results are easily reproducible (i.e., all necessary datasets, code, and evaluation procedures must be accessible and documented)
 - Datasets:** DISCOVERYBENCH is released at <https://github.com/allenai/discoverybench/tree/main/discoverybench>.
 - Code (Baseline Models):** Code for Discovery Agents are provided in the repository, at: <https://github.com/allenai/discoverybench/tree/main/agents>. A CLI is available to run the discovery agents on the benchmark.
 - Evaluation Procedures:** Please follow our main paper for the details of our evaluation process. The code to run eval on a single instance of our benchmark is provided at: <https://github.com/allenai/discoverybench/tree/main/eval>. A CLI and some example scripts have been provided as well.
3. **Accessibility:** The following are accessibility items on the submission checklist:
 - Links to access the benchmark:** The link to access the benchmark is provided in the main submission (<https://github.com/allenai/discoverybench/tree/main/discoverybench>).
 - Any data should use open and widely used formats. Simulation environments should explain how they can be used:** Our data are stored in widely accessible standard formats (e.g., JSON, CSV), with the structure described in Appendix F.
 - Long-term preservation.** Code and data are provided on GITHUB. All aspects will be publicly available for a long term.
 - Explicit Licence:** Our benchmark is licensed using ODC-BY and the associated code is licensed with APACHE 2.0, as included in the GITHUB repository.
 - Structured Metadata for a dataset:** Our dataset is also available as the HuggingFace dataset: <https://huggingface.co/datasets/allenai/discoverybench>. Structured Metadata will be available once we finalize our work after addressing the reviewers’ comments, if any.
 - A persistent dereferenceable identifier (e.g., a code repository such as GitHub):** The repository for our benchmark is: <https://github.com/allenai/discoverybench>.

B Limitations

We currently filtered domains and tasks that required forecasting, simulation, or very specific modeling (species distribution, infection spread, astrophysics equations for exoplanets) in the benchmark as they were very time-consuming to replicate as well as discover hypotheses. As a result, we discarded more papers focused on natural and physical sciences compared to social sciences, which we plan to include in future benchmarks.

We currently do not tackle the challenge of understanding and processing massive datasets, such as the 8.92 petabytes data from the Cancer Genome Atlas (<https://portal.gdc.cancer.gov>) or the extensive brain data from the Allen Institute (<https://alleninstitute.org/division/brain-science>). While the potential to discover new insights from such vast data volumes is significant, ensuring these findings are robust and not subject to p -hacking remains unaddressed by our current methods.

We currently do not handle multi-modal data and complex pipelines, such as those needed for analyzing satellite and other geospatial data relevant to climate science and astronomy data. This would involve multiple stages of data processing, the use of various tools, and managing workflow complexities, for example, analyzing thousands of species patterns combined with satellite data to study habitats. So we do not incorporate workflows like those of EarthRanger (<https://www.earthranger.com>).

Ethical Considerations There could be many potential societal consequences of systems tuned on our proposed benchmark since it involves using LLMs, such as policy misuse, legal ramifications, and false discovery. On the positive side, our proposed benchmark can advance the rate of discovery, leading to an improved standard of living and social well-being.

C Data collection for DB-REAL

For data-first approach, replication took 15 to 40 person-hours for each NLS-related paper and up to 90 person-hours for the GBIF dataset, where specialized domain knowledge and tools led to higher complexity. All papers replicated in the NLS dataset were included, while less than half of the papers in specialized datasets like GBIF and WBOD were added to DISCOVERYBENCH.

Citation/Repositories for DB-REAL: List of scientific works from where we have replicated our gold workflows and hypotheses:

1. Sociology: [53, 3, 1, 45, 11]
2. Biology: [8, 39]
3. Economics: [35, 4, 48, 38, 34]
4. Engineering: [2]
5. Meta-science: [17]
6. Humanities: [7, 6, 29, 31, 46, 12, 37, 22, 20, 5, 10, 36, 37]

All assets come under CC license or open licenses.

