MLRC-BENCH: Can Language Agents Solve Machine Learning Research Challenges?

Yunxiang Zhang $^{\alpha*}$ Muhammad Khalifa $^{\alpha}$ Shitanshu Bhushan $^{\alpha}$ Grant D Murphy $^{\alpha}$ Lajanugen Logeswaran $^{\beta}$ Jaekyeom Kim $^{\beta}$ Moontae Lee $^{\beta\gamma}$ Honglak Lee $^{\alpha\beta}$ Lu Wang $^{\alpha}$ University of Michigan $^{\alpha}$ LG AI Research $^{\beta}$ University of Illinois at Chicago $^{\gamma}$

Abstract

We introduce MLRC-BENCH, a benchmark designed to quantify how effectively language agents can tackle challenging Machine Learning (ML) Research Competitions, with a focus on open research problems that demand novel methodologies. Unlike prior work, e.g., AI Scientist [36], which evaluates the end-to-end agentic pipeline by using LLM-as-a-judge, MLRC-BENCH measures the key steps of proposing and implementing novel research methods and evaluates them with rigorous protocol and objective metrics. Our curated suite of 7 competition tasks reveals significant challenges for LLM agents. Even the best-performing tested agent (gemini-exp-1206 under MLAB [20]) closes only 9.3% of the gap between baseline and top human participant scores. Furthermore, our analysis reveals a misalignment between the *LLM-judged* innovation and their actual performance on cutting-edge ML research problems. MLRC-BENCH is a dynamic benchmark, which is designed to continually grow with new ML competitions to encourage rigorous and objective evaluations of AI's research capabilities. Our leaderboard and code are publicly available at https://huggingface.co/spaces/launch/MLRC_Bench.

1 Introduction

Evaluating large language model (LLM) research agents [2, 32, 36] has so far been restricted to two directions. One involves tasking the agent with end-to-end scientific discovery—proposing a research idea, writing implementation code, running experiments, and eventually producing a full paper as done by AI Scientist [36]. One issue with such evaluation is the lack of reliable baseline method that enables *objective* evaluation of the proposed approach. The second direction, on the other hand, evaluates the agent's ability to produce code that solves a Kaggle-style machine learning (ML) engineering competition, skipping idea proposal and paper writing altogether [20, 7]. While evaluation in this case is straightforward, this setup rarely demands genuine research *novelty* beyond existing methods. Consequently, neither of these setups paints a full picture of whether LLM research agents can design research ideas that are both novel and effective, which we aim to address here.

Competitions at ML conferences and workshops provide a valuable testbed for evaluating research agents by assessing both novelty and effectiveness against established baselines. Unlike Kaggle-style contests, these challenges address unresolved and important problems recognized by the ML community. In addition, public leaderboards facilitate objective comparisons to human experts. If an algorithm truly outperforms known baselines, improvements in benchmark scores provide a reliable signal as to the effectiveness of the proposed method.

Therefore, this paper introduces **MLRC-BENCH** as a benchmark to evaluate the ability of LLM-based research agents to propose and implement novel methods. Drawing on tasks from recent ML conference competitions, MLRC-BENCH enables the evaluation of both **novelty** and **effectiveness**

^{*} Correspondence to yunxiang@umich.edu

Table 1: 7 MLRC-BENCH tasks representing cutting-edge machine learning research. For each competition, we show the venue where the competition is held, research area, data modality, performance metric, along with the constraints presented to the agents, including maximum allowed runtime and GPU memory based on our hardware configurations. Detailed task descriptions are given in Appendix A.

Competition	Venue	Research Area	Modality	Iodality Metric		GPU Memory
LLM Merging [47]	NeurIPS 2024	Efficient LLM	Text Accuracy, ROUGE		1 hour	48 GB
Backdoor Trigger Recovery [51]	NeurIPS 2024	LLM Safety	Text REASR, Recall		0.5 hour	48 GB
Temporal Action Localisation [18]	ECCV 2024 Workshop	Multimodal Perception	Video, mAP Audio		0.5 hour	16 GB
Rainfall Prediction [17]	NeurIPS 2023	AI for Science	Satellite Data	Critical Success Index	0.5 hour	48 GB
Machine Unlearning [48]	NeurIPS 2023	Data Privacy	Image	Forgetting Quality, Accuracy	0.5 hour	16 GB
Next Product Recommendation [26]	KDD Cup 2023	Recommendation System	Text	Mean Reciprocal Rank	0.5 hour	16 GB
Cross-Domain Meta Learning [6]	NeurIPS 2022	Few-Shot Learning	Image	Accuracy	3.5 hours	16 GB

of research agents' ideas compared to a reliable baseline method and the top human solution. In particular, it emphasizes objective metrics on tasks such as LLM merging [47] and machine unlearning [48], closely mirroring ongoing research challenges. Moreover, the challenges in MLRC-BENCH can dynamically grow by incorporating new competitions from future ML conferences.

We curate MLRC-BENCH starting with 7 competition tasks as shown in Table 1. We pick tasks that involve novel and high-impact problems, spanning areas including LLM safety, multimodal perception, and few-shot learning. Our experimental findings reveal that even the best-performing tested LLM agents, such as gemini-exp-1206 [42] under the MLAB [20] scaffolding, closes only 9.3% of the gap between baseline and top human participant score. Additionally, our analysis highlights a poor correlation between the novelty judged by LLM and practical effectiveness of agents' solutions, questioning the reliability of LLM-as-a-judge for research idea evaluation. These results underscore the limitations of current AI research agents in generating and implementing innovative ML solutions, providing a crucial benchmark for future advancements.

