LLM-BASED MULTI-AGENT BLACKBOARD SYSTEM FOR INFORMATION DISCOVERY IN DATA SCIENCE

Anonymous authors

000

001

002 003 004

006

008

009

010

011

012

013

014

016

017

018

019

021

023

025

026

027

028

030

032

033

034

037

038

040 041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

The rapid advancement of Large Language Models (LLMs) has opened new opportunities in data science, yet their practical deployment is often constrained by the challenge of discovering relevant data within large data lakes. Existing methods struggle with this: single-agent systems are quickly overwhelmed by large, heterogeneous files in the data lakes, while multi-agent systems designed based on a master-slave paradigm depend on a rigid central controller that requires precise knowledge of each sub-agent's capabilities. To address these limitations, we propose a novel multi-agent communication paradigm inspired by the blackboard architecture for traditional AI models. In this framework, a central agent posts requests to a shared blackboard, and autonomous subordinate agents-either responsible for a partition of the data lake or general information retrieval-volunteer to respond based on their capabilities. This design improves scalability and flexibility by eliminating the need for a central coordinator to have prior knowledge of all sub-agents expertise. We evaluate our method on three benchmarks that require explicit data discovery: KramaBench and modified versions of DS-Bench and DA-Code to incorporate data discovery. Experimental results demonstrate that the blackboard architecture substantially outperforms baselines, including RAG and the master–slave multi-agent paradigm, achieving between 13% to 57% relative improvement in end-to-end task success and up to a 9% relative gain in F1 score for data discovery over the best-performing baselines across both proprietary and open-source LLMs. Our findings establish the blackboard paradigm as a scalable and generalizable communication framework for multi-agent systems.

1 Introduction

The recent developments in Large Language Models (LLMs) have introduced new paradigms for data science workflows, enabling natural language-based approaches to data interpretation, transformation, and analysis (Jing et al., 2025; Huang et al., 2024; Hong et al., 2025; Wang et al., 2025). Existing work, however, typically assumes an idealized setting in which relevant datasets are already curated and provided to the model—an assumption that diverges substantially from the practical challenges encountered in real-world data science (Lai et al., 2025). In practice, a substantial fraction of effort is devoted to locating the appropriate data within large and heterogeneous data lakes, often comprising thousands of loosely organized files—a process that constitutes a major bottleneck before any downstream analysis can be performed (Xu et al., 2021). We argue that this stage of data discovery is both a critical and underexplored challenge for applying LLMs effectively.

Previous work on data science tasks that require discovery¹ from a data lake has primarily relied on single-agent systems in which an LLM is given access to all candidate files within its context window and is then asked to solve the problem (Lai et al., 2025). This method suffers from several limitations. First, it is not scalable: as the number of files grows, fitting them into the limited context window of an LLM becomes infeasible. Second, the heterogeneity of files poses a challenge, as a single agent may struggle to effectively analyze, interpret, and integrate diverse forms of information. Third, such systems lack robustness to noise, since the presence of many irrelevant files can overwhelm the model and degrade both reasoning quality and precision. One may argue that Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Kim et al., 2024; Salemi & Zamani, 2024b) provides a solution by choosing a subset of files in the data lake; However, current retrieval techniques are known to perform poorly on tabular and domain-specific data, which are pervasive in data science applications (Yu et al., 2025; Ji et al., 2025; Huang et al., 2022; Gu et al., 2025).

¹Some example tasks are shown in Figure 13 and 14 in Appendix D. They require computing or aggregating information from raw data within a large data lake, where the specific source files are not pre-identified.

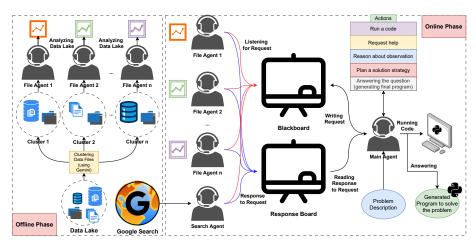


Figure 1: Overview of the blackboard multi-agent system for information discovery in data science. In this framework, the main agent does not assign tasks to subordinate agents. Instead, it posts requests to the blackboard, and subordinate agents autonomously decide whether to respond based on their expertise. The main agent then uses the responses to the request to solve the given task.

An alternative approach explores multi-agent systems, which frequently adopt a master–slave paradigm (Li et al., 2024; Han et al., 2025; Xu & Peng, 2025). In this setting, a single controller (e.g., orchestration agent) assigns subtasks to a set of subordinate agents that then execute the specified actions. While conceptually straightforward, this architecture has several drawbacks. First, this master-slave paradigm limits the agents' autonomy: subordinate agents are forced to execute instructions from the coordinator even when they lack sufficient information or hold outdated or erroneous information. Second, the central controller must maintain an accurate model of each agents capabilities to assign tasks, an assumption that is often unrealistic when agents have only partial or evolving knowledge of the problem space. Finally, when multiple agents possess overlapping expertise, the controller faces an inherent assignment ambiguity, making task routing difficult.

Inspired by the blackboard architecture that with substantial impact on traditional AI systems since the 1980s (Erman et al., 1980; Botti et al., 1995), we adopt a new communication paradigm for LLM multi-agent systems. In this paradigm, a central agent remains responsible for solving the overall task, similar to the master-slave paradigm. However, rather than assigning subtasks to specific agents, the central agent posts a request on a shared *blackboard* that describes the task or information needed, as shown in Figure 1. Subordinate agents monitoring the blackboard can independently decide whether they possess the capability, knowledge, or interest to contribute to solving the task. This design shifts decision-making from a single coordinator to a distributed model whose agents autonomously determine their participation, enabling more flexible collaborations. This differs from the conventional shared-memory paradigm in multi-agent systems. In shared-memory (Sagirova et al., 2025), agents perform assigned tasks based on information in the shared memory, effectively being asked to execute tasks determined by a central coordinator. Conversely, *in the blackboard architecture, there is no task assignment; instead, requests are broadcast on the blackboard, and each agent retains full autonomy to decide whether to participate in solving the task or not.*

While the blackboard architecture can be applied broadly within multi-agent frameworks, its application to data science with data discovery is particularly compelling and underexplored. As shown in Figure 1, the data lake can be partitioned into smaller clusters, e.g., based on similarity, homogeneity, or any criteria that facilitate efficient handling, each assigned to a subordinate agent responsible for understanding and processing that subset. The main agent, which is tasked with solving the given problem, posts requests on the blackboard specifying the data or general information required. Subordinate agents with the relevant knowledge or capability then autonomously volunteer to respond. This design ensures that each sub-agent manages only a subset of files or web-based information, enhancing scalability compared to approaches that require all data to be loaded into the main agent's prompt. Importantly, the main agent does not need prior knowledge of sub agents knowledge or capability to solve the task, simplifying coordination and improving flexibility in large-scale data lake environments. Here, the main agent's role is primarily to describe the information it requires and define tasks for the sub agents, without directly managing or assigning tasks to them.

We conduct experiments on three datasets for data science tasks that require an explicit information discovery phase. KramaBench (Lai et al., 2025) is a recently released benchmark designed for this purpose and, to the best of our knowledge, the only publicly available dataset that directly evaluates data discovery in data science. In addition, we repurpose DS-Bench (Jing et al., 2025) and DA-Code (Huang et al., 2024) by introducing a data discovery component, thereby making them more challenging than their original formulations. Experimental results across these datasets demonstrate that the proposed blackboard architecture consistently outperforms strong baselines, including RAG and the master-slave multi-agent framework, achieving 13% to 57% relative improvement over the best performing baseline in end-to-end problem solving depending on the backbone LLM. Notably, this improvement is observed across both proprietary and open-source LLMs, highlighting the generalizability of the approach. Furthermore, our method also surpasses baselines in data discovery performance, yielding up to a 9% relative gain in F1 score for correctly identifying relevant files from the data lake. These results underscore the effectiveness of the blackboard architecture as a communication paradigm for multi-agent systems in data science.

2 Problem Formulation

Let $\mathbb{D}=\{d_i\}_{i=1}^N$ denote a data lake consisting of N distinct data files, each containing information potentially completely or partially relevant to answering a data science question q (some examples of these questions are shown in Figures 13 and 14 in Appendix D). The objective of this work is to design a generative system π_s that, given the query q and the data lake \mathbb{D} as input, produces a program $p \sim \pi_s(q; \mathbb{D})$ in response. When executed (e.g., using a Python interpreter in this paper), this program p retrieves, loads, and processes the appropriate data from the data lake \mathbb{D} and solve the given problem in the question q to compute the answer. To evaluate the generated program p, we assume the existence of an evaluation function $\mu_{\text{generation}}$ that executes p to produce an output o_p , compares o_p with the ground-truth response y_q , and assigns a corresponding score. In addition, we assume a metric $\mu_{\text{retrieval}}$, given the program p and the ground-truth files \mathbb{D}_q , assigns a score reflecting the performance in discovering the correct data sources.

3 LLM-BASED MULTI-AGENT BLACKBOARD SYSTEM

This section introduces an alternative communication paradigm for LLM-based multi-agent systems inspired by blackboard systems (Erman et al., 1980), distinct from the widely used master–slave architecture. As outlined in §1, blackboard-based multi-agent systems provide several advantages over the master-slave approach. Here, rather than directly assigning tasks to sub-agents, the main agent posts its requests (i.e., sub-tasks for which it requires assistance) on a shared blackboard, which functions as a broadcast channel accessible to all other agents. Each helper agent independently evaluates whether it can respond to a request, considering its own capabilities, availability, cost, and other factors. If an agent decides to contribute, it writes its response to the corresponding request, and the main agent then decides whether to use or ignore the provided information. *This way, all agents in the system retain full autonomy over their actions, and no centralized controller forces them to execute a specific task*. While the blackboard paradigm is applicable to a wide range of multi-agent systems, we focus on data science tasks that require data discovery, where its characteristics are particularly advantageous, as discussed in §1. The remainder of this section details our method and its design for data science problems that require information discovery.

Overview: An overview of our proposed method is presented in Figure 1. The system π_s operates over the data lake $\mathbb D$ by first partitioning $\mathbb D$ into C clusters of related files. Each cluster $\mathbb D_i$ is assigned to a file agent π_{f_i} , which is responsible for handling, loading, processing, and retrieving information from the files within its cluster. In addition, a search agent π_s is included to retrieve external information from the web that may be required to solve the problem. The overall system π_s is composed of a main agent π_m , which is responsible for solving the query q, and a set of C+1 helper agents $\Pi_{\text{helper}}=\{\pi_{f_i}\}_{i=1}^M\cup\{\pi_s\}$ that provide specialized assistance. The query q is presented to π_m , which iteratively selects an action $a\in\mathbb A$ from the action space $\mathbb A$, executes it, and observes the resulting outcome. Among its actions, the main agent may interact with a blackboard β , a shared communication medium where it can post a request r without addressing a specific sub-agent. The helper agents Π_{helper} continuously monitor the blackboard, determine

whether they can address a posted request, and, if so, provide their outputs on the corresponding response board β_r . These responses are then collected and made available to π_m , which incorporates them into its decision-making process.² The main agent is limited to at most T sequential actions (including blackboard interactions) to solve the query q, ultimately producing a program p in python programming language that computes the final answer to q.

Clustering Data Lake: There are multiple approaches for partitioning the data lake into clusters; applying clustering algorithms over file representations, random partitioning, or other heuristic methods. For simplicity, we do not utilize file content and instead rely solely on file names. Specifically, the file names are provided to an LLM—Gemini-2.5-Pro³—which using the prompt shown in Figure 5, clusters the files into categories based only on their names.⁴ An example of this clustering is provided in Figure 12 in Appendix D, where the model successfully groups related files together. For instance, it clusters all files originating from the National Interagency Fire Center into a category labeled "NIFC Wildfire Statistics." The number of automatically derived clusters for each dataset is reported in Table 3 in Appendix A.