D Data Generation for DB-SYNTH

For leaves, we use different sampling strategies based on the data type. Specifically, for categorical nodes, we sample instances with replacement from the range of allowed values, whereas for numeric, we first select a distribution (e.g., normal) and its parameters based on the specified range and then perform sampling. For each subsequent level in \mathcal{F} , we create new columns for nodes by simply executing their pandas expressions⁹. To recover from any execution errors, we additionally use a self-refine [30] approach to generate new pandas expressions guided by the execution error logs. Finally, to mimic real-world challenges in data collection, we probabilistically perturb each instance

⁹The expression is guaranteed to only have variables already generated due to the bottom-up construction.

$x \in \mathbf{x}_i$ by adding noise or dropping values to create missing data¹⁰. After generation, D contains a column for each node in \mathcal{F} .

E Datasheets

E.1 Motivation

- **For what purpose was the dataset created?** DISCOVERYBENCH is created to help assess large language models’ (LLMs) ability to automate the search and verification of hypotheses purely from a set of provided datasets.
- **Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?** Authors belong to the Allen Institute for AI, OpenLocus, and the University of Massachusetts Amherst. The data collection is part of research efforts conducted by the Allen Institute for AI.
- **Who funded the creation of the dataset?** Allen Institute for AI.

E.2 Collection Process

- **How was the data associated with each instance acquired?** Our goal is to replicate the scientific process undertaken by researchers to search for and validate a hypothesis from one or more datasets. We focus on six scientific domains where data-driven research is the cornerstone of scientific progress: sociology, biology, humanities, economics, engineering, and meta-science. Our data collection follows either a **data-first** or **code-first** approach. Each instance has been manually implemented and verified by the authors for solvability.

E.3 Uses

- **Has the dataset been used for any tasks already?** We use this benchmark to evaluate LLM’s ability to search and verify hypotheses purely from a set of datasets.
- **Are there tasks for which the dataset should not be used?** We do not expect the community members to use this data to train models that can aggravate p -hacking.

E.4 Distribution and Maintainance

- **How will the dataset will be distributed?** We distribute this benchmark via our GITHUB repository: <https://github.com/allenai/discoverybench> and <https://huggingface.co/datasets/allenai/discoverybench>.
- **How can the owner/curator/manager of the dataset be contacted?** For any benchmark-related queries, please contact: bodhisattwam@allenai.org. For any code-related discussions, please raise an issue in GITHUB: <https://github.com/allenai/discoverybench>.

F Composition of DISCOVERYBENCH

F.1 Metadata structure

- **id:** An identifier for the metadata.
- **domain:** The broad field of study or area of research.
- **workflow_tags:** A set of keywords summarizing the main processes or techniques used in the replication implementation. They provide an overview of the methodological approach and facilitating the identification of relevant analytical techniques.
- **domain_knowledge:**
 - Contextual information or insights related to the dataset, explaining how certain behaviors or variables can be interpreted within the field of study.

¹⁰Each value is noised independently; therefore, each row has sufficient true data useful for discovery.

- It helps open avenues to think in directions that LLM might not have considered otherwise, broadening the understanding of the field.
- **datasets:** Contains detailed information about the datasets used, including:
 - **name:** The name or filename of the dataset.
 - **description:** A summary of the dataset’s contents and the type of data it includes.
 - **max_depth:** The maximum hierarchical level of nested data structures within the dataset, indicating the complexity of the data.
 - **columns:** Detailed descriptions of each column in the dataset, including:
 - * **name:** The column’s name or header.
 - * **description:** Explanation of the data contained in the column and its significance.
 - * **depth:** The hierarchical level of the column within the dataset, indicating its structural position.
- **hypotheses:** Statements or predictions being tested, divided into:
 - **main:** Primary hypotheses that are central to the discovery task.
- **workflow:** A step-by-step description of the replication process followed to validate the hypotheses, outlining the methods and procedures used from data preparation to final analysis. Some of the workflows and sub-workflows are high-level and thus the same for different queries as they follow the same implementation leading to a range of hypotheses.
- **queries:** Goals related to each hypothesis, each including:
 - **qid:** A unique identifier for the goal for a given true/gold hypothesis.
 - **difficulty:** Categorization of the difficulty. Structurally defined for DB-SYNTH using the semantic tree definition.
 - **true_hypothesis:** The hypothesis being tested through the goal. This defines the primary statement or prediction under investigation.
 - **relevant_cols:** Columns from the dataset that are relevant to answering the query, indicating the specific data points that can be used in the analysis. Only appears for DB-SYNTH.
 - **target_col:** The column being predicted or the dependent variable in the analysis. Only appears for DB-SYNTH.
 - **question_type:** The type of question being asked categorizing the nature of the inquiry.
 - **question:** The discovery goal.