Our contributions can be summarized as below: 1) We introduce MLRC-BENCH, a dynamic benchmark suite curated from ML conference competitions, featuring open research problems that are both impactful and objectively measurable, and that demand the development of novel methodologies; 2) We conduct large-scale, objective evaluations for a wide array of frontier LLMs with representative agent scaffoldings, highlighting their inability to propose and implement innovative solutions with notable performance gains; 3) We pinpoint the flaws in subjective evaluations of LLM-based research agents, by showing that the LLM-judged idea novelty is misaligned with empirical effectiveness.

2 Related Work

While there are recent benchmarks that focus on code generation in machine learning domain, they do not always require methodological innovation. Works like MLAgentBench [20] and MLE-Bench [7] evaluate agents on Kaggle-style ML tasks but prioritize code implementation over novel research contributions. MLE-Bench requires the final submission to be a CSV file, limiting the data modality. Broader benchmarks such as ScienceAgentBench [13] and DiscoveryBench [38] span multiple scientific domains but lack granularity for ML-specific challenges, while CHIME [16] and OpenD5 [14] target auxiliary tasks like literature review or hypothesis generation. DSBench [27] and AAAR-1.0 [34] extend evaluations to data science and general R&D workflows but still fall short of addressing cutting-edge ML research innovation. RE-Bench [50] and MLGym [39] provide collections of ML research task environments, but their problems either fail to represent most recent research directions (e.g., image classification with CIFAR-10 [30]), or only cover a narrow range of research domains. For example, 6 out of 7 tasks in RE-bench are related to language modeling.

Besides, existing benchmarks often fail to specify computation constraints (e.g. runtime and GPU memory limit), which are important to encourage efficient yet effective solutions. The differences between our benchmark and existing work on automating ML research workflow with LLM agents are summarized in Table 3 of Appendix B.

Unlike DiscoPOP [35] and DA-Code [21], which focus on function-level coding or data science, MLRC-BENCH requires repository-level code comprehension and generation, thus better reflecting the capabilities required to work with realistic and complex research codebases. This design enables substantial project-level coding tasks, such as generating and editing multiple files [37] and reusing existing utility functions [33]. Moreover, repository-level design supports multi-agent collaboration, where specialized agents for tasks like literature survey, idea brainstorming, code writing, and performance evaluation can iteratively work together [19].

Automated end-to-end research workflows like The AI Scientist [36, 52] and MLR-Copilot [32] rely largely on subjective reviews of papers or research proposals for evaluating success. In parallel, ResearchAgent [2] iteratively refines ideas through multi-agent feedback, and Chain-of-Idea-Agent [31] organizes literature into progressive chains to stimulate ideation. However, it remains unclear how subjectively evaluated "novel" ideas translate into actual performance gains. In contrast, we explicitly investigate how such subjective assessments of novelty or idea quality align—or fail to align—with measurable performance improvements.

3 MLRC-BENCH

3.1 Task Selection

MLRC-BENCH prizes high-quality competitions that are both non-trivial and reproducible. To form our dataset, we screen competitions held at recent machine learning, natural language processing, data mining, and computer vision conferences or workshops using the following criteria:

- **Novel Research-Focused:** The tasks should require genuine methodological innovation, rather than being solvable through purely brute-force or superficial engineering approaches, such as exhaustive search for hyperparameters or features without any theoretical motivation or problem understanding.
- Non-Trivial: The problem must involve complexity so that it will not be solved by simply applying standard ML algorithms, e.g., calling the XGBoost classifier [11] on a new dataset or prompt engineering with LLMs.
- **Feasible:** Starter code, data splits, and evaluation procedures must be publicly available so that researchers, either human or agentic AI, can reproduce the experiments while keeping computational costs manageable.

These tasks represent a diverse landscape of applied ML research. As shown in Table 1, topics range from LLM safety to multimodal perception, ensuring that benchmarks stress multiple facets of algorithmic creativity. Besides, these tasks are cutting-edge. We select competitions released in recent years, whose problems remain actively researched, preventing trivial or purely off-the-shelf solutions. In addition, we have the following two considerations for MLRC-BENCH design.

Continual Updates. MLRC-BENCH is designed to evolve over time. We will continue to expand the benchmark with tasks from recent conference competitions and retire older ones once model performance saturates. This ensures the benchmark remains aligned with the frontier of machine learning research, enabling continuous tracking of agent progress on fresh and unsolved challenges. In parallel, we have released templates and detailed instructions to facilitate community contributions.³ All competitions follow a uniform and standardized code structure (detailed in Section 3.2) to streamline research and accelerate progress on LLM agents.

²While numerous competitions have been hosted at recent ML/AI conferences, only a limited subset was selected due to factors such as certain challenges being considered solved with the rapid advancement of foundation models, missing or irreproducible evaluation data/code, qualitative rather than quantitative evaluation, insufficient recency, or limited emphasis on algorithmic innovation.

³https://tinyurl.com/MLRC-Task-Template

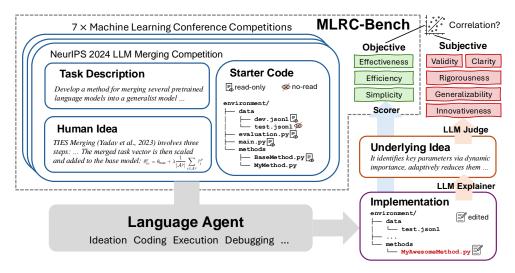


Figure 1: Overview of MLRC-BENCH and the evaluation pipeline.