3.1 MAIN AGENT

The primary role of the main agent is to solve the problem in collaboration with the helper agents. The main agent follows the ReAct framework (Yao et al., 2023), where at each step t, given the query q and the history of actions and observations \mathbb{H}_{t-1} , it first reasons about what is the best next action and selects an action from a predefined action space, executes the action, observes the outcome, and appends the resulting observation to update the history \mathbb{H}_t . The prompt used by the main agent is shown in Figure 6 in Appendix B. The agent selects one of the following predefined actions in each step, executes them, and observe their outcomes:

- *Planning:* In this action, the LLM decomposes the problem into smaller sub-problems and outlines a plan for addressing each of them. This action has no external effect on the environment but serves as an internal reasoning step to guide the LLM's problem-solving process. In response, the system simply acknowledges the proposed plan and instructs the LLM to proceed.
- *Reasoning:* In this action, the LLM focuses on a specific aspect of the problem and explains its reasoning, analysis, or interpretation of the available observations and steps taken so far in this process. Similar to the planning step, this action has no external effect on the environment but functions as an internal reasoning mechanism to guide the LLM's problem-solving process. In response, the system simply acknowledges the reasoning and prompts the LLM to continue.
- Executing Code: In this action, the agent generates python code, which is executed using a python interpreter. If the code runs successfully, the resulting outputs are returned to the agent for observation; otherwise, the agent receives the corresponding error messages. This action enables the agent to explore the problem interactively, inspect data files, and experiment with them to gain a deeper understanding of their content and structure and how to process them.
- Requesting Help: In this action, the agent formulates a request for assistance from the sub-agents, specifying, for example, the types of data files or information needed, or the resources required to apply a tool or solve a sub-problem. This request is posted on the blackboard β for visibility by the helper agents. Once the sub-agents respond, if they respond, their responses on the response board β_r are collected and provided back to the main agent as the outcome of this action for observation and further use in its decision-making process.
- Answering: In this action, the agent concludes the problem-solving process by generating a final program that produces the answer to the query. This action terminates the process, and the output of this step constitutes the final program p generated by the system to address the problem.

²Responses are not written back to the blackboard β to avoid dependencies where one sub-agent's output could influence the behavior of others negatively. Instead, all responses are directed exclusively to the response board β_r , ensuring independent operation of sub-agents and exclusive access by the main agent π_m .

³Available at: https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-pro

⁴This method represents just one possible approach to clustering, chosen for simplicity; more scalable alternatives could equally be employed in real world scenarios.

⁵In this work, the inputs, outputs of the model, and observations are appended directly to the prompt of the LLM, formatted according to its chat-based input template.

3.2 Helper Agents

In a data science, information discovery can typically be categorized into two tasks: (1) identifying the specific files that contain the data necessary to the problem, and (2) retrieving general knowledge about concepts relevant to the problem, such as domain-specific terms or details of particular algorithms and methods. To support these, our framework employs two types of helper agents:

File Agent: Handling all the files in a data lake with a single agent is not feasible for several reasons: it typically involve a large number of files, many of which are lengthy and may exceed the agents context window; the files span diverse topics, which can confuse the agent and hinder effective reasoning; and accessing and processing all files simultaneously can be computationally expensive and inefficient, leading to unnecessary overhead and slower problem-solving. For these reasons, in our framework each file agent is assigned responsibility for a subset of data files determined to be relevant, as described earlier in the clustering procedure. In an offline phase, the file agent π_{f_i} takes as input a subset of the data lake \mathbb{D}_i and operates through a two-step procedure. In the first step, the agent selects a subset⁶ (or all) of the files to examine their content. The contents of them are presented to the agent for inspection (details of presentation are in Appendix C). In the second step, after observing the selected files, the agent reasons about and analyzes them, learning how they are structured, what pre-processing or transformations may be required, and how they should be processed in general. An example of such an analysis is provided in Figure 11 in Appendix D. Then, in the online phase, the agent listens for requests from the main agent. Upon receiving a request, based on the analysis it did earlier, it determines whether it can contribute to answering it. If so, the agent generates a detailed plan specifying which files in \mathbb{D}_i are relevant, how they should be loaded in Python code, what libraries to use, the steps required for data processing, and samples from the data. The prompt used to guide the file agent is shown in Figure 7 in Appendix B.

Search Agent: Certain data science problems require task-specific knowledge about algorithms or domain expertise that the LLM may not possess. To address this, we design a web-search agent that retrieves relevant information from a search engine. This agent operates according to the prompt shown in Figure 8 in Appendix B. Given a request r posted on the blackboard β , the agent first determines whether it is capable of addressing the request. It is specifically restricted to general web-based information retrieval and does not respond to requests involving access to local files or datasets. If the agent determines that the request can be answered, it enters an iterative search process with a maximum of $T_{\text{search}} = 3$ steps. At each step t, the agent generates a set of queries \mathbb{Q}_t , which are submitted to a search engine—in this work, Google Custom Search Engine⁷—to retrieve k=3 webpage per query. The content of the webpages are then extracted using beautifulsoup library⁸ to be presented to the search agent. The extracted documents are then evaluated by the agent to determine whether they provide sufficient information to answer the request. If so, the agent generates a response to the request, which is posted to the response board β_r . If the information is insufficient, a new set of queries is generated to continue gathering relevant data from the web.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Benchmarks: To the best of our knowledge, KramaBench is the only public benchmark for data science problems that explicitly incorporates a data discovery phase, which we adopt in our evaluation. In addition, we repurpose two existing datasets, DS-Bench (Jing et al., 2025) and DA-Code (Huang et al., 2024), to include in this phase. Specifically, we manually filtered out all questions that do not require any data file for answering, as well as those that lack sufficient hints for data discovery. After filtering, we aggregated all remaining files across questions into a unified data lake, such

⁶When filenames indicate multiple files containing the same type of data over different time periods, the agent does not need to inspect all of them to infer the structure; a small representative sample is sufficient.

⁷We use Google Custom Search Engine, configured to exclude all websites associated with the datasets used in this paper to prevent data leakage: https://developers.google.com/custom-search

⁸Available at: https://pypi.org/project/beautifulsoup4/

⁹For example, questions that request the computation of a data science metric on a column without specifying the structure or content of the relevant file.

that the model must perform discovery to identify relevant files at inference time. In this setup, only the question and the data lake are provided to the model, requiring it to identify the relevant files to answer the question, following the same protocol as KramaBench. Further details on this filtering process, along with dataset statistics in Table 3, are provided in Appendix A.

Evaluation: To evaluate the generated programs, we execute each and compare its output against the ground-truth reference for the corresponding question. For each dataset, we adopt its standard evaluation protocol. For KramaBench, we use the official evaluation script provided in its repository. For DA-Code, we likewise rely on the official evaluation script released by its authors. To DS-Bench, we use the original evaluation method, in which an LLM serves as the judge. The generated programs output is compared against the reference answer using Gemini-2.5-Pro as the judge LLM, with the evaluation prompt shown in Figure 4 in Appendix A, producing a binary score.

Inference Setup: We set the maximum actions of the main agent to T=10. We use nucleus sampling (Holtzman et al., 2020) with a temperature of 0.1 for more deterministic inference and default value for other hyperparameters. Proprietary models are accessed via Vertex AI,¹² while open-source models are served by vLLM.¹³ At each step, we cap the number of generated tokens at 8,192. We use Gemini-2.5-Pro and -Flash (Gemini-Team, 2025), and Claude-4-Opus (Anthropic, 2025) as the proprietary and Qwen3-Coder¹⁴ with 30 billion parameters (Qwen-Team, 2025) as the open-source LLMs. Experiments are conducted on 2 NVIDIA A100 (80GB VRAM) GPUs.

Baselines: To evaluate our method against alternative approaches for solving data science problems involving data discovery, we compare it with the following baselines:

- **DS-GRU:** We adopt the only existing baseline (to the best of our knowledge) for data discovery in data science problems, which appends all available files directly into the LLM prompt and attempts to solve the problem (Lai et al., 2025). This baseline uses a self-correction loop that retries when errors occur in generated codes. For details, we refer the reader to Lai et al. (2025).
- Retrieval-Augmented Generation (RAG): This retrieves the top 5 files¹⁵ based on the file names and contents (the method for presenting a file content to the LLM is explained in Appendix C) from the data lake using E5-large¹⁶ (Wang et al., 2022), a 330M-parameter embedding model and use it to solve the problem. It then follows the same procedure as the main agent described in Section 3.1, with two key modification: 1) the retrieved files contents and addresses are presented directly to the LLM in the prompt and 2) the general help-request action is replaced with a restricted action that only allows direct requests to the search agent. This design isolates the effect of substituting the file discovery mechanism with RAG, enabling a controlled study of its impact on performance. The prompt used for this baseline is shown in Figure 10 in Appendix B.
- *Master-Slave:* This baseline follows the same procedure as the main agent described in Section 3.1. The key difference is that, instead of posting requests on the blackboard, the agent directly invokes sub-agents (consisting of the search agent and the file agents as explained in Section 3.2) based on their description by referencing their names and assign task to them. The prompt used for this baseline is shown in Figure 9 in Appendix B.

4.2 EMPIRICAL FINDINGS

Main Results: We conduct our experiments on the datasets described in Section 4.1 using our method and the baselines. The results are presented in Table 1. These results demonstrate that our method, the Blackboard System, outperforms all baselines on average across all the datasets. Specifically, the Blackboard System surpasses the DS-GRU, RAG and Master-Slave approaches on all three

Available at: https://github.com/mitdbg/KramaBench

¹¹Available at: https://github.com/yiyihum/da-code

¹²Available at: https://cloud.google.com/vertex-ai?hl=en

¹³Available at: https://docs.vllm.ai/en/latest/

¹⁴Available at: https://huggingface.co/Qwen/Qwen3-Coder-30B-A3B-Instruct

¹⁵This number is chosen based on the average number of files required to solve the problems (1.6) and the length of the context window of the backbone LLMs used in this paper.

¹⁶Available at: https://huggingface.co/intfloat/e5-large-v2

Table 1: Results on the KramaBench, DS-Bench, and DA-Code benchmarks. The best results for each LLM are highlighted in **bold**. The KramaBench categories are abbreviated: Arc. (Archaeology), Ast. (Astronomy), Bio. (Biomedical), Env. (Environment), Leg. (Legal), and Wild. (Wildfire).

	Method	LLM		KramaBench						DS-	DA-	Average
	Method		Arc.	Ast.	Bio.	Env.	Leg.	Wild.	Average	Bench	Code	(macro)
(1) (2) (3)	DS-GRU RAG Master-Slave	Qwen3- Coder	0.00% 0.00% 0.00%	1.80% 3.16% 3.55%	2.11% 4.99% 3.39%	1.15% 0.54% 7.77 %	3.27% 6.19% 8.90 %	13.54% 16.93% 21.79%	3.64% 5.30% 7.56%	0.00% 6.32% 7.55%	0.00% 0.00% 0.00%	1.21% 3.87% 5.03%
(4)	Blackboard		0.00%	7.69%	7.85%	4.47%	6.36%	23.97%	8.39%	14.22%	1.11%	7.90%
(5) (6) (7)	DS-GRU RAG Master-Slave	Gemini 2.5 Flash	0.00% 16.66% 16.66%	7.83% 3.57% 3.16% 3.57%	0.09% 13.98% 13.98%	10.93% 28.57% 17.46% 22.92%	12.46% 10.97% 21.75% 27.09 %	13.34% 33.67% 25.80% 41.04 %	7.44% 17.90% 16.46% 21.01%	5.53% 22.92% 26.48% 28.06%	0.00% 2.75% 0.55%	4.32% 14.52% 14.49%
(8)	Blackboard	<u> </u>									0.55%	
(9) (10) (11)	DS-GRU RAG Master-Slave	Gemini 2.5 Pro	25.00% 33.33% 33.33%	6.69% 8.47% 8.47%	10.64% 32.53% 24.74%	27.47% 31.36% 32.81%	5.94% 25.55% 34.64%	39.36% 38.32% 58.98%	19.18% 28.26% 32.16%	3.95% 27.27% 34.38%	0.00% 0.00% 5.49%	7.71% 18.51% 24.01%
(12)	Blackboard		33.33%	17.95%	36.83%	39.31%	34.92%	62.88%	37.53%	38.73%	9.34%	28.53%
(13) (14) (15)	DS-GRU RAG Master-Slave	Claude 4 Opus	8.33% 33.33% 33.33%	1.38% 11.52% 8.69%	1.90% 23.42% 32.28%	8.14% 31.61% 39.16 %	9.80% 31.80% 44.08 %	23.14% 45.80% 48.35%	8.78% 29.58% 34.31%	3.55% 35.57% 45.84%	0.00% 3.85% 2.75%	4.11% 23.00% 27.63%
(16)	Blackboard		33.33%	18.69%	45.31%	34.35%	42.48%	50.06%	37.37%	49.80%	7.14%	31.43%

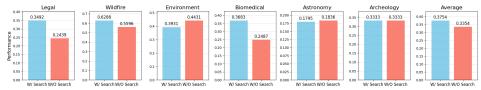


Figure 2: Performance of Blackboard System w/ and w/o search agent (Gemini 2.5 Pro).

datasets and achieves similar or higher performance in 4 out of 6 categories on KramaBench. Furthermore, we observe that the Blackboard System consistently outperforms the baselines regardless of the backbone LLM, highlighting its robustness and generalizability. We attribute this improvement to the design of the Blackboard System, where tasks are not explicitly assigned to helper agents; instead, each agent autonomously decides whether to participate based on its capabilities. This self-selection enhances both problem-solving efficiency and data discovery performance.