F.2 Directory structure for DB-REAL

There may be more than one query per metadata. The train split contains 14 metadata files and 25 queries. The test split contains 144 metadata files and 239 queries. Metadata folders with the same prefixes use the same underlying dataset with either a subset or a preprocessed version. When dealing with a full dataset (i.e., nls_raw), the task becomes substantially harder due to the data preparation required.

```

|-test
|---archaeology
|---introduction_pathways_non-native_plants
|---meta_regression
|---meta_regression_raw
|---nls_incarceration
|---nls_raw
|---nls_ses
|---requirements_engineering_for_ML_enabled_systems
|---worldbank_education_gdp
|---worldbank_education_gdp_indicators
|-train
|---evolution_freshwater_fish
|---immigration_offshoring_effect_on_employment
|---nls_bmi
|---nls_bmi_raw

```

E.3 Directory structure for DB-SYNTH

There is one query per metadata. The train split contains 551 metadata files (queries), the dev split contains 153 metadata files (queries), and the test split contains 200 metadata files (queries).

```
| -test
| ---ancient-languages_*. *
| ---artificial-ecosystems_*. *
| ---astronomy_*. *
| ---board-games_*. *
| ---coding-competitions_*. *
| ---digital-artistry_*. *
| ---futuristic-technology_*. *
| ---impressionist-art_*. *
| ---machine-learning_*. *
| ---molecular-gastronomy_*. *
| ---neuroscience_*. *
| ---philosophical-debates_*. *
| ---robotics_*. *
| -train
| ---adventure-travel_*. *
| ---ancient-architecture_*. *
| ---ancient-astronomy_*. *
| ---aviation_*. *
| ---biodiversity-conservation_*. *
| ---cryptic-puzzles_*. *
| ---cryptocurrency_*. *
| ---culinary-arts_*. *
| ---cybersecurity_*. *
| ---environmental-activism_*. *
| ---fashion-design_*. *
| ---fine-arts_*. *
| ---literary-classics_*. *
| ---marine-biology_*. *
| ---marine-conservation_*. *
| ---medieval-literature_*. *
| ---musical-therapy_*. *
| ---photography_*. *
| ---robotic-explorers_*. *
| ---solar-power_*. *
| ---space-tourism_*. *
| ---steampunk-culture_*. *
| ---theater-productions_*. *
| ---underwater-archaeology_*. *
| ---urban-gardening_*. *
| ---vintage-automobiles_*. *
| ---virtual-reality_*. *
```

G Discovery Agent

The command `discovery_agent.py` is used with various options to customize its behavior for discovery tasks. Below are the options explained:

- **Usage:** `discovery_agent.py [OPTIONS] QUERY` – Executes the discovery agent with specified options.
- **Options:**
 - `-agent_type [coder|react]`: Specifies the type of agent to use for discovery. The default type is `coder`. Options include `coder` for code-related tasks and `react` for reactive tasks.

- `-model_name` TEXT: Sets the model to be used. The default is `gpt-4o`. Available models include `gpt-4-turbo`, `llama-3-70b-chat`, `claude-3-opus`, and `gemini-pro`. An exhaustive list is available in `config/model_config.json`.
- `-api_config` TEXT: Path to the API configuration file. The default path is `config/api_config.json`.
- `-log_file` TEXT: Specifies the path to the log file where operations details are stored.
- `-metadata_path` TEXT: Path to the metadata file. This option is required.
- `-metadata_type` [`real|synth`]: Specifies the type of metadata, where `real` stands for actual metadata and `synth` for synthetic. This option is required.
- `-add_domain_knowledge`: Includes domain-specific knowledge in the query processing.
- `-add_workflow_tags`: Includes workflow tags in the query to enhance context.
- `-help`: Displays the help message and exits, showing all available command options.

H Evaluation

Explain about evaluation in a line and then explain the CLI usage here.