Data Contamination Mitigation. Although most academic competitions require top teams to submit solution reports, the actual implementation code may not be publicly available, reducing the risk of such code appearing in the pretraining data of LLMs. As previously discussed, we will regularly update the benchmark with suitable recent competitions, ensuring that it assesses an agent's genuine research capabilities rather than its capacity for mere data retrieval and replication.

Computational Constraints. MLRC-BENCH specifies explicit computational constraints, such as runtime and GPU memory limits, closely mirroring real-world competition scenarios. These standardized constraints ensure fairness by preventing agents from leveraging excessive computational resources to gain undue advantage, thereby incentivizing efficient and innovative solutions.

3.2 Task Environment

MLRC-BENCH offers a modular, agent-agnostic environment for specifying and automatically evaluating agents on research tasks. As shown in Figure 1, for each task, we provide:

- Task Description. A detailed account of the research problem, including essential terminology, data format, desired model outputs, and constraints (e.g., limitations of model size or training time).
- Starter Code. Refactored from official competition repositories, it contains: 1) a simple, baseline model for comparison; 2) a python environment with the necessary ML frameworks/packages; 3) scripts for training, inference, and offline or online evaluation; 4) train, development and test data splits. Training data may not be available for some competitions.
- **Human Idea.** Insights from state-of-the-art papers or top-participant solution reports are included. Agents can optionally utilize these ideas to refine or inspire their own solutions.

Task-Agonistic Starter Code Structure. Because our primary goal is to focus on method development, we simplify ML experimentation by refactoring each competition starter kit into a standardized, well-organized format, comparable to common ML research project layouts (Figure 1). The resulting codebase allows users to launch experiments with a single command: python main.py -method my_awesome_method -phase dev/test, which applies the specified method to the task and evaluates the result in both development and test phases. To ensure a fair comparison and preserve the integrity of evaluations, the repository enforces file permission management: agents may only modify the methods/ directory (where the algorithmic logic resides in MyMethod.py), while evaluation scripts remain read-only. Additionally, files containing held-out test data are invisible to agents during development phase.

Development and Test Splits. We prioritize preventing overfitting by providing explicit development and test splits for each competition. Agents can choose to refine their implementations based on the development set and then submit their best-performing solution to a hidden test set. Wherever possible, we use the original competition test set (via local evaluation or online leaderboard API). Otherwise, we partition the existing development data into custom dev and test sets, reproduce the top human solution if available, and evaluate it on our new test split for a valid comparison.

3.3 Objective Evaluation Metrics

MLRC-BENCH supports objective evaluation based on three dimensions. We measure **Effectiveness** by the *performance metric* (e.g., accuracy) defined by the competition organizer, **Efficiency** by the solution *runtime* during training (if applicable) and inference, and **Simplicity** in terms of *logical lines of code (LLoC)* [40], inspired by standard practice in software estimation [40]. LLoC excludes comments and blank lines, focusing on executable statements. This metric, while imperfect, offers a rough gauge of code complexity and maintainability [3]. For better readability, we refer to LLoC as "lines of code" throughout this paper.

Main Metric: Relative Improvement to Human Solution. Quantitative performance comparisons across competitions can be tricky, as each task may differ significantly in its intrinsic difficulty, and the official baseline may be weaker or stronger. To address this, we use the *Relative Improvement to Human Solution* as our main leaderboard metric that convert each raw performance score s_{agent} into a normalized score s_{agent} , using a linear transformation [5, 50]. Therefore, the score of the baseline solution will be 0, and the top human solution in competition is set to 100. Formally, the normalization is computed as:

$$s_{\rm agent}^{'} = \frac{s_{\rm agent} - s_{\rm baseline}}{s_{\rm top_human} - s_{\rm baseline}} \times 100(\%)$$

3.4 Evaluation Protocol

Our evaluation protocol is designed to prevent AI agents from test set overfitting. Agents will submit their implementation in the form of an edited codebase, particularly within their proposed method in the methods/ directory. Specifically, in a single trial, an agent can iteratively modify the codebase multiple times. We store snapshots of the codebase immediately after each change. Whenever an execution occurs on the development set, we record the resulting metrics and the name of evaluated method for that snapshot. At the end of this iterative development phase, we pick the snapshot with the best development performance i.e., **Effectiveness**. We then evaluate the method contained in that snapshot on the test set for our final result.⁴ This approach strictly follows standard ML practice and ensures reproducible experimentation. Future work may explore more sophisticated multi-objective selection criteria that additionally weigh runtime (**Efficiency**) or lines of code (**Simplicity**) of implementations.

4 Experiments and Results

To evaluate the capability of LLM agents in solving ML research tasks, we conduct comprehensive experiments across different agent scaffoldings and language models. Each agent trial is conducted either on a single NVIDIA Quadro RTX 8000 GPU with 48GB of memory (for llm-merging, backdoor-trigger and rainfall-pred tasks) or a Tesla V100 GPU with 16GB memory (for all other tasks), determined by the size of the base model used in each task. Unless otherwise specified, we perform 8 trials⁵ per configuration and report the best attempt.