File Discovery Performance: To analyze the effectiveness of different methods in data discovery, we report recall, precision, and F1-score for the file discovery task, i.e., identifying the correct files required to answer each question. The results of this experiment, using Gemini 2.5 Pro as the backbone LLM, are presented in Table 2. The results in this table indicate that the blackboard system achieves the highest recall, precision, and F1-score compared to all baselines, both on average and across the three datasets. In particular, for KramaBench, the blackboard system attains the highest F1-score in 4 out of 6 domains. We attribute this improvement to the design of the blackboard system, where the main agent does not directly assign requests to specific file agents, as in the master—slave setup. Instead, each file agent independently decides whether it can contribute based on its capabilities and data holdings, leading to more accurate and comprehensive file discovery.

Effect of Web Search (Search Agent) on the Performance: We observed that in some cases the backbone LLM lacks the necessary domain-specific knowledge or familiarity with specialized algorithms to fully understand and solve the problem. To address this limitation, the inclusion of a search agent that can retrieve relevant external information may be beneficial. To evaluate this, we compare the blackboard system with and without the search agent. The results on KramaBench, shown in Figure 2 using Gemini 2.5 Pro as the backbone LLM, demonstrate that incorporating the search agent improves the average performance of the blackboard system. Further analysis reveals that when the main agent encounters unfamiliar concepts, it issues requests to obtain such information from the web. In these cases, the search agent typically responds by retrieving the required knowledge, thereby enabling the main agent to continue solving the problem effectively. Illustrative examples of this behavior are provided in Figures 13 and 14 in Appendix D, highlighting the importance of the search agent in scenarios where external domain knowledge is essential.

Table 2: File discovery performance, obtained using Gemini-2.5-Pro as the LLM. The best results are highlighted in **bold**. The KramaBench categories are abbreviated: Arc. (Archaeology), Ast. (Astronomy), Bio. (Biomedical), Env. (Environment), Leg. (Legal), and Wild. (Wildfire).

	Method	Metric	KramaBench							DS-	DA-	Average
			Arc.	Ast.	Bio.	Env.	Leg.	Wild.	Average	Bench	Code	(macro)
(1)	RAG	recall precision F1	0.875 1.000 0.916	0.125 0.125 0.125	0.666 0.666 0.629	0.3506 0.450 0.332	0.127 0.133 0.105	0.238 0.452 0.301	0.396 0.471 0.401	0.035 0.047 0.034	0.257 0.456 0.307	0.229 0.324 0.247
(2)	Master- Slave	recall precision F1	0.916 0.930 0.913	0.5138 0.750 0.577	0.648 0.722 0.674	0.382 0.500 0.389	0.444 0.494 0.450	0.567 0.642 0.576	0.578 0.673 0.596	0.323 0.503 0.358	0.546 0.767 0.584	0.482 0.647 0.513
(3)	Black- board	recall precision F1	0.916 1.000 0.944	0.576 0.733 0.618	0.648 0.722 0.674	0.604 0.703 0.588	0.383 0.302 0.304	0.464 0.603 0.495	0.598 0.677 0.603	0.402 0.584 0.438	0.600 0.837 0.643	0.533 0.699 0.561

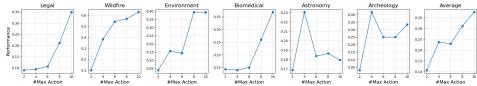


Figure 3: Performance of Blackboard System with various maximum actions by the main agent.

Effect of Number of Main Agent's Actions on the Performance: To examine the impact of the maximum number of actions available to the main agent, we vary this parameter across 2, 4, 6, 8, 10 and evaluate the blackboard system on KramaBench using Gemini 2.5 Pro as the backbone LLM. The results, presented in Figure 3, indicate that increasing the action budget consistently improves the average performance of the system. This trend aligns with intuition: a larger exploration budget allows the agent to more thoroughly analyze the problem, consider alternative strategies, better investigate the solution space, and generate a better program that answers the question.

Case Studies: To qualitatively analyze the blackboard system—specifically how it formulates requests and how this process improves the generated program—we present several case studies:

- Writing Request on the blackboard: An example of a request posted by on the blackboard is shown in Figure 15 in Appendix D. In this case, the main agent, given the data science question, formulates a request that specifies the likely column names and data formats needed to solve the problem, along with some guidance for interpretation. In response, several helper agents (3 out of 8 in this example) chose to contribute. Although the relevant files were distributed across different clusters managed by different file agents, each responding agent independently provided the file addresses, code snippets for loading the data, and explanations of the data structure along with suggested preprocessing steps. Collectively, these responses covered all the ground-truth files required to answer the question. This case study demonstrates how the main agent can effectively leverage the blackboard mechanism to discover and integrate necessary information.
- Comparing Generated Program by Blackboard System with Master-Slave System: To study this further, we present an example of programs generated by the Blackboard system and the Master-Slave system in Figure 16 in Appendix D. In this case, the Blackboard agent achieved a better solution because it accurately interpreted the prompt and selected the correct data files. Specifically, it identified that the patients Age was located in the mmc1.xlsx file and, more importantly, that the requested APP-Z score was in the mmc7.xlsx file. In contrast, the Master-Slave agent misinterpreted the request and instead used a general protein abundance score (APP_log2_abundance) from the wrong file, mmc2.xlsx. This critical error in data selection led the Master-Slave agent to produce an incorrect result of 74, while the Blackboard agents precise data discovery and reasoning yielded the correct answer of 60.

5 RELATED WORK

LLMs for Data Science: Specialized benchmarks have emerged to evaluate LLMs in data science. DS-1000 (Lai et al., 2023), ARCADE (Yin et al., 2023), DataSciBench (Zhang et al., 2025),

and DSEval (Zhang et al., 2024) assess the translation of natural language instructions into correct implementations, distinguishing them from broader programming benchmarks such as SWE-Bench (Jimenez et al., 2024), ML-Bench (Tang et al., 2025), and BigCodeBench (Zhuo et al., 2025). While most assume that the relevant data files are pre-specified, recent efforts address multi-step reasoning: DS-Bench (Jing et al., 2025) and BLADE (Gu et al., 2024) evaluate implementation planning, and ScienceAgentBench (Chen et al., 2025) and BixBench (Mitchener et al., 2025) focus on integrating domain knowledge. These benchmarks, however, still overlook the practical challenge of discovering relevant data within large, heterogeneous repositories—a gap addressed by KramaBench (Lai et al., 2025), which explicitly evaluates data discovery. Building on this, we study how agents can autonomously identify and leverage the correct data sources for end-to-end analysis.

Applications of LLMs in data science have evolved from single-turn code generation to interactive, tool-augmented agents that exploit models specialized for code, including GPT (Brown et al., 2020), CodeGen (Rubavicius et al., 2025), StarCoder (Li et al., 2023), and Code Llama (Rozière et al., 2024). While few-shot prompting (Brown et al., 2020) remains effective, state-of-the-art approaches adopt agentic or multi-agentic frameworks that combine iterative reasoning with external tool use. ReAct (Yao et al., 2023) pioneered the interleaving of reasoning and action, later extended to execution environments (Chen et al., 2019). Toolformer (Schick et al., 2023) and Gorilla (Patil et al., 2024) explicitly train LLMs to call APIs, a capability critical for tasks relying on specialized libraries. Self-correction is a another key feature: frameworks like Self-Debug (Chen et al., 2024) and Reflexion (Shinn et al., 2023) refine generated code using execution feedback. To further enhance reliability, many systems integrate RAG (Lewis et al., 2020; Salemi et al., 2025; Salemi & Zamani, 2025; 2024a) to retrieve documentation or code examples, reducing hallucinations and ensuring up-to-date library use. Additionally, multi-agent master-slave frameworks, such as AutoKaggle (Li et al., 2024), have demonstrated promising results in addressing these challenges.

Blackboard Systems: The blackboard system is a seminal architectural model from classical AI, developed for complex problems that require incremental and opportunistic reasoning. It was implemented in the Hearsay-II speech understanding (Erman et al., 1980) and is characterized by three components: (1) a global, hierarchical data structure (the blackboard) that maintains the current state of the solution; (2) independent specialist modules, known as knowledge sources, which monitor the blackboard and contribute partial solutions; and (3) a control mechanism that opportunistically determines which knowledge source to activate next (Nii, 1986). Following successful applications in domains such as sonar interpretation with the HASP/SIAP system (Nii et al., 1982), the architecture evolved to incorporate more sophisticated control strategies. Inspired by this paradigm, we adapt the blackboard architecture for multi-agent communication: rather than a central controller assigning tasks, all agents operate autonomously, responding to requests posted on the blackboard. A central main agent then leverages the information contributed by sub-agents to solve the problem.

6 Conclusions & Future Work

We addressed the critical challenge of data discovery in large, heterogeneous data lakes, a key bottleneck for applying LLMs in data science. We introduced a novel multi-agent communication paradigm based on the blackboard architecture, which replaces rigid centralized task assignment with a flexible, decentralized model of agent collaboration. Extensive experiments on three data science benchmarks demonstrate that our framework consistently outperforms strong baselines, including RAG and the master–slave paradigm, achieving up to 57% relative improvement in end-to-end task success and a 9% relative gain in data discovery accuracy. These results highlight the importance of communication architecture in multi-agent systems and establish the blackboard paradigm as a scalable, flexible, and effective solution for complex data science workflows.

Future work could extend the blackboard architecture and paradigm beyond data science, as the proposed approach is general and applicable to a wide range of multi-agent systems and domains. Another promising direction is to investigate more adaptive strategies for data partitioning among agents, enabling the system to better handle dynamic and evolving data environments. Ultimately, our findings point toward a broader path for developing more capable, scalable, and autonomous multi-agent AI systems for real-world data analysis applications.

REFERENCES

- Anthropic. System card: Claude opus 4 & claude sonnet 4. https://www-cdn.anthropic.com/4263b940cabb546aa0e3283f35b686f4f3b2ff47.pdf, May 2025. Also referenced in the official release blog post: https://www.anthropic.com/news/claude-4.
- V. Botti, F. Barber, A. Crespo, E. Onaindia, A. Garcia-Fornes, I. Ripoll, D. Gallardo, and L. Hernández. A temporal blackboard for a multi-agent environment. Data & Knowledge Engineering, 15(3):189-211, 1995. ISSN 0169-023X. doi: https://doi.org/10.1016/0169-023X(95)00007-F. URL https://www.sciencedirect.com/science/article/pii/0169023X9500007-F.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- Xinyun Chen, Chang Liu, and Dawn Song. Execution-guided neural program synthesis. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=HlgfOiAqYm.
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=KuPixIqPiq.
- Ziru Chen, Shijie Chen, Yuting Ning, Qianheng Zhang, Boshi Wang, Botao Yu, Yifei Li, Zeyi Liao, Chen Wei, Zitong Lu, Vishal Dey, Mingyi Xue, Frazier N. Baker, Benjamin Burns, Daniel Adu-Ampratwum, Xuhui Huang, Xia Ning, Song Gao, Yu Su, and Huan Sun. Scienceagentbench: Toward rigorous assessment of language agents for data-driven scientific discovery. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=6z4YKr0GK6.
- Lee D. Erman, Frederick Hayes-Roth, Victor R. Lesser, and D. Raj Reddy. The hearsay-ii speech-understanding system: Integrating knowledge to resolve uncertainty. *ACM Comput. Surv.*, 12(2): 213–253, June 1980. ISSN 0360-0300. doi: 10.1145/356810.356816. URL https://doi.org/10.1145/356810.356816.
- Gemini-Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025. URL https://arxiv.org/abs/2507.06261.
- Ken Gu, Ruoxi Shang, Ruien Jiang, Keying Kuang, Richard-John Lin, Donghe Lyu, Yue Mao, Youran Pan, Teng Wu, Jiaqian Yu, Yikun Zhang, Tianmai M. Zhang, Lanyi Zhu, Mike A Merrill, Jeffrey Heer, and Tim Althoff. BLADE: Benchmarking language model agents for data-driven science. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), Findings of the Association for Computational Linguistics: EMNLP 2024, pp. 13936–13971, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. findings-emnlp.815. URL https://aclanthology.org/2024.findings-emnlp.815/.
- Ken Gu, Zhihan Zhang, Kate Lin, Yuwei Zhang, Akshay Paruchuri, Hong Yu, Mehran Kazemi, Kumar Ayush, A. Ali Heydari, Maxwell A. Xu, Girish Narayanswamy, Yun Liu, Ming-Zher Poh, Yuzhe Yang, Mark Malhotra, Shwetak Patel, Hamid Palangi, Xuhai Xu, Daniel McDuff, Tim Althoff, and Xin Liu. Radar: Benchmarking language models on imperfect tabular data, 2025. URL https://arxiv.org/abs/2506.08249.

Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, and Zhaozhuo Xu. Llm multi-agent systems: Challenges and open problems, 2025. URL https://arxiv.org/abs/2402.03578.

- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=rygGQyrFvH.
- Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Ceyao Zhang, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Li Zhang, Lingyao Zhang, Min Yang, Mingchen Zhuge, Taicheng Guo, Tuo Zhou, Wei Tao, Robert Tang, Xiangtao Lu, Xiawu Zheng, Xinbing Liang, Yaying Fei, Yuheng Cheng, Yongxin Ni, Zhibin Gou, Zongze Xu, Yuyu Luo, and Chenglin Wu. Data interpreter: An LLM agent for data science. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational Linguistics:* ACL 2025, pp. 19796–19821, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/2025.findings-acl.1016. URL https://aclanthology.org/2025.findings-acl.1016/.
- Junjie Huang, Wanjun Zhong, Qian Liu, Ming Gong, Daxin Jiang, and Nan Duan. Mixed-modality representation learning and pre-training for joint table-and-text retrieval in OpenQA. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 4117–4129, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.303. URL https://aclanthology.org/2022.findings-emnlp.303/.
- Yiming Huang, Jianwen Luo, Yan Yu, Yitong Zhang, Fangyu Lei, Yifan Wei, Shizhu He, Lifu Huang, Xiao Liu, Jun Zhao, and Kang Liu. DA-code: Agent data science code generation benchmark for large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 13487–13521, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.748. URL https://aclanthology.org/2024.emnlp-main.748/.
- Xingyu Ji, Parker Glenn, Aditya G. Parameswaran, and Madelon Hulsebos. Target: Benchmarking table retrieval for generative tasks, 2025. URL https://arxiv.org/abs/2505.11545.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.
- Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming Zhang, Xinya Du, and Dong Yu. DSBench: How far are data science agents from becoming data science experts? In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=DSsSPrORZJ.
- To Eun Kim, Alireza Salemi, Andrew Drozdov, Fernando Diaz, and Hamed Zamani. Retrieval-enhanced machine learning: Synthesis and opportunities, 2024. URL https://arxiv.org/abs/2407.12982.
- Eugenie Lai, Gerardo Vitagliano, Ziyu Zhang, Sivaprasad Sudhir, Om Chabra, Anna Zeng, Anton A. Zabreyko, Chenning Li, Ferdi Kossmann, Jialin Ding, Jun Chen, Markos Markakis, Matthew Russo, Weiyang Wang, Ziniu Wu, Michael J. Cafarella, Lei Cao, Samuel Madden, and Tim Kraska. Kramabench: A benchmark for ai systems on data-to-insight pipelines over data lakes, 2025. URL https://arxiv.org/abs/2506.06541.
- Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Wen-tau Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: a natural and reliable benchmark for data science code generation. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.

- Raymond Li, Loubna Ben allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia LI, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Joel Lamy-Poirier, Joao Monteiro, Nicolas Gontier, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Ben Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason T Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Urvashi Bhattacharyya, Wenhao Yu, Sasha Luccioni, Paulo Villegas, Fedor Zhdanov, Tony Lee, Nadav Timor, Jennifer Ding, Claire S Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro Von Werra, and Harm de Vries. Starcoder: may the source be with you! *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=KoFoq41haE. Reproducibility Certification.
- Ziming Li, Qianbo Zang, David Ma, Jiawei Guo, Tuney Zheng, Minghao Liu, Xinyao Niu, Yue Wang, Jian Yang, Jiaheng Liu, Wanjun Zhong, Wangchunshu Zhou, Wenhao Huang, and Ge Zhang. Autokaggle: A multi-agent framework for autonomous data science competitions, 2024. URL https://arxiv.org/abs/2410.20424.
- Ludovico Mitchener, Jon M Laurent, Benjamin Tenmann, Siddharth Narayanan, Geemi P Wellawatte, Andrew White, Lorenzo Sani, and Samuel G Rodriques. Bixbench: a comprehensive benchmark for llm-based agents in computational biology, 2025. URL https://arxiv.org/abs/2503.00096.
- H. Penny Nii. The blackboard model of problem solving and the evolution of blackboard architectures. *AI Magazine*, 7(2):38, Jun. 1986. doi: 10.1609/aimag.v7i2.537. URL https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/537.
- H. Penny Nii, Edward A. Feigenbaum, and John J. Anton. Signal-to-symbol transformation: Hasp/siap case study. AI Magazine, 3(2):23, Jun. 1982. doi: 10.1609/aimag. v3i2.368. URL https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/368.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive APIs. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=tBRNC6YemY.
- Qwen-Team. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2024. URL https://arxiv.org/abs/2308.12950.
- Rimvydas Rubavicius, Antonio Valerio Miceli-Barone, Alex Lascarides, and Subramanian Ramamoorthy. Conversational code generation: a case study of designing a dialogue system for generating driving scenarios for testing autonomous vehicles, 2025. URL https://arxiv.org/abs/2410.09829.
- Alsu Sagirova, Yuri Kuratov, and Mikhail Burtsev. Shared memory for multi-agent lifelong pathfinding, 2025. URL https://openreview.net/forum?id=9DrPvYCETp.

- Alireza Salemi and Hamed Zamani. Towards a search engine for machines: Unified ranking for multiple retrieval-augmented large language models. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, pp. 741–751, New York, NY, USA, 2024a. Association for Computing Machinery. ISBN 9798400704314. doi: 10.1145/3626772.3657733. URL https://doi.org/10.1145/3626772.3657733.
- Alireza Salemi and Hamed Zamani. Evaluating retrieval quality in retrieval-augmented generation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, pp. 2395–2400, New York, NY, USA, 2024b. Association for Computing Machinery. ISBN 9798400704314. doi: 10.1145/3626772.3657957. URL https://doi.org/10.1145/3626772.3657957.
- Alireza Salemi and Hamed Zamani. Learning to rank for multiple retrieval-augmented models through iterative utility maximization. In *Proceedings of the 2025 International ACM SI-GIR Conference on Innovative Concepts and Theories in Information Retrieval (ICTIR)*, ICTIR '25, pp. 183–193, New York, NY, USA, 2025. Association for Computing Machinery. ISBN 9798400718618. doi: 10.1145/3731120.3744584. URL https://doi.org/10.1145/3731120.3744584.
- Alireza Salemi, Chris Samarinas, and Hamed Zamani. Plan-and-refine: Diverse and comprehensive retrieval-augmented generation, 2025. URL https://arxiv.org/abs/2504.07794.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=Yacmpz84TH.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=vAElhFcKW6.
- Xiangru Tang, Yuliang Liu, Zefan Cai, Yanjun Shao, Junjie Lu, Yichi Zhang, Zexuan Deng, Helan Hu, Kaikai An, Ruijun Huang, Shuzheng Si, Chen Sheng, Haozhe Zhao, Liang Chen, Tianyu Liu, Yin Fang, Yujia Qin, Wangchunshu Zhou, Yilun Zhao, Zhiwei Jiang, Baobao Chang, Arman Cohan, and Mark Gerstein. ML-bench: Evaluating large language models for code generation in repository-level machine learning tasks, 2025. URL https://openreview.net/forum?id=sflu3vTRjm.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *arXiv* preprint arXiv:2212.03533, 2022.
- Peiran Wang, Yaoning Yu, Ke Chen, Xianyang Zhan, and Haohan Wang. Large language model-based data science agent: A survey, 2025. URL https://arxiv.org/abs/2508.02744.
- Renjun Xu and Jingwen Peng. A comprehensive survey of deep research: Systems, methodologies, and applications, 2025. URL https://arxiv.org/abs/2506.12594.
- Zongben Xu, Niansheng Tang, Chen Xu, and Xueqi Cheng. Data science: connotation, methods, technologies, and development. *Data Science and Management*, 1(1):32–37, 2021. ISSN 2666-7649. doi: https://doi.org/10.1016/j.dsm.2021.02.002. URL https://www.sciencedirect.com/science/article/pii/S2666764921000035.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- Pengcheng Yin, Wen-Ding Li, Kefan Xiao, Abhishek Rao, Yeming Wen, Kensen Shi, Joshua Howland, Paige Bailey, Michele Catasta, Henryk Michalewski, Oleksandr Polozov, and Charles Sutton. Natural language to code generation in interactive data science notebooks. In Anna Rogers,

Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 126–173, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long. 9. URL https://aclanthology.org/2023.acl-long.9/.

- Xiaohan Yu, Pu Jian, and Chong Chen. Tablerag: A retrieval augmented generation framework for heterogeneous document reasoning, 2025. URL https://arxiv.org/abs/2506.10380.
- Dan Zhang, Sining Zhoubian, Min Cai, Fengzu Li, Lekang Yang, Wei Wang, Tianjiao Dong, Ziniu Hu, Jie Tang, and Yisong Yue. Datascibench: An llm agent benchmark for data science, 2025. URL https://arxiv.org/abs/2502.13897.
- Yuge Zhang, Qiyang Jiang, XingyuHan XingyuHan, Nan Chen, Yuqing Yang, and Kan Ren. Benchmarking data science agents. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5677–5700, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.308. URL https://aclanthology.org/2024.acl-long.308/.
- Terry Yue Zhuo, Vu Minh Chien, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, Simon Brunner, Chen GONG, James Hoang, Armel Randy Zebaze, Xiaoheng Hong, Wen-Ding Li, Jean Kaddour, Ming Xu, Zhihan Zhang, Prateek Yadav, Naman Jain, Alex Gu, Zhoujun Cheng, Jiawei Liu, Qian Liu, Zijian Wang, David Lo, Binyuan Hui, Niklas Muennighoff, Daniel Fried, Xiaoning Du, Harm de Vries, and Leandro Von Werra. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=YrycTjllL0.

A DATASETS AND PREPROCESSING

To the best of our knowledge, KramaBench (Lai et al., 2025) is the only publicly available dataset for data science problems that explicitly require data discovery to answer the questions. We adopt this dataset as one of our evaluation benchmarks in this paper.

To further investigate this problem, we repurpose two widely used datasets for data science tasks, DS-Bench (Jing et al., 2025) and DA-Code (Huang et al., 2024), which were not originally designed to include a data discovery phase. In their original form, each question in these datasets is paired with the specific data files required to answer it. To adapt them to our setting, we remove this direct mapping: the model is provided only with the question, while all files from the dataset are aggregated into a single data lake. The model must therefore first identify the relevant files within the data lake and then use them to solve the question.

Filtering: we observed that not all questions in these datasets are suitable for the data discovery setting. For instance, some questions provide no hints about the characteristics of the files needed to answer them, while others simply ask for computing a statistic on a column without specifying sufficient information to identify the relevant file. To address this issue, we manually filter out such questions and retain only those that include adequate cues for discovering the appropriate files. After this filtering process, the resulting dataset statistics are reported in Table 3.