The command `discovery_eval.py` is used to evaluate the outputs generated by the discovery agent. Below are the detailed descriptions of the command options:

- **Usage:** `discovery_eval.py [OPTIONS] QUERY` – Executes the evaluation agent with specified options and a query.
- **Options:**
 - `-gold_hypo` TEXT: Specifies the gold standard hypothesis for comparison. This field is required.
 - `-gold_workflow` TEXT: Specifies the gold standard workflow to be used as a reference during evaluation.
 - `-pred_hypo` TEXT: Specifies the predicted hypothesis generated by the discovery agent. This field is required.
 - `-pred_workflow` TEXT: Specifies the predicted workflow generated by the discovery agent.
 - `-metadata_path` TEXT: Specifies the path to the metadata file that is utilized during evaluation. This field is required.
 - `-metadata_type` [`real|synth`]: Determines the type of metadata used in the evaluation, where `real` indicates actual metadata and `synth` indicates synthetic metadata. This field is required.
 - `-eval_output_path` TEXT: Specifies where the evaluation results should be saved.
 - `-help`: Displays the help message and exits, detailing all available command options.

I Experiments

For GPT-based models, we use OpenAI API (<https://platform.openai.com/docs/models>), and for Llama3, we used Together API (<https://docs.together.ai/docs/inference-models>)

J Evaluator Prompts

We provide below the exact prompts used for our GPT-4 based evaluation of the generated hypothesis against the gold hypothesis.

Listing 1 Decomposition Prompt to obtain sub-hypotheses from a hypothesis.

```
decomposition_prompt = f"""\n    Given a set of dataset columns, a ground-truth hypothesis, and the\n    analysis workflow used, your task is to extract the set of sub-hypotheses\n    that are present in the hypothesis such that each sub-hypothesis covers a\n    separate context, is self-sufficient, and operates on a coherent set of 3\n    dimensions: Context, Variables, and Relations.\n\n    Here are the definitions for these dimensions:\n\n    - Contexts: Boundary conditions that limit the scope of a sub-hypothesis.\n    E.g., "for men over the age of 30", "in Asia and Europe", or "None" if\n    there is no boundary condition specified.\n\n    - Variables: Known concepts that interact in a meaningful way under a\n    given context to produce the sub-hypothesis. E.g., gender, age, income, or\n    "None" if there is no interacting variable.\n\n    - Relations: Interactions between a given set of variables under a given\n    context to produce the sub-hypothesis. E.g., "quadratic relationship",\n    "inversely proportional", piecewise conditionals, or "None" if there is no\n    interacting relationship.\n\n    Make sure to only use the information present in the hypothesis and the\n    workflow. Do not add any new information.\n    If no sub-hypotheses can be extracted, return an empty list.\n\n    Here is the metadata for the task:\n    ```json\n    {{\n      "datasets": {dataset_metadata},\n      "hypothesis": "{hypothesis}",\n      "workflow": "{workflow}"\n    }}\n    ```\n\n    Return your answer as a JSON object in the following format:\n    ```json\n    {{\n      "sub_hypo": [\n        {{\n          "text": the sub-hypothesis in natural language,\n          "context": a short text description of the context of the\n          sub-hypothesis,\n          "variables": a list of columns involved in the sub-hypothesis,\n          "relations": a short text description of the relationship between\n          the variables of the sub-hypothesis,\n          "explanation": a short text explanation for the breakdown of the\n          sub-hypothesis\n        }}\n        ...]\n      }\n    }}\n    ```\n    """
```

Listing 2 Matching prompt to match contexts of two sub-hypotheses.

```
matching_prompt = f"""
    Given a gold hypothesis, a gold context, a predicted hypothesis, and a
    predicted context, your task is
    to determine if the predicted context semantically matches the
    ground-truth context.

    Here is the definition for Context: Boundary conditions that limit the
    scope of a sub-hypothesis. E.g., “for men over the age of 30”, “in Asia
    and Europe”, or "None" if there is no boundary condition specified.

    If the predicted context matches the gold context, return true, otherwise
    return false.

    Here is the metadata for the task:
    ```json
 {{
 "gold_hypothesis": "{gold_hypothesis}",
 "gold_context": "{gold_context}",
 "predicted_hypothesis": "{pred_hypothesis}",
 "predicted_context": "{pred_context}"
 }}
    ```

    Return your answer as a JSON object in the following format:
    ```json
 {{
 "match": true or false
 }}
 """
```