4.1 Agent Scaffolding Comparison

In addition to allow agents to directly propose and implement ideas, we investigate whether providing AI-generated or human-sourced ideas can enhance agent performance. Due to computational cost

⁴Concretely, we execute the command python main.py -method best_dev_method -phase test.

⁵The number of trials is limited to 8 due to budget constraints on API usage.

Table 2: For each competition and agent, we report the test-phase *relative improvement to the human solution*. Best performing agent in each task is highlighted in **bold**. Our results indicate that providing additional ideas, whether sourced from AI or humans, does not consistently yield performance improvements. The best-performing configuration—gemini-exp-1206 under MLAB—achieves only 9.3% of the human-level improvement over baseline on average, underscoring the inherent difficulty of these research tasks. See Table 4 in Appendix D for *absolute improvements to baseline solution*.

Agent	temporal -action-loc	llm -merging	meta -learning	product -rec	t rainfall -pred	machine -unlearning	backdoor -trigger	r Avg
MLAB (gemini-exp-1206)	-0.5	5.0	-1.1	0.1	43.1	5.6	12.9	9.3
MLAB (llama3-1-405b-instruct)	0.5	-1.0	-4.9	0.0	31.5	6.2	11.5	6.3
MLAB (o3-mini)	0.3	-1.0	-4.9	0.1	25.1	3.6	6.2	4.2
MLAB (claude-3-5-sonnet-v2)	0.8	5.0	-4.9	3.0	14.6	-94.7	39.9	-5.2
MLAB (gpt-4o)	0.3	2.0	-4.9	0.6	47.5	-18.0	10.4	5.4
w/ Human Idea	0.5	-1.0	-4.9	2.2	12.3	6.8	8.8	3.5
w/ CoI-Agent Idea (o1)	0.4	-1.0	-4.9	0.1	39.4	11.8	4.0	7.1

limits, we follow the practice of MLE-Bench [7] to evaluate a commonly used model for agentic tasks, GPT-40 [22], under three scaffolding configurations:

- MLAB: We chose a general-purpose MLAB [20] as our primary agent for evaluation. It is a ReAct-style [53] agent that alternates between thinking (such as reflection, research planning, fact-checking) and taking actions (file system operations or Python script execution) to implement methods.
- CoI-Agent Idea + MLAB: We augment MLAB with ideas generated by Chain-of-Ideas (CoI) [31], an LLM-based agent that structures relevant literature into progressive chains to enhance ideation. We choose a strong model, OpenAI's o1 [23], as the backbone for the ideation agent to generate more creative ideas.
- **Human Idea + MLAB**: To investigate whether agents can achieve notable performance gain given the right direction to work on, we provide MLAB with reliable human ideas extracted from state-of-the-art papers and top-performing competition participants' reports.

For all tasks, MLAB agents are allowed a maximum of 50 steps and 5 hours per trial, except for Rainfall Prediction, where the limits are extended to 100 steps and 10 hours due to the longer baseline training time. As shown in Table 2, incorporating additional ideas—whether generated by AI or proposed by humans—does not consistently improve performance. This underscores the challenges agents face not only in generating high-quality ideas but also in effectively implementing them, even when the ideas originate from humans.

4.2 Model Comparison

Taking MLAB as our major scaffold, we evaluate five prominent LLMs: Claude 3.5 Sonnet v2 [1], gemini-exp-1206 [42]⁷, Llama 3.1 405B Instruct [15], o3-mini-high [41] and GPT-4o (2024-11-20) [22]. The results in Table 2 show varying success rates across models and tasks. Among all models, Gemini-exp-1206 performs the best overall, closing 9.3% of the gap between baseline and top human participant scores. Claude 3.5 Sonnet V2 performs the best on most tasks but fails significantly on machine unlearning. A case study of its final solution (Appendix C.1) suggests that the failure stems from treating the removal of unwanted data and the preservation of useful knowledge as separate objectives, rather than jointly optimizing both to maintain model performance. Furthermore, our cost-effectiveness evaluation in Appendix C.2 identifies Llama 3.1 405B with MLAB as a Pareto-optimal setup balancing API cost and performance.

Table 2 also reveals that agent solutions' performance gains remain modest compared to human solutions in many cases, if not degrading baseline performance. There are, however, a few notable

⁶ Another candidate agent framework, AIDE [24] is tailored for Kaggle-style tasks with single-file solutions and tabular outputs, making it less generalizable to our repository-level coding with diverse output formats (e.g., images, tensors). Its tree-search-based design with performance-based rewards is promising but requires careful adaptation, which we leave to future work.

⁷Gemini-exp-1206 is renamed to Gemini 2.0 Pro afterwards.

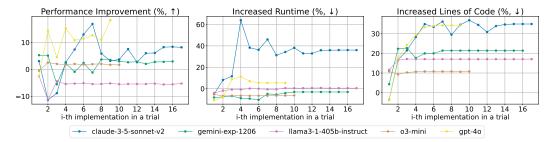


Figure 2: We track the percentages of changes of performance, runtime, and lines of code compared to baseline across iterative refinement of implementations within a trial of LLM-based MLAB agent on the development set. Performance improvement is the higher the better, while increased runtime and lines of code are the lower the better. These figures show the averaged metrics across all tasks. For results breakdown on each task, please refer to Figure 9 and 10 in Appendix D. The takeaway is that agents tend to over-refine their solutions over time, leading to increased complexity and runtime without proportional performance gains.

exceptions. For instance, MLAB (gpt-4o) achieves a score of 47.5 on the rainfall prediction task, likely because similar solutions (e.g., variants of U-Net [43]) are readily available online. In the backdoor-trigger task, the baseline GCG method [56] performs poorly—essentially making random predictions—thereby lowering the bar for agents to surpass it with more meaningful solutions. This substantial gap highlights *the current limitations of AI agents in generating novel, effective methods*, underscoring the need for further advances to match or surpass human-led research efforts.