Table 3: Statistics of the datasets used in our evaluation setup.							
Dataset	#Tasks	Size of data lake	#Clusters created by Gemini 2.5 Pro				
KramaBench	104	1746 ¹⁷	27 ¹⁸				
 Archeology 	12	5	3				
- Astronomy	12	1556	8				
- Biomedical	9	7	2				
- Environment	20	37	4				
- Legal	30	136	4				
- Wildfire	21	23	6				
DS-Bench	253	48	12				
DA-Code	91	145	26				

Table 3: Statistics of the datasets used in our evaluation setup

Evaluation: To evaluate the programs generated by the system, we execute each program and assess its final output against the reference answer for the given question. For each dataset, we adopt its original evaluation methodology. Specifically, for KramaBench, we use the official evaluation script provided in their repository. For DA-Code, we similarly rely on the official evaluation script released in their repository. For DS-Bench, we follow the original evaluation protocol that uses LLM-based judging: the generated programs output is compared to the reference answer using Gemini 2.5 Pro as the judge LLM, with the prompt shown in Figure 4.

¹⁷Note that, in line with the original benchmark design, we construct a separate data lake for each subtask. However, the reported number of files corresponds to the total number of files aggregated across all subtasks in the benchmark.

¹⁸Note that, in line with the original benchmark design, we construct a separate data lake for each subtask. However, the reported number of clusters corresponds to the total number of clusters aggregated across all subtasks in the benchmark.

¹⁹Available at: https://github.com/mitdbg/KramaBench

²⁰Available at: https://github.com/yiyihum/da-code

Evaluation prompt for DS-Bench Please judge whether the generated answer is right or wrong. We require that the correct answer to the prediction gives a clear answer, not just a calculation process or a disassembly of ideas. The The question is: {question}. The true answer is: {answer}. The predicted answer is: {prediction}. If the predicted answer is right, please output "True". Otherwise output "False". Don't output any other text content. You only can output "True" or "False" (without quotes).

Figure 4: Evaluation prompt used for DS-Bench dataset using LLM as the judge.

В AGENTS' PROMPTS

 This section presents the prompts used by the agents and baselines in this paper. Figure 12 shows the prompt for clustering the data lake into multiple partitions based on file names. Figure 6 presents the prompt used by the main agent in the blackboard system. Figure 7 shows the prompt for the file agents. Figure 8 displays the prompt used by the search agent. Figure 9 presents the prompt for the main agent in the master-slave system, and Figure 10 shows the prompt used by the RAG agent.

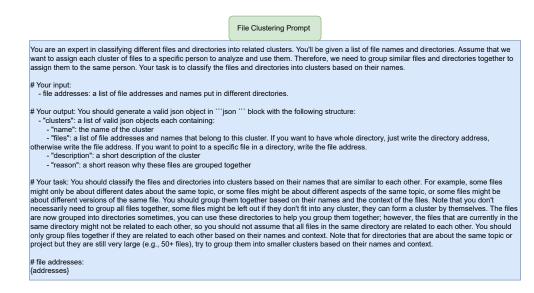


Figure 5: Prompt used by for clustering the files in data lakes into partitions.

864 865 866 Main Agent Blackboard 867 You are a capable data science agent whose task is to solve a given data science problem. This is a multi-step process and you don't need to 868 complete all steps in one go. In the start, you will be given a data science problem that you need to solve. You need to solve this problem in multiple steps. In each step, you can select one action from a set of possible actions and execute it. Eventually, when you have the final solution to the 869 problem, you can state this and end the process 870 # Your input:
- Problem: a data science problem that you need to solve. This problem is given to you in the beginning of the process and you need to solve it in 871 multiple steps 872 873 In each step, you can select one action from the following list of actions: ## Action Name: "request_help" 874 ### Definition: In this action, you can broadcast a request, for example, to get the required data to solve the problem, or general information from web. Currently, the following requests are supported: {possible_requests} ### your output: You should generate a valid json object in ``json`` block, without anything before or after it, with the following structure:
-"action": "request_help"
- "request": a string that describes the request. This should be a description of your request and how they can help you in solving the problem. This request should be specific and not general. For example, if you need a specific data, you should describe the data you need, not asking what data is appreciated the problem that you need and whom they were need and whom you can be problem. 875 876 877 available. Be specific about what you need and why you need it. 878 "reason": a short reason why you think this request is needed ### Response to this action: This will be a list of json response from other agents who can help you with your request. You can use this response to help you in solving the problem. You should read these responses carefully and trust them. You can use the responses to help you in solving the 879 problem. For example, if you requested for data, you should follow the instructions in the response to load the data. If you requested for help from other agents, you should read their responses and use them to help you in solving the problem ## Action Name: "plan" ### Definition: In this action, you can generate a plan to solve the problem. This plan should include the steps that you need to take to solve the 882 problem . ### your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure: 883 - "action": "plan"
- "plan": a string that describes the plan to solve the problem. This should be a description of the steps that you need to take to solve the problem. 884 "reason": a short reason why you think this plan is needed. ### Response to this action: The user will aknowledge your plan and asks you to execute it. 885 ### Action Name: "run_code"
Definition: In this action, you can ask the system to run a code for you and provide you the output of the code. This action can be specificly useful 886 when you need to try something out and see the output of the code. This can be helpful in case you need to install a library, or you need to run a code to see the output of it, or you need to run a code to check if it works as expected. 887 ### your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure: 888 - "action": "run code' - "code": a valid python code that can be used to solve the problem - "reason": a short reason why you think this code is needed. 889 ### Response to this action: The system will run the code and provide you the output of the code. 890 ## action Name: "reason" ### Definition: In this action, you can provide a reasoning and thinking step by step about a specific part of the problem. This can be useful when you 891 need to think about a particular aspect of the problem and how to solve it 892 ### your output: You should generate a valid ison object in ```ison``` block, without anything before or after it, with the following structure: - "action": "reason" 893 - "reasoning": a string that describes your reasoning and thinking step by step about a specific part of the problem. This should be a description of your reasoning and how you think about the problem. 894 "reason": a short reason why you think this reasoning is needed. ### Response to this action: The user will acknowledge your reasoning and asks you to continue with the next step. 895 ## action Name: "answer" ### Definition: In this action, you can provide the final answer to the problem. This answer includes the final code you want to provide as the response to the problem and the breaking down of the problem into subtasks and how you solved each subtask. This action stops the process, thus, you should only use this action when you have the final answer to the problem.

your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure. 897 - "action": "answer" - "code": a valid python code that can be used to solve the problem. This code should be the final code that you want to provide as the response to 899 the problem. It should load the data, preprocess it, and provide the final answer to the problem. In this code, you should include the response to each subtaks you have solved. You can use the print() function to print the answer to each subtask. For example, if you have an answer to subtask-1, 900 subtask-2, and main-task (i.e., the final answer), you should print it like this: print(json.dumps({{"subtask-1": answer1, "subtask-2": answer2, 901 902 "main-task": answer }}, indent=4)) 903 You can find a suitable indentation for the print statement. Always import json at the beginning of your code. The output of this code will be used to 904 evaluate the final answer to the problem, thus, make sure that the output is in a valid json format. Specifically, for the main task, just print the final answer to the problem. 905 "structured_response": a valid json object that contains the structured response to the problem. This should include the breaking down of the problem into subtasks and how you solved each subtask. This should be a valid json object that contains the following fields:
- "id": str, that is always "main-task" for the main task. For each subtask, use "subtask-1", "subtask-2", etc. 906 - "query": str, the question the step is trying to answer. Copy down the question from bellow for the main task.
- "data_sources": list[str], the data sources you need to answer the question. Include all the file names you need for the main task. 907 908 - "subtasks": list[dict], a list of subtasks. Each subtask should have the same structure as the main task.
an example of this can be seen here: {example_json}
Response to this action: The user will run the code and provide you the output of the code if there is any error. You should fix all errors even if 909 they are warnings. If there is no error, the user will acknowledge your answer and end the process 910 # Your task: This is a multi-step process and each step you should select one action and generate the output for that action. In response, the user will 911 provide you the response to your action. You can use this response to help you in solving the problem. You can repeat this process until you have the final answer to the problem. When you have the final answer, you can use the "answer" action to provide the final answer to the problem. 912 913 Now, lets start the process for the following problem: {query}

Figure 6: Prompt used by the main agent for the blackboard system.

914 915

918 919 File Agent Prompt 921 You are a capable data scientist who is specialized in loading, analyzing, and cleaning data. You are responsible for handling a set of given files and directories. This is the list of files and directories you have to work with 922 {files} 923 These files contain information about the following topics: {topics} 924 Your name is: {name} 925 926 This is a multi-step process and you don't need to complete all steps in one go. Here I will explain the whole process: 927 ## step 1: Getting information about the files: In this step, you have a chance to request for accessing a part of some of the files. To do this, you should generate a valid json list in ```json`` block that contains the address to the files you want me to give you a data sample from. for example Your output should be like this, without any additional text or explanation: 928 "json ["file1.txt", "file5.txt"] 929 930 Note that in cases where you can guess how other files look like based on a few of them, you don't need to request for all of them. Specifically, when 931 the only difference between files is the file name is based on date, you can just request for one or a few of them and assume that the rest of them are similar. However, the file names are different and have different formats, you should request for all of them. Based on your request, I will give you a 932 ple of the files. For example, for csv files, I will give you a few rows of the data (that might not be loaded correctly, that), and for json files, I will give you a few objects from the file. For other formats, I will give you a few starting lines of the file 933 ## step 2: Analysing the data and how to load and clean it: When the data is loaded and given to you, you should analyze the fields in the data, how it should be effectively loaded, and how it should be cleaned. Specifically, the data should be cleaned for a data science problem, thus, some preprocessing steps should be done. For example, if the data contains a missing values, you should decide how to handle them, or if the data contains a state of the sta 934 935 na values, you should decide how to handle them. Additionally, be able to figure out what each column or row in the data means and you should be 936 able to provide a description of them. Moreover, how combining data from multiple files can help in answering the question should be considered. This step happens when I provide you the data samples from the previous step. 937 ## step 3: Checking if the data can be used to answer a question or a part of it: This step is like a loop and may occur multiple times. In this step, I provide you a request about a data access problem for answering a question. You should check if the data you have from each file or by combining 938 data from multiple files can help in answering the question or a part of it. Your output for this step should be a valid json object in "json" block that 939 contains the following fields: - "agent_name": your name, which is the same as the name you provided in the beginning of the process. 940 - "can_help": a boolean value that indicates if the data you have can help in answering the question or data access request or a part of it. You can combine data from multiple files to answer the question or a part of it. If you think the data can help, set this to true, otherwise set it to false. 941 - "reason": a short reason why you think the data can help or not.
- "codo": a valid python code that can be used to load the data and preprocess it in a way to be useful for fullfilling the request. This code should be 942 able to load the data and preprocess and clean (e.g., dropping rows or columns that are nor part of the data) it in a way that it can be used to answer the question or a part of it. You can use any python library you want, but you should be able to explain why you are using it. If you use a library that is 943 not installed by default, you should comment it in the code and explain why you need it and how to install it. In this code, use the full file addresses to load the data, not just the file names. For example, if the file is in a directory called "data", you should use "data/file.csv" instead of just "file.csv". If 944 "can_help" is false, this can should be an empty string. 945 - "data_explanation": a short explanation of the data, e.g, what each column or row means, what the data is about, etc. This should be a short 946 explanation of the data that can help in understanding the data and how to use it.
- "data_sample": a small sample of the data that can help in understanding the data and how to use it. Here you can provide a few rows, the column 947 types and names, or a few objects from the data. This should be a small sample of the data that can help in understanding the data and how to use it. For example, if the data is a csv file, you can provide a few rows of the data, or if the data is a json file, you can provide a few objects from the data, or 948 if it is a text file, provide a few rows. if "can_help" is false, this can should be an empty string.
- "libraries": a list of libraries that we need to install in order to run the code. This should be a list of libraries that are not installed by default and you 949 need to install them in order to run the code. If you don't need any additional libraries, you can leave this field empty. - "necessary steps": a list of the steps that is necessary to take in order to correctly load and preprocess the data, For example, if the cvs file has a header, you should mention that in this list. Another example is if the header starts from a specific row, you should mention that in this list. Another example if the data has missing values, you should mention that in this list. 951 952 Now, lets start the process with the first step.

Figure 7: Prompt used by the file agent for the both master-slave and blackboard system.