---

---

**Listing 3** Prompt for variable alignment between two sub-hypotheses.

---

```
main_context = f"""
 You are going to compare two natural-language hypotheses HypoA and HypoB
 accompanied with optional workflows: WorkflowA for HypoA and WorkflowB for
 HypoB.
 Both the hypotheses answer the natural language query "QUERY" over the
 dataset(s) described by dataset description(s) and column description(s)
 below.
 Compare HypoA and HypoB in terms of three aspects: Contexts, Variables,
 and Relations.
 E.g., for the hypothesis "From 1995 to 2009, the number of sandhill cranes
 around the tundra (Indigilka River) surged by an astounding ~10X":
 * Contexts refer to the stratification of the data under which the given
 hypothesis is True. E.g., "For all women", "From 1995 to 2009".
 * Variables refer to the set of variables (either dependent or independent)
 that are mentioned in the hypothesis. E.g., number of sandhill cranes,
 location.
 * Relations refer to the form of relation between the variables. E.g.,
 "surged by ~10x".

 Answer the following questions for a given pair of hypotheses, HypoA and
 HypoB, along with an explanation grounded on the QUERY and the DATASET(S).

 Here is the metadata for the task:
    ```json
    {{
      "datasets": {datasets_json},
      "query": {query},
      "HypoA": {gold_hypo},
      "WorkflowA": {gold_workflow},
      "HypoB": {gen_hypo},
      "WorkflowB": {gen_workflow}
    }}
    ```

 {variable_question}"""
variable_question = """\
 Question: For both HypoA and HypoB, what are the different variables found
 in the hypotheses? \
 Return your answer as a JSON object in the following format:
    ```json
    {{
      "sizeA": num of variables used in HypoA
      "sizeB": num of variables used in HypoB
      "intersection": num of variables common in HypoA and HypoB. Use *fuzzy
      matching* to determine intersection, accounting for paraphrases or
      slightly different surface forms
      "explanation": a short text explanation about the variables
    }}```
    Answer: """
```

Listing 4 Prompt for relationship alignment between two sub-hypotheses.

```
main_context = f"""
    You are going to compare two natural-language hypotheses HypoA and HypoB
    accompanied with optional workflows: WorkflowA for HypoA and WorkflowB for
    HypoB.
    Both the hypotheses answer the natural language query "QUERY" over the
    dataset(s) described by dataset description(s) and column description(s)
    below.
    Compare HypoA and HypoB in terms of three aspects: Contexts, Variables,
    and Relations.
    E.g., for the hypothesis "From 1995 to 2009, the number of sandhill cranes
    around the tundra (Indigilka River) surged by an astounding ~10X":
    * Contexts refer to the stratification of the data under which the given
    hypothesis is True. E.g., "For all women", "From 1995 to 2009".
    * Variables refer to the set of variables (either dependent or independent)
    that are mentioned in the hypothesis. E.g., number of sandhill cranes,
    location.
    * Relations refer to the form of relation between the variables. E.g.,
    "surged by ~10x".

    Answer the following questions for a given pair of hypotheses, HypoA and
    HypoB, along with an explanation grounded on the QUERY and the DATASET(S).

    Here is the metadata for the task:
    ```json
 {{
 "datasets": {datasets_json},
 "query": {query},
 "HypoA": {gold_hypo},
 "WorkflowA": {gold_workflow},
 "HypoB": {gen_hypo},
 "WorkflowB": {gen_workflow}
 }}
    ```

    {variable_question}"""
dimension_question = """
    Question: Does HypoB exhibit the same relation as HypoA?
    Compare using the following example hierarchy of relationships (based on
    specificity): \
    "there exists a relationship" > "positive relationship" > "positive AND
    (linear OR quadratic)" > "positive AND linear."
    Options: A) very similar B) similar but general than HypoA C) different
    Return your answer as a JSON object in the following format:
    ```json
 {{
 "answer": one of the options from A) very similar B) similar but general
 than HypoA C) different
 "explanation": a short text explanation about the relationship comparison
 }}```
 Answer: """
```

---