4.3 Solution Development Analysis

Figure 2 illustrates how performance, runtime, and code complexity evolve as agents iteratively refine their implementations within a single trial. Three key trends emerge: (1) GPT-40 and Claude gradually improve their performance through refinement, while other models plateau after a few iterations; (2) runtime consistently increases, probably because models are exploring more complicated solutions over time, which may naturally conflate with better solutions; and (3) code size expands over time, reflecting increasingly complex solutions that do not yield proportional performance gains. Together, these trends suggest that agents tend to over-refine their solutions, resulting in more complex and time-consuming implementations without further performance improvements.

We further analyze the agent traces (gemini-exp-1206 with MLAB) in Appendix F to understand its behavioral patterns and limitations. Two key insights emerge. First, a significant portion of action errors—11.5% of all steps—stem from incorrect tool arguments, where the model either hallucinates or misidentifies expected parameter names. This highlights the need to *improve the agent's tool-usage capabilities*. Second, the agent was able to fix only 17.2% of the errors encountered during code execution, revealing the *challenge of self-debugging in complex, repository-level ML codebases* [12].

4.4 Subjective Evaluation with LLM-as-a-Judge

Our benchmark also facilitates the investigation on whether LLM-as-a-judge style evaluations can reliably assess the quality of research ideas by comparing subjective (here as judged by o1) and objective metrics. As shown in Figure 1, we first prompt an LLM to explain each implementation's underlying idea. We choose OpenAI's o1 [23] model here because of its superior reasoning and coding capabilities. We then apply the rubric from prior work [2] to have the o1 model assign a 1–5 Likert score on five dimensions: validity, clarity, rigorousness, generalizability, and innovativeness. These scores are all the higher the better. The prompts used for this evaluation are shown in Appendix E. Furthermore, we examine how the presence of code influences the assessments through two settings. (1) Without Code, in which the LLM judges only access the task description and the proposed idea; and (2) With Code, in which the judges also see the code implementation. We then compute Spearman's correlation [46] for each pair of objective and subjective metrics, using data from all valid implementations that include test-phase scores.

⁸We prompt the MLAB agent to include detailed comments in the code to enable faithful explanations.

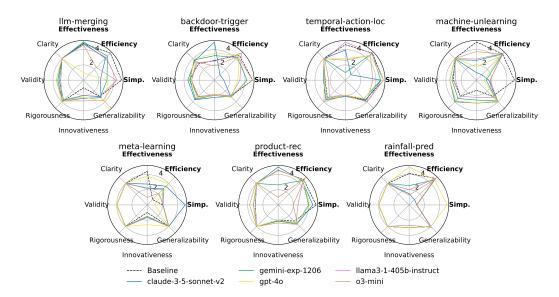


Figure 3: Radar plots of objective and subjective evaluations for agent-generated solutions across seven research tasks. Each dimension is normalized on a 1–5 scale, where higher values indicate better performance. *Objective* metrics include **Effectiveness**, **Efficiency**, and **Simplicity** (**Simp.**), which are highlighted in **bold**. The rest are *subjective* metrics, assessed by prompting o1 as a judge. Notably, more effective solutions identified by agents tend to be more complex and time-consuming (e.g., in backdoor trigger recovery). Additionally, overlapping scores in subjective dimensions suggest that LLM-based evaluation struggles to distinguish the research capabilities of different models.

Figure 3's radar plots provide a holistic view of agent performance across both subjective and objective dimensions. The plots show that while agents occasionally produce effective solutions, they often struggle to balance other criteria such as efficiency and simplicity. For instance, on the backdoor-trigger task, Claude 3.5 Sonnet V2 scores well on effectiveness but poorly on efficiency and simplicity, suggesting that agentgenerated solutions tend to be more complex and time-consuming. Notably, agents generally underperform compared to the baseline when evaluated using objective metrics. However, when subjective metrics are used and judged by LLMs, they often receive more favorable ratings. This discrepancy highlights a risk of overly optimistic conclusions when relying solely on subjective evaluations.

Figures 4 and 7 (in Appendix D) illustrate the correlation heatmaps for both settings. The overall correlations remain weak. For example, there is a near-zero correlation (-0.06) between innovativeness



Figure 4: Correlation heatmap between objective (x-axis) and subjective (y-axis) metrics for agent-generated solutions across all tasks. Code is included when prompting the LLM to evaluate subjective dimensions. No strong correlation is observed, suggesting that LLM-judged subjective metrics may not reliably indicate empirical performance gains.

and effectiveness, implying that an agent's ability to generate novel ideas, as judged by an LLM, does not necessarily equate to success in practical tasks. Consequently, our finding indicates that *LLM-based evaluations alone are not a reliable proxy for real-world research impact*. While LLM agents can certainly assist in generating creative ideas, relying solely on LLM-based evaluations to gauge agents' progress in improving machine learning research may lead to misinterpretations. This again highlights the importance of employing objective metrics to ensure that proposed solutions are not only novel but also effective.