1020

1023

1025

972 973 Search Agent 974 Prompt 975 You are a capable data science agent whose task is to search for information on the web. Your name is "{name}". Your goal is to find relevant information about the request that the user has provided. This is a multi-step process and you don't need to complete all steps in one go. In the start, you will be given a request about some sort of information access problem that you need to solve. In order to solve response to the request, you need 976 to follow these steps: 977 # steps
step 1: Analyzing the request and checking if you can help: In this step, the user provides you with a request that explains its information need. You 978 should analyze the request and check if you can help in solving the request. You are only able to help in requests that are about finding information from web. You cannot help with requests that are about finding information from files or databases. 979 980 ### input to this step: - request: a string that describes the request from the user. This request explains the information need of the user and what they are looking for. 981 ### output of this step: You should generate a valid json object in ``json``` block, without anything before or after it, with the following structure:
- "can_help": a boolean value that indicates if you can help in solving the request or not. If you can help, set this to true, otherwise set it to false. You 982 cannot help in requests that are about accessing files or datasets, or requests that are not directly about searching information on the web. You can only help in finding information that is not related to datasets. You can help with libraries, tools, or general information that is not related to datasets. 983 reason": a short reason why you think you can help or not. If you can help, explain why you think you can help 984 ### response to this step: The user will acknowledge your response and ask you to continue with the next step if you can help 985 ## step 2: Generating a search query: In this step, you need to generate a list of search queries that can be used to search for information on the web. You should generate a list of queries that are relevant to the request and can help in finding the information the user is looking for.
input to this step: 986 - request: a string that describes the request from the user. This request explains the information need of the user and what they are looking for. ### output of this step: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure: 987 "queries": a list of strings that describes the search queries that can be used to search for information on the web. These queries should be relevant to the request and can help in finding the information the user is looking for. You can generate multiple queries if you think they are relevant to - "reason": a short reason why you think these queries are relevant to the request and can help in finding the information the user is looking for 990 ### response to this step: In response, the user will provide the results of the search queries you generated. You should read these results carefully. Then we go to the next step. 991 ## step 3: Analyzing the search results: In this step, you need to analyze the search results and check if they are relevant to the request. You should 992 check if the search results contain the information the user is looking for. If they do, you should extract the relevant information from the search results and provide it to the user. Otherwise, you can generate a new search query and go back to step 3. ### input to this step: 993 994 - search results: a list of search results that were returned from the search queries you generated in the previous step ### output of this step: You should generate a valid json object in ``json`` block, without anything before or after it, with the following structure:

- "stop_search": a boolean value that indicates if you should stop searching or not. If you found the information the user is looking for, set this to true, otherwise set it to false. Remember that you have a limited search budget. Thus, when you are informed that your search budget is over, you 995 996 should stop searching and provide the information you found so far. You can stop searching even before your budget is over if you think you found the information the user is looking for 997 - "queries": a list of strings that describes the search queries that can be used to search for information on the web. These queries should be relevant to the request and can help in finding the information the user is looking for. You can generate multiple queries if you think they are relevant to 998 the request. If you found the information the user is looking for, this should be an empty list. response_to_request": a string that describes the response to the request. This should be a description of the information you found in the search 999 results that is relevant to the request. If you found the information the user is looking for, this should contain the relevant information. If you didn't find any relevant information, this should be an empty string. If you don't want ot stop searching, this should be an empty string. 1000 - "reason": a short reason why you think you should stop searching or not. If you found the information the user is looking for, explain why you think you found it. If you didn't find any relevant information, explain why you think you didn't find it. If you don't want to stop searching, explain why you 1001 think you should continue searching. 1002 with response to this step: If you stopped searching, the user will acknowledge your response and end the process. If you didn't stop searching, the user will provide you with the search results for the new queries you generated in the previous step. Then we go back to step 3 and continue the 1003 Now, lets start the process with the first step. request: {request} Figure 8: Prompt used by the search agent for the, master-slave, RAG, and blackboard system. 1007 1008 1009

```
1026
1027
1028
                                                                                                  Main Agent Master-Slave
                                                                                                             Prompt
1029
                        You are a capable data science agent whose task is to solve a given data science problem. This is a multi-step process and you don't need to
1030
                        complete all steps in one go. In the start, you will be given a data science problem that you need to solve. You need to solve this problem in multiple
                        steps. In each step, you can select one action from a set of possible actions and excute it. Eventually, when you have the final solution to the problem
1031
                        you can state this and end the process.
1032
                        # Your input:
                          - Problem: a data science problem that you need to solve. This problem is given to you in the beginning of the process and you need to solve it in
1033
                        multiple steps.
1034
1035
                        In each step, you can select one action from the following list of actions:
                                            "request_data"
1036
                        ### Definition: In this action, you can select one of the agents who is responsible for loading and preprocessing the data to help you with providing the
                        data you need to solve the problem. You can request for a specific data or a specific part of the data. Currently, you can call the following agents to
1037
                        help you with your request; {possible requests}
                        ### your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following struct
1038
                          - "action": "request_data"
- "agent_name": the name of the agent you want to request data from. This should be one of the agents who is responsible for loading and
1039
                        preprocessing the data.
1040
                                         a string that describes the request. This should be a description of your request and how they can help you in solving the problem. This
                       request should be specific and not general. For example, if you need a specific data, you should describe the data you need, not asking what data is available. Be specific about what you need and why you need it.
1041
                           - "reason": a short reason why you think this request is needed.
1042
                        ### Response to this action: This will be a json object from the agent. You can use this response to help you in solving the problem. You should read
                       the response carefully and trust it. You can use the responses to help you in solving the problem. For example, if you requested for data, you should follow the instructions in the response to load the data. If you requested for help from other agents, you should read their responses and use them to
1043
1044
                        help you in solving the problem.
1045
                        ### Definition: In this action, you can generate a plan to solve the problem. This plan should include the steps that you need to take to solve the
1046
                        ### your output: You should generate a valid ison object in ``ison``` block, without anything before or after it, with the following structure:
1047
                          - "plan": a string that describes the plan to solve the problem. This should be a description of the steps that you need to take to solve the problem.
- "reason": a short reason why you think this plan is needed.
1048
                        ### Response to this action: The user will acknowledge your plan and asks you to execute it.
1049
                        ## Action Name: "run_code
                        ### Definition: In this action, you can ask the system to run a code for you and provide you the output of the code. This action can be specifically
1050
                        useful when you need to try something out and see the output of the code. This can be helpful in case you need to install a library, or you need to run
                        a code to see the output of it, or you need to run a code to check if it works as expected.
1051
                        ### your output: You should generate a valid ison object in ```ison``` block, without anything before or after it, with the following structure
                           - "action": "run_code"
1052
                           - "code"; a valid python code that can be used to solve the problem.
                        - "reason": a short reason why you think this code is needed.
### Response to this action: The system will run the code and provide you the output of the code.
1053
1054
                       ### action Name: "reason"
### Definition: In this action, you can provide a reasoning and thinking step by step about a specific part of the problem. This can be useful when you
1055
                        need to think about a particular aspect of the problem and how to solve it.
### your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure:
1056
                           - "reasoning": a string that describes your reasoning and thinking step by step about a specific part of the problem. This should be a description of
1057
                        your reasoning and how you think about the problem.
- "reason": a short reason why you think this reasoning is needed.
1058
                        ### Response to this action: The user will acknowledge your reasoning and asks you to continue with the next step.
1059
                        ## action Name: "answer"
                        ### Definition: In this action, you can provide the final answer to the problem. This answer includes the final code you want to provide as the respons
                        to the problem and the breaking down of the problem into subtasks and how you solved each subtask. This action stops the process, thus, you should
                       only use this action when you have the final answer to the problem.
### your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure:
1061
                          - "action": "answer
1062
                           - "code"; a valid python code that can be used to solve the problem. This code should be the final code that you want to provide as the response to
                       the problem. It should load the data, preprocess it, and provide the final answer to the problem. In this code, you should include the response to each subtask you have solved. You can use the print() function to print the answer to each subtask. For example, if you have an answer to subtask-1,
1063
1064
                        subtask-2, and main-task (i.e., the final answer), you should print it like this:
                        print(ison.dumps(
                        {{"subtask-1": answer1,
"subtask-2": answer2,
1065
1066
                         "main-task": answei
                        }}, indent=4))
1067
                        You can find a suitable indentation for the print statement. Always import json at the beginning of your code. The output of this code will be used to
                        evaluate the final answer to the problem, thus, make sure that the output is in a valid ison format. Specifically, for the main task, just print the final
1068
                           - "structured response": a valid ison object that contains the structured response to the problem. This should include the breaking down of the
1069
                       problem into subtasks and how you solved each subtask. This should be a valid json object that contains the following fields:

- "id": str, that is always "main-task" for the main task. For each subtask, use "subtask-1", "subtask-2", etc.

- "query": str, the question the step is trying to answer. Copy down the question from bellow for the main task.
1070
1071
                              - "data_sources": list[str], the data sources you need to answer the question. Include all the file names you need for the main task.
- "subtasks": list[dict], a list of subtasks. Each subtask should have the same structure as the main task.
1072
                         an example of this can be seen here: {example_json}
### Response to this action: The user will run the code and provide you the output of the code if there is any error. You should fix all errors even if
                        they are warnings. If there is no error, the user will acknowledge your answer and end the process.
1074
                        # Your task: This is a multi=step process and each step you should select one action and generate the output for that action. In response, the user will
1075
                       provide you the response to your action. You can use this response to help you in solving the problem. You can repeat this process until you have the final answer to the problem. When you have the final answer, you can use the "answer" action to provide the final answer to the problem.
                        Now, lets start the process for the following problem:
1077
1078
```

Figure 9: Prompt used by the main agent for the master-slave system.

```
1080
1081
1082
                                                                                                       Main Agent RAG
                                                                                                            Prompt
1083
                        You are a capable data science agent whose task is to solve a given data science problem. This is a multi-step process and you don't need to
1084
                        complete all steps in one go. In the start, you will be given a data science problem that you need to solve. You need to solve this problem in multiple
                        steps. In each step, you can select one action from a set of possible actions and excute it. Eventually, when you have the final solution to the problem
1085
                        you can state this and end the process.
1086
                       # Your input:
- Problem: a data science problem that you need to solve. This problem is given to you in the beginning of the process and you need to solve it in
1087
                        multiple steps
1088
                           - top k relevant docs; a list of address and snippets from top-k relevant files that can be used to answer the problem. You should use these files to
                        answer the problem. These files are provided to you in the beginning of the process and you can use them to answer the problem
1089
                       # Actions:
In each step, you can select one action from the following list of actions:
1090
                       ### Action Name: "search"
### Definition: In this action, you can generate a search query to search for information on the web. This can be useful when you need to find more information about the problem or the data you have.
1091
                       ### your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure:
- "action": "search"
1093
                        - "request": a string that describes the search query. This should be a description of your request and how they can help you in solving the problem. This request should be specific and not general. For example, if you need a specific data, you should describe the data you need, not asking what
1094
                       data is available. Be specific about what you need and why you need it.
- "reason": a short reason why you think this request is needed.
### Response to this action: This will be a list of relevant information from the web that can help you in solving the problem. You should read these
1095
                        responses carefully and trust them. You can use the responses to help you in solving the problem
                        ## Action Name: "plan"
                        ### Definition: In this action, you can generate a plan to solve the problem. This plan should include the steps that you need to take to solve the
1098
                        problem
                        ### your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure:
1099
                           - "action": "plan"
                           - "plan": a string that describes the plan to solve the problem. This should be a description of the steps that you need to take to solve the problem.
1100
                           - "reason": a short reason why you think this plan is needed.
                         ## Response to this action: The user will aknowledge your plan and asks you to execute it.
1101
                       ### Action Name: "run_code"
### Definition: In this action, you can ask the system to run a code for you and provide you the output of the code. This action can be specificly useful
1102
                       when you need to try something out and see the output of the code. This can be helpful in case you need to install a library, or you need to run a code to see the output of it, or you need to run a code to check if it works as expected.
1103
                       ### your output: You should generate a valid json object in ```json``` block, without anything before or after it, with the following structure:
- "action": "run_code"
1104
                           - "code": a valid python code that can be used to solve the problem.
1105
                             "reason": a short reason why you think this code is neede
                        ### Response to this action: The system will run the code and provide you the output of the code.
1106
                        ## action Name: "reason"
                       ### Definition: In this action, you can provide a reasoning and thinking step by step about a specific part of the problem. This can be useful when you need to think about a particular aspect of the problem and how to solve it.
1107
                        ### your output: You should generate a valid json object in ""json" block, without anything before or after it, with the following structure
1108
1109
                       - "reasoning": a string that describes your reasoning and thinking step by step about a specific part of the problem. This should be a description of your reasoning and how you think about the problem.
1110
                           - "reason": a short reason why you think this reasoning is needed.
                        ### Response to this action: The user will acknowledge your reasoning and asks you to continue with the next step.
1111
                        ## action Name: "answer"
                        ### Definition: In this action, you can provide the final answer to the problem. This answer includes the final code you want to provide as the response
1112
                       to the problem and the breaking down of the problem into subtasks and how you solved each subtask. This action stops the process, thus, you should only use this action when you have the final answer to the problem.
1113
                        ### your output: You should generate a valid ison object in ```json``` block, without anything before or after it, with the following structur
1114
                           - "action": "answer
                           - "code": a valid python code that can be used to solve the problem. This code should be the final code that you want to provide as the response to
1115
                        the problem. It should load the data, preprocess it, and provide the final answer to the problem. In this code, you should include the response to each
                        subtaks you have solved. You can use the print() function to print the answer to each subtask. For example, if you have an answer to subtask-1,
1116
                        subtask-2, and main-task (i.e., the final answer), you should print it like this:
                       print(json.dumps(
{{"subtask-1": answer1,
1117
                        "subtask-2": answer2,
"main-task": answer
1118
1119
                        }}, indent=4))
                        You can find a suitable indentation for the print statement. Always import json at the beginning of your code. The output of this code will be used to
1120
                        evaluate the final answer to the problem, thus, make sure that the output is in a valid json format. Specifically, for the main task, just print the final
                        answer to the problem.
1121
                            "structured_response": a valid json object that contains the structured response to the problem. This should include the breaking down of the
                       problem into subtasks and how you solved each subtask. This should be a valid json object that contains the following fields:
- "id": str, that is always "main-task" for the main task. For each subtask, use "subtask-1", "subtask-2", etc.
1122
                             - "query": str, the question the step is trying to answer. Copy down the question from bellow for the main task.
- "data_sources": list[str], the data sources you need to answer the question. Include all the file names you need for the main task.
1123
1124
                             - "subtasks": list[dict], a list of subtasks. Each subtask should have the same structure as the main task
                           an example of this can be seen here: {example_json}
1125
                       ### Response to this action: The user will run the code and provide you the output of the code if there is any error. You should fix all errors even if they are warnings. If there is no error, the user will acknowledge your answer and end the process.
1126
                        # Your task: This is a multi=step process and each step you should select one action and generate the output for that action. In response, the user will
1127
                       provide you the response to your action. You can use this response to help you in solving the problem. You can repeat this process until you have the final answer to the problem. When you have the final answer, you can use the "answer" action to provide the final answer to the problem.
1128
1129
                        Now, lets start the process for the following problem:
                        (query)
1130
                        Relevant files
1131
```