⁹We find that removing the code leads to similar correlation results and does not significantly affect the conclusion we make.

4.5 Inference-Time Scaling on ML Research Tasks

Increasing inference-time compute via repeated sampling [7, 8, 4] has been shown to boost LLM performance on reasoning and coding tasks. Here we explore how LLM research agents scale with more inference-time compute on both the idea and solution spaces. We sample 4 ideas for each task from CoI-Agent [31] and repeat MLAB agent for 8 trials to implement each idea into code.

Figure 5 plots pass@k [9], i.e., the probability that at least one of k trials converges to a successful implementation, defined as the agent closes at least 5% of the gap between baseline and top human participant scores (Relative Improvement to Human, Section 3.3). Our results show that providing high-quality ideas enhances an agent's ability to generate meaningful solutions when given multiple attempts, and human ideas appear to be more effective than those produced by AI. Furthermore, under a fixed inference budget, we did not observe a significant difference between allocating resources to idea exploration versus exploitation. For example, there is no significant pass@k difference between using 4 ideas with 2 trials per idea, 2 ideas with 4 trials per idea, and 1 idea with 8 trials per idea. This phenomenon likely occurs because once a highquality idea is identified, the performance gains from additional trials tend to plateau, resulting in diminishing returns despite further exploita-

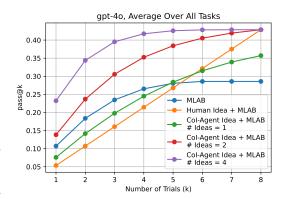


Figure 5: We measure Pass@k as we scale the number of trials and ideas, running MLAB for eight trials per idea. The total inference-time computes are equivalent among these points: k=4 for one-idea line, k=2 for two-idea line, k=1 for four-idea line, and k=4 for the remaining lines. For results breakdown on each task, please refer to Figure 8 in Appendix D. Our results indicate that 1) providing high-quality ideas—especially human-generated ones—significantly boosts an agent's success rate across multiple attempts, 2) while varying the balance between idea exploration and exploitation under a fixed budget yields similar outcomes due to diminishing returns from repeated trials.

tion. We hypothesize that performance-informed tree-search that navigates the vast space of possible solutions [24, 29] or allocating more computational resources [7] could offer more promising scaling properties.

5 Conclusion and Limitations

MLRC-BENCH draws upon the rigor of conference competitions to provide a scalable, objective, and realistic benchmark for evaluating LLM agents in proposing and implementing novel algorithms that advance research on impactful topics. Our benchmark features modular tasks, objective evaluation metrics, tamper-proof evaluations, and ongoing updates as new suitable competitions become available. Our results show that MLRC-BENCH presents a significant challenge for state-of-the-art LLMs and agent scaffoldings. Our analysis highlights the misalignment between the LLM-judged innovation and their actual performance on cutting-edge ML research problems. MLRC-BENCH will evolve alongside the rapid pace of ML research and continuously support the pursuit of AI-assisted or automated scientific discovery.

Our work has certain limitations. First, the benchmark currently covers only seven competitions, though we plan to expand it with new tasks as outlined in Section 3.1. Second, our evaluation focuses solely on a general-purpose MLAB agent for code implementation, as some other agent frameworks are currently unsuitable for our repository-level coding tasks (see Section 4.1). Nonetheless, recent advancements in LLM agents, including improved debugging tools [54], inference-time scaling techniques [49, 29], and reinforcement learning-based fine-tuning [55, 25, 10, 45], have the potential to achieve substantially better performance at a lower cost than human researchers.

Acknowledgements

This work is supported by LG AI Research. We thank Xinnuo Li and Aryan Sharma for their help with the survey and collection of research competitions. We also thank LAUNCH lab members for their useful feedback.

References

- [1] AI Anthropic. Claude 3.5 sonnet model card addendum. Claude-3.5 Model Card, 3:6, 2024.
- [2] Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. Researchagent: Iterative research idea generation over scientific literature with large language models. *CoRR*, abs/2404.07738, 2024.
- [3] Kaushal Bhatt, Vinit Tarey, Pushpraj Patel, Kaushal Bhatt Mits, and Datana Ujjain. Analysis of source lines of code (sloc) metric. *International Journal of Emerging Technology and Advanced Engineering*, 2(5):150–154, 2012.
- [4] Bradley C. A. Brown, Jordan Juravsky, Ryan Saul Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *CoRR*, abs/2407.21787, 2024.
- [5] Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, Ilya Sutskever, and Jeffrey Wu. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. In Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net, 2024.
- [6] Dustin Carrión-Ojeda, Hong Chen, Adrian El Baz, Sergio Escalera, Chaoyu Guan, Isabelle Guyon, Ihsan Ullah, Xin Wang, and Wenwu Zhu. Neurips'22 cross-domain metadl competition: Design and baseline results. In Pavel Brazdil, Jan N. van Rijn, Henry Gouk, and Felix Mohr, editors, ECML/PKDD Workshop on Meta-Knowledge Transfer, 23 September 2022, Grenoble, France, volume 191 of Proceedings of Machine Learning Research, pages 24–37. PMLR, 2022.
- [7] Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Lilian Weng, and Aleksander Madry. Mle-bench: Evaluating machine learning agents on machine learning engineering. *CoRR*, abs/2410.07095, 2024.
- [8] Lingjiao Chen, Jared Quincy Davis, Boris Hanin, Peter Bailis, Ion Stoica, Matei Zaharia, and James Zou. Are more llm calls all you need? towards scaling laws of compound inference systems, 2024.
- [9] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021.
- [10] Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Haofen Wang, Jeff Z. Pan, Wen Zhang, Huajun Chen, Fan Yang, Zenan Zhou, and Weipeng Chen. Research: Learning to reason with search for llms via reinforcement learning. *CoRR*, abs/2503.19470, 2025.

- [11] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In Balaji Krishnapuram, Mohak Shah, Alexander J. Smola, Charu C. Aggarwal, Dou Shen, and Rajeev Rastogi, editors, *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pages 785–794. ACM, 2016.
- [12] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [13] 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. Scienceagent-bench: Toward rigorous assessment of language agents for data-driven scientific discovery. *CoRR*, abs/2410.05080, 2024.
- [14] Jiangshu Du, Yibo Wang, Wenting Zhao, Zhongfen Deng, Shuaiqi Liu, Renze Lou, Henry Peng Zou, Pranav Narayanan Venkit, Nan Zhang, Mukund Srinath, Haoran Zhang, Vipul Gupta, Yinghui Li, Tao Li, Fei Wang, Qin Liu, Tianlin Liu, Pengzhi Gao, Congying Xia, Chen Xing, Cheng Jiayang, Zhaowei Wang, Ying Su, Raj Sanjay Shah, Ruohao Guo, Jing Gu, Haoran Li, Kangda Wei, Zihao Wang, Lu Cheng, Surangika Ranathunga, Meng Fang, Jie Fu, Fei Liu, Ruihong Huang, Eduardo Blanco, Yixin Cao, Rui Zhang, Philip S. Yu, and Wenpeng Yin. Llms assist NLP researchers: Critique paper (meta-)reviewing. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 5081–5099. Association for Computational Linguistics, 2024.
- [15] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The llama 3 herd of models. CoRR, abs/2407.21783, 2024.
- [16] Shanghua Gao, Ada Fang, Yepeng Huang, Valentina Giunchiglia, Ayush Noori, Jonathan Richard Schwarz, Yasha Ektefaie, Jovana Kondic, and Marinka Zitnik. Empowering biomedical discovery with AI agents. *CoRR*, abs/2404.02831, 2024.
- [17] Aleksandra Gruca, Federico Serva, Llorenç Lliso, Pilar Rípodas, Xavier Calbet, Pedro Herruzo, Jiří Pihrt, Rudolf Raevskyi, Petr Šimánek, Matej Choma, Yang Li, Haiyu Dong, Yury Belousov, Sergey Polezhaev, Brian Pulfer, Minseok Seo, Doyi Kim, Seungheon Shin, Eunbin Kim, Sewoong Ahn, Yeji Choi, Jinyoung Park, Minseok Son, Seungju Cho, Inyoung Lee, Changick Kim, Taehyeon Kim, Shinhwan Kang, Hyeonjeong Shin, Deukryeol Yoon, Seongha Eom, Kijung Shin, Se-Young Yun, Bertrand Le Saux, Michael K Kopp, Sepp Hochreiter, and David P Kreil. Weather4cast at neurips 2022: Super-resolution rain movie prediction under spatio-temporal shifts. In Marco Ciccone, Gustavo Stolovitzky, and Jacob Albrecht, editors, *Proceedings of the NeurIPS 2022 Competitions Track*, volume 220 of *Proceedings of Machine Learning Research*, pages 292–313. PMLR, 28 Nov–09 Dec 2022.

- [18] Joseph Heyward, João Carreira, Dima Damen, Andrew Zisserman, and Viorica Patraucean. Perception test 2024: Challenge summary and a novel hour-long videoqa benchmark. *CoRR*, abs/2411.19941, 2024.
- [19] Dong Huang, Qingwen Bu, Jie M. Zhang, Michael Luck, and Heming Cui. Agentcoder: Multi-agent-based code generation with iterative testing and optimisation. *CoRR*, abs/2312.13010, 2023.
- [20] Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Mlagentbench: Evaluating language agents on machine learning experimentation. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024.* OpenReview.net, 2024.
- [21] 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, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 13487–13521. Association for Computational Linguistics, 2024.
- [22] Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll L. Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, and Dane Sherburn. Gpt-40 system card. CoRR, abs/2410.21276, 2024.
- [23] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, and Ilge Akkaya. Openai o1 system card. CoRR, abs/2412.16720, 2024.
- [24] Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and Yuxiang Wu. AIDE: ai-driven exploration in the space of code. *CoRR*, abs/2502.13138, 2025.