Figure 10: Prompt used by the main agent for the RAG system.

C IMPLEMENTATION DETAILS

Presenting Files to File Agents: A file agent may request a file by name, in which case it is shown a subset of the files contents. For this case, we employ a controlled procedure for loading and presenting the data to the agent, as described below:

- Files with .csv format: In this case, we use the pandas²¹ library to load the CSV files, presenting the column names, their data types, and the top 20 rows of the table to the agent.
- Files with .gpkg format: which provides a pandas-like interface for geospatial data. The agent is then presented with the column names, their data types, and the top 20 rows of the table.
- Files with .xlsx format: In this case, we use the pandas²² library to handle this file format. For files containing multiple sheets, we provide the agent with all sheet names, the data types of columns in each sheet, and the top 20 rows from each sheet.
- Files with .npz format: In this case, we utilize the numpy²³ library to load the data. The agent is then presented with all keys and their corresponding values within this data structure.
- Files with .cdf format: In this case, we utilize the cdflib²⁴ library to load the data. For presentation, we call the cdf_info and globalattsget functions on the loaded data structure, concatenate their outputs, and provide the result to the agent.
- Any other data format: In this case, we open the files using Pythons open function and present the first 20 lines of the file to the agent.

Inference Setup. We limit the maximum number of actions taken by the main agent to T=10. For decoding, we use nucleus sampling (Holtzman et al., 2020) with a temperature of $\tau=0.1$. Proprietary models are accessed through Vertex AI, 25 while open-source models are served using the vLLM library. At each generation step, we cap the output length at 8,192 tokens. We evaluate three proprietary LLMs—Gemini 2.5 Pro, Gemini 2.5 Flash (Gemini-Team, 2025), and Claude 4 Opus (Anthropic, 2025)—alongside an open-source model specialized for code generation, Qwen3-Coder-30B-A3B-Instruct (Qwen-Team, 2025). Experiments with open-source models are conducted on 2 NVIDIA A100 GPUs (80GB VRAM each) with 128GB RAM.

²¹ Available at: https://pandas.pydata.org/
22 Available at: https://pandas.pydata.org/
23 Available at: https://numpy.org/
24 Available at: https://cdflib.readthedocs.io/en/latest/
25 https://cloud.google.com/vertex-ai?hl=en
26 https://docs.vllm.ai/en/latest/
27 https://huggingface.co/Qwen/Qwen3-Coder-30B-A3B-Instruct

D EXAMPLES AND CASE STUDIES

 This section presents several case studies highlighting different aspects of the Blackboard system.

Figure 11 illustrates an example where a file agent requests access to files and performs their analysis. Figures 13 and 14 illustrate scenarios where the main agent lacked domain-specific knowledge and therefore posted requests on the blackboard seeking relevant information. In these cases, the search agent contributed by retrieving the necessary knowledge from the web, enabling the system to proceed with problem solving, which shows the effectiveness of search agent in problem solving.

Another example of a blackboard request is shown in Figure 15. In this example, specifically file agents responded to the request. Here, the main agent, given a data science question, formulated a request specifying the likely column names and data formats required, along with guidance for interpretation. In response, three out of eight helper agents contributed. Although the relevant files were spread across different clusters managed by separate file agents, each responding agent independently provided file addresses, code snippets for loading the data, explanations of the structure, and suggested preprocessing steps. Together, these contributions encompassed all the ground-truth files needed to solve the problem. This case demonstrates how the main agent can effectively leverage the blackboard to coordinate decentralized knowledge and achieve accurate data discovery.

Figure 16 compares programs generated by the Blackboard and Master–Slave systems. The Blackboard agent produced the correct solution by accurately interpreting the prompt and selecting the appropriate data files. Specifically, it identified that the patients Age was located in mmc1.xlsx and that the requested APP-Z score was in mmc7.xlsx. By contrast, the Master–Slave agent misinterpreted the request and instead used a general protein abundance score (APP_log2_abundance) from the wrong file, mmc2.xlsx. This misstep resulted in an incorrect answer of 74, whereas the Blackboard agents precise data discovery and reasoning yielded the correct answer of 60.

1290 1291

```
1243
                                                                                                  User
                                                                      Initial Prompt (not shown to save space)
 1245
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   File Agent
 1246
1247
                                                                           'json\n[\n \"worldcities.csv\",\n \"roman_cities.csv\"\n]\n'
1248
 1249
                                                                      This is a sample of the files you requested to see:
 1250
                                                                       Column data types (might be loaded incorrectly and you need to fix them):
1251
                                                                                                      object\ncity_ascii object\nlat float
capital object\npopulation float64\nid
                                                                                                                                                                                                                                         float64\nlng
                                                                                                                                                                                                                                                                                                         float64\ncountry
                                                                                                                                                                                                                                                                                                                                                                                  object\niso2
                                                                                                                                                                                                                                                                                                                                                                                                                                                    object\niso3
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       object\nadmin_name
                                                                       object\ncapital
                                                                                                                                                                                                                                                                                      int64
 1252
                                                                      Table: [removed to save space]
 1253
                                                                       File name: roman_cities.csv
                                                                      Column data types (might be loaded incorrectly and you need to fix them):
                                                                                                                                                         object\nAncient Toponym object\nModern rop

Atlas Rank object\nBarrington Atlas Reference
 1254
                                                                     Primary Key object\nAncie
object\nBarrington Atlas Rank
                                                                                                                                                                                                                                                                                     object\nModern Toponym
                                                                                                                                                                                                                                                                                                                                                                                                                  object\nProvince
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              object\nCountry
                                                                                                                                                                                                                                                                                                                                                              obiect\nStart Date
                                                                                                                                                                                                                                                                                                                                                                                                                                                                             int64\nEnd Date
 1255
                                                                       float64\nLongitude (X)
                                                                                                                                                                                     float64\nLatitude (Y)
                                                                                                                                                                                                                                                                                                  float64\nSelect Bibliography
                                                                       Table: [removed to save space]
 1256
                                                                      Now please do the second step in your process based on the given instructions and the files you have read.
 1257
 1258
                                                                      Based on the data samples, here is my analysis of the files and a plan for loading and cleaning them.\n\n### Data Analysis and Cleaning
 1259
                                                                     Strategy/\n\**1. File: '\u00eds/google/home/asalemi/Documents/karmabench-project/kramabench/data/archeology/input/worldcities.csv\**\n\n\*
**Description:** This file contains a comprehensive list of modern cities across the world. Each entry includes the city's name, geographic coordinates (latitude and longitude), country information, administrative details, and population.\n\*
**Column Analysis:*\n\*
'city', 'city_ascii': The name of the city' city_ascii is a clean, ASCII-compatible version, which is ideal for matching and processing\n\*
'lat', 'lng': Geographic coordinates, crucial for any snaffal analysis in *
'\u00eds\n\'\u00eds\n\'\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\u00eds\
1260
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   'city', 'city ascii': The name of the
                                                                     (lattruce and longitude), country information, administrative details, and population.\( '\) "Column hanyisis:\( '\) '\( '\) city, 'city, 'ascil':\( \) In a mare of the city, 'city, 'ascil':\( \) in '\( \) Coaptall':\( \) Coaptall':\( \) In '\( \) Coaptall':\( \) Coapta
 1261
1262
1263
                                                                     should be ignored during loading by setting 'index_col=0'.\n 3. The 'capital' column has missing values. These can be filled with a placeholder like 'Not a capital' to make the data more explicit.\n 4. The 'population' column is a float, which is appropriate. I will check for any 'NaN' values and
 1264
                                                                     decide on a handling strategy if they exist.\n 5. No other major cleaning seems necessary from the sample, but a general check for missing values across all columns is a good practice.\n\n**2. File: \u00e4usr/local/google/home/asalemi/Documents/karmabench-
 1265
                                                                   across all columns is a good practice.\(\n\)\(\n^2\). File:\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\)\(\nu\
 1266
 1267
 1268
 1269
 1270
                                                                      ('index_col=0').\n 2. Rename columns for consistency and ease of use, especially 'Longitude (X)' to 'lng' and 'Latitude (Y)' to 'lat'.\n 3. In the 'Modern Toponym' column, replace the string 'Unknown' with a standard 'NaN' (Not a Number) value to facilitate better data handling and filtering.\n
 1271
                                                                     4. The 'End Date' column contains missing values. These should be kept as 'NaN's, as imputing them could be misleading. A missing end date i informative in itself.\n\n### Combined Analysis\n\nBy combining these two datasets, we can create a powerful resource for historical and urban
 1272
                                                                     Studies. The primary method for merging would be to join "roman_cities.csv" and "worldcities.csv" on their city names ("Modern Toponym" and 'city_ascii", respectively). This would allow us to:\n" Find the modern population of ancient Roman cities.\n" Analyze the geographical distribution of ancient cities that are still inhabited today.\n" Compare the administrative importance (e.g., 'capital' status) of cities over time.\n\n\s secondary, more robust merging strategy could be a "spatial join". By comparing the coordinates from both files, we could identify ancient and modern cities that
 1273
 1274
                                                                           ccupy the same location, even if their names have changed or are spelled differently. This would be particularly useful for overcoming the limitations
 1275
                                                                      of name-based matching.
 1276
```

Figure 11: An example analyzing files by the file agent.