- [25] Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. Searchr1: Training llms to reason and leverage search engines with reinforcement learning. CoRR, abs/2503.09516, 2025.
- [26] Wei Jin, Haitao Mao, Zheng Li, Haoming Jiang, Chen Luo, Hongzhi Wen, Haoyu Han, Hanqing Lu, Zhengyang Wang, Ruirui Li, Zhen Li, Monica Xiao Cheng, Rahul Goutam, Haiyang Zhang, Karthik Subbian, Suhang Wang, Yizhou Sun, Jiliang Tang, Bing Yin, and Xianfeng Tang. Amazon-m2: A multilingual multi-locale shopping session dataset for recommendation and text generation. *arXiv preprint arXiv:2307.09688*, 2023.
- [27] 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 to becoming data science experts? *CoRR*, abs/2409.07703, 2024.
- [28] Sayash Kapoor, Benedikt Stroebl, Zachary S. Siegel, Nitya Nadgir, and Arvind Narayanan. AI agents that matter. CoRR, abs/2407.01502, 2024.
- [29] Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. Tree search for language model agents. *CoRR*, abs/2407.01476, 2024.
- [30] Alex Krizhevsky et al. Learning multiple layers of features from tiny images. 2009.
- [31] Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang, Yifei Xin, Ronghao Dang, Deli Zhao, Yu Rong, Tian Feng, and Lidong Bing. Chain of ideas: Revolutionizing research via novel idea development with LLM agents. *CoRR*, abs/2410.13185, 2024.
- [32] Ruochen Li, Teerth Patel, Qingyun Wang, and Xinya Du. Mlr-copilot: Autonomous machine learning research based on large language models agents. *CoRR*, abs/2408.14033, 2024.
- [33] Dianshu Liao, Shidong Pan, Xiaoyu Sun, Xiaoxue Ren, Qing Huang, Zhenchang Xing, Huan Jin, and Qinying Li. \$\mathbf{A^{3}}\sa3-codgen: A repository-level code generation framework for code reuse with local-aware, global-aware, and third-party-library-aware. *IEEE Trans. Software Eng.*, 50(12):3369–3384, 2024.
- [34] Renze Lou, Hanzi Xu, Sijia Wang, Jiangshu Du, Ryo Kamoi, Xiaoxin Lu, Jian Xie, Yuxuan Sun, Yusen Zhang, Jihyun Janice Ahn, Hongchao Fang, Zhuoyang Zou, Wenchao Ma, Xi Li, Kai Zhang, Congying Xia, Lifu Huang, and Wenpeng Yin. AAAR-1.0: assessing ai's potential to assist research. *CoRR*, abs/2410.22394, 2024.
- [35] Chris Lu, Samuel Holt, Claudio Fanconi, Alex J. Chan, Jakob N. Foerster, Mihaela van der Schaar, and Robert Tjarko Lange. Discovering preference optimization algorithms with and for large language models. *CoRR*, abs/2406.08414, 2024.
- [36] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI scientist: Towards fully automated open-ended scientific discovery. *CoRR*, abs/2408.06292, 2024.
- [37] Hongru Ma, Yanjie Liang, Jiasheng Si, Weiyu Zhang, Hongjiao Guan, Chaoqun Zheng, Bing Xu, and Wenpeng Lu. Srlcg: Self-rectified large-scale code generation with multidimensional chain-of-thought and dynamic backtracking. *arXiv* preprint arXiv:2504.00532, 2025.
- [38] Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Bhavana Dalvi Mishra, Abhijeetsingh Meena, Aryan Prakhar, Tirth Vora, Tushar Khot, Ashish Sabharwal, and Peter Clark. Discoverybench: Towards data-driven discovery with large language models. *CoRR*, abs/2407.01725, 2024.
- [39] Deepak Nathani, Lovish Madaan, Nicholas Roberts, Nikolay Bashlykov, Ajay Menon, Vincent Moens, Amar Budhiraja, Despoina Magka, Vladislav Vorotilov, Gaurav Chaurasia, Dieuwke Hupkes, Ricardo Silveira Cabral, Tatiana Shavrina, Jakob N. Foerster, Yoram Bachrach, William Yang Wang, and Roberta Raileanu. Mlgym: A new framework and benchmark for advancing AI research agents. *CoRR*, abs/2502.14499, 2025.

- [40] Vu Nguyen, Sophia Deeds-Rubin, Thomas Tan, and Barry Boehm. A sloc counting standard. In *Cocomo ii forum*, volume 2007, pages 1–16. Citeseer, 2007.
- [41] OpenAI. Openai o3-mini system card, January 2025.
- [42] Sundar Pichai. Introducing gemini 2.0: our new ai model for the agentic era, 2024. Accessed: 2025-01-29.
- [43] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells III, and Alejandro F. Frangi, editors, Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015 - 18th International Conference Munich, Germany, October 5 - 9, 2015, Proceedings, Part III, volume 9351 of Lecture Notes in Computer Science, pages 234–241. Springer, 2015.
- [44] Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can llms generate novel research ideas? A large-scale human study with 100+ NLP researchers. *CoRR*, abs/2409.04109, 2024.
- [45] Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *CoRR*, abs/2503.05592, 2025.
- [46] C Spearman. The proof and measurement of association between two things. *American Journal of Psychology*, 15:72–101, 1904.
- [47] Derek Tam, Margaret Li, Prateek Yadav, Rickard Brüel Gabrielsson, Jiacheng Zhu, Kristjan Greenewald, Mikhail Yurochkin, Mohit Bansal, Colin Raffel, and Leshem Choshen. Llm merging: Building llms efficiently through merging. In *NeurIPS 2024 Competition Track*, 2024.
- [48] Eleni Triantafillou, Peter Kairouz, Fabian Pedregosa, Jamie Hayes, Meghdad Kurmanji, Kairan Zhao, Vincent Dumoulin, Júlio C. S. Jacques Júnior, Ioannis Mitliagkas, Jun Wan, Lisheng Sun-Hosoya, Sergio Escalera, Gintare Karolina Dziugaite, Peter Triantafillou, and Isabelle Guyon. Are we making progress in unlearning? findings from the first neurips unlearning competition. *CoRR*, abs/2406.09073, 2024.
- [49] Xingyao Wang. Sota on swe-bench verified with inference-time scaling and critic model