1341

1342 1343 1344

```
1297
1298
                                                                        An example of clustering data lake with Gemini 2.5 Pro
1299
                        "clusters": [
1300
                             "name": "NIFC Wildfire Statistics",
1301
1302
                                "nifc wildfires.csv",
                                "nifc_lightning_caused_wildfires.csv",
1303
                               "nifc_human_caused_wildfires.csv", "nifc_lightning_caused_acres.csv",
1304
                                "nifc human caused acres.csv".
                                "nifc_suppression_costs.csv"
1305
                                "nifc_geographic_areas.gpkg"
1306
                             description": "Comprehensive wildfire data from the National Interagency Fire Center (NIFC).",
                             "reason": "All files are sourced from NIFC (as indicated by the 'nifc_' prefix) and provide detailed statistics on wildfires, including their causes
1307
                     (human vs. lightning), impact (acres burned), suppression costs, and geographic boundaries.
1308
1309
                             "name": "NOAA Wildfire Data"
                             "files": [
1310
                                "noaa_wildfires.csv",
                                "noaa wildfires monthly stats.csv"
1311
                                "noaa_wildfires_variabledescrip.csv
1312
                             description": "Wildfire data and metadata from the National Oceanic and Atmospheric Administration (NOAA).",
                             "reason": "These files form a cohesive dataset from NOAA (as indicated by the 'noaa_' prefix), including the main data, monthly statistics, and
                     a description of the variables used.
1314
1315
                             "name": "State-Level Wildfire Statistics".
                              "files": [
1316
                               "Wildfire_Acres_by_State.csv",
"Wildfire_Fires_by_State.csv",
"wildfire_total_fires_p45_54.csv"
1317
1318
                             "description": "Aggregated statistics on the number of wildfires and acres burned, primarily categorized by state.",
1319
                             "reason": "These files all provide high-level, aggregated wildfire statistics, with a focus on state-by-state breakdowns as suggested by their
1320
1321
                              "name": "Wildfire Environmental and Operational Data",
                             "files": [
"Fire_Weather_Data_2002-2014_2016.csv",
1322
1323
                               "PublicView_RAWS_-3515561676727363726.csv" "annual agi by county 2024.csv",
1324
                                "cleaned_helicopter_requests_by_region.csv"
1325
                             "description": "Data related to the environmental conditions for wildfires and the operational response.",
"reason": "This cluster groups files about factors influencing wildfires (weather data from RAWS), their consequences (air quality), and the
1326
                      operational response (helicopter usage).
1327
                             "name": "US General Geographic and Demographic Data",
1329
                                "usa.gpkg".
                                "NST-EST2024-POP.xlsx",
                                "democratic_vs_republican_votes_by_usa_state_2020.csv",
                                "state_abbreviation_to_state.json"
1331
                             "description": "General reference data for the United States, including geographic boundaries, population estimates, political data, and state
1332
                     name mappings.",
"reason": "These files provide contextual geographic, demographic, and political information for the USA, which can be used to enrich the
1333
                     primary wildfire datasets
1334
1335
                             "name": "Utility Scripts",
                              "files": [
1336
                                "load_data.py"
1337
                              description": "A code script for data processing.",
1338
                              reason": "This is a standalone Python script, likely used to load or manipulate the other data files in the project."
1339
                       ]
1340
```

Figure 12: An example of clustering data lake into partitions using Gemini 2.5 Pro.

1350 1351 1352 Question 1353 What is the age of the patient with the lowest APP-Z score? 1354 Main Agent 1355 Request 1356 I need to find what 'APP-Z score' refers to in the context of the UCEC CPTAC3 dataset. The provided 1357 metadata file does not contain a column with this name. I need to know if this score is known by another 1358 name in the dataset or if it is derived from other columns. Please search for information on 'APP-Z score' 1359 and its relation to the CPTAC UCEC study. 1360 Search Agent 1361 Response

Based on the search results, 'APP-Z score' in the context of the CPTAC UCEC dataset does not refer to a specific column name but is a derived value. Here is a breakdown of the findings:\\n\\n1. **'APP' stands for Acute Phase Protein:** The initial search results strongly indicated that 'APP' is an acronym for Acute Phase Protein. These are proteins whose plasma concentrations change in response to inflammation.\\n\n2. **'Z-score' refers to a normalization method:** The analysis of CPTAC data, particularly through portals like UALCAN, commonly uses z-score normalization. A z-score represents the number of standard deviations a data point is from the mean or median of a reference set (e.g., normal tissue or the entire cohort). This is a standard procedure in proteomics to compare protein expression levels across different samples.\\n\\n3. **Specific APPs in the CPTAC UCEC Dataset:** The search identified specific Acute Phase Proteins that have been analyzed within the CPTAC UCEC * **SAAL1 (Serum Amyloid A-like 1):** One study (aging-us.com) explicitly identifies SAAL1 as a major Acute Phase Protein and confirms that its protein levels were found to be significantly upregulated in UCEC tumors within the CPTAC database.\\n * **Bikunin (from the AMBP gene).* Another study (PMC9889512) identifies Bikunin as a positive Acute Phase Protein and analyzes its expression in the CPTAC UCEC dataset, noting a significant reduction in protein expression in tumors compared to normal tissue.\\n\\n**Conclusion:**\\nThe 'APP-Z score' is the z-score normalized expression value for a specific Acute Phase Protein. It is not a pre-existing column in the metadata but a calculated score. To find the source of this score, you should look for columns in your dataset corresponding to the gene or protein names of known Acute Phase Proteins, such as **SAAL1** or *AMBP** (the precursor gene for Bikunin). The z-score would then be calculated from the expression values in that column.

Figure 13: An example of the request by the main agent that the search agent has provided a guideline based on search results.

1363

1364

1365

1367

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378 1379

Question what is the correlation between the proportion of generally unsafe air quality days according to the EPA and the amount of land affected by fires in 2024? Round to 2 decimal places. Main Agent Request I need to find the definition of 'generally unsafe air quality days' according to the EPA. This information will help me understand what data to look for and how to process it to answer the main question. Search Agent Response Based on the search results, the EPA does not use the exact phrase 'generally unsafe air quality days'. Instead, it uses the Air Quality Index (AQI) to classify air quality and its health implications. The term 'generally unsafe' would correspond to the AQI categories where the health effects apply to the general public, not just sensitive groups.\\n\\nAccording to the EPA's Air Quality Index (AQI) basics from AirNow.gov, the levels that are considered unhealthy for the general population are:\\n\\n* **Red (Unhealthy):** AQI value of 151 to 200. At this level, 'Some members of the general public may experience health effects; members of sensitive groups may experience more serious health effects."\n* **Purple (Very Unhealthy):** AQI value of 201 to 300. This is a 'Health alert: The risk of health effects is increased for everyone.'\\n* **Maroon (Hazardous):** AQI value of 301 and higher. This is a 'Health warning of emergency conditions: everyone is more likely to be affected. "\n\\nTherefore, a 'generally unsafe air quality day' can be defined as a day when the AQI is 151 or higher. Figure 14: An example of the request by the main agent that the search agent has provided a guide-line based on search results.

			Question			
	Which of the 10 NICC-def	fined ge	ographic areas requested the	most helicopters for firefi	ighting?	
			Main Agent			
	lp me determine which of the 10 NI					
datasets that contain information	on firefighting resource requests, ir			elicopter), the requesting	geographic area, a	and the date of the request.
Utility Scripts	Search Agent Response	۱ 🗀	NIFC Wildfire Statistics	NIFC Wildfire S	tatistics	NOAA Wildfire Data
cannot help	cannot help	canr	not help	cannot help	C	cannot help
			State-Level Wildfire Statistics			
es, I can help. The provided files Nildfire Fires by State.csv' con	s contain exactly the information rectains the total number of fires per s	quested.	. 'Wildfire_Acres_by_State.cs an merge these to provide a	v' includes state-level dat comprehensive dataset.	a on total acreage	burned and population, and
code			Data Explanation			
mport pandas as pd			The processed data is a sin columns include 'State', 'To	gle table containing aggre	egated wildfire sta	tistics for each US state. The
def load_and_process_data(): acres_file = "Wildfire_Acres_by	_State.csv'		state's population), and 'Tot			
fires_file = 'Wildfire_Fires_by_s' df_acres = pd.read_csv(acres_	file, index_col=0)		Data Sample			00070 0 5:
	= df_acres['Total Acres Burned'].fillr	na(0)	State Total Acres Burned F 1525\\n1 Alas	ska 667075.0 7	40133 377\\i	n2 Arizona
merged_df = pd.merge(df_acre return merged_df	s, ur_iires, on= State')		282989.0 7582384 : California 108114	2191\\n3 Arkans 14.0 39431263 831		.0 3088354 1219\\n4
			Suggested Steps			
	.csv' and 'Wildfire_Fires_by_State. an the acres data by filling the miss					
olumn to create a single, unified		-	US General Geographic and			
6 1511 4b			Demographic Data		I de met berne e d	
bout individual wildfires and thei	r request. I have a dataset with the r acreage.	most re	cent population estimates for	each US state. However	, i do not nave a d	ataset containing information
code			Data Explanation			
mport pandas as pd						states from April 1, 2020, to Ju nt population estimate for 2024
ef load_population_data(): file_path = '/NST-EST2024-PO	P.xlsx'		Data Sample			
df_pop = pd.read_excel(file_path,			state Pop_2024\\n0 7582384\\n3 Arkansas	Alabama 5157699\\n 3088354\\n4 California		0133\\n2 Arizona
skiprows=3, skipfooter=1, # Skip the note	at the bottom		Suggested Libraries			
header=None)			pandas, openpyxl			
Pop 2021', 'Pop 2022', 'Pop 20	c_Area', 'April_2020_Base', 'Pop_2 23', 'Pop_2024']		Suggested Steps			
df_pop['Geographic_Area'] = d	raphic_Area'].str.startswith('.', na=F f_pop['Geographic_Area'].str.lstrip('	'.')				as a header. 2) The footer row
df_recent_pop = df_pop[['state			containing notes should be spread across multiple rows	in the original file. 4) Ro	ws corresponding	to national or regional
df_recent_pop['Pop_2024'] = d return df_recent_pop	f_recent_pop['Pop_2024'].astype(ir	nt)	leading dot ('.') which must	be removed for clean mer	ging. 6) Select the	ne state name column contains e state column and the column
			for the most recent year's p Wildfire Environmental	opulation (2024).		
es, I can help. The file `cleaned_	helicopter_requests_by_region.csv	r` contair	and Operational Data ns the exact information requ	ired: the total number of h	nelicopter requests	aggregated by NICC-defined
ographic areas.			Data Evalenation			
code			Data Explanation			
nport pandas as pd			by geographic area. It conta	ins two columns: 'Region	', which lists the n	
ef get_helicopter_requests_by_re	egion():		geographic area, and 'Total	riencopier Requests, wh	ion is the count of	requests for that region.
file_path = '/cleaned_helicopter df = pd.read csv(file_path, inde	requests_by_region.csv'		Data Sample Region Total Helicopter Re	nuests\\n2 Great	Basin Area	232\\n6
df_sorted = df.sort_values(by=" scending=False)			Northwest Area Southern California Area	185\\n5 Norther	n Rockies Area Iorthern California	112\\n9
return df_sorted			Southwest Area	65\\n8 39\\n0	Southern Area Alaska Area	
			Rocky Mountain Area Eastern Area	39\\n0 9\	Alaska Area	19\\n1

Figure 15: An example of the generated request by the blackboard system.

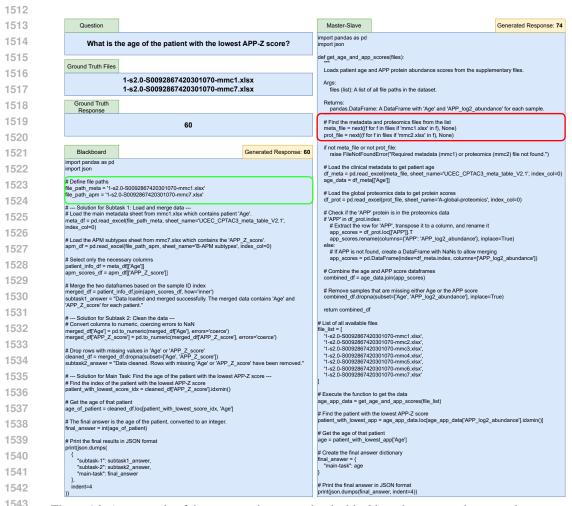


Figure 16: An example of the generated program by the blackboard system and master-slave system. The green box highlights where the blackboard system correctly selected the relevant files from the data lake, while the red box indicates where the master-slave system made an incorrect selection.

E LARGE LANGUAGE MODEL USAGE FOR WRITING

In this paper, we employ LLMs—specifically Gemini and ChatGPT—as general-purpose writing tools. Draft text is provided to these models, which are then asked to improve the writing by correcting grammatical errors and refining the structure. The edited text is then verified and edited if needed. The use of LLMs in this paper is limited strictly to text refinement. They were not employed for tasks such as generating any new content or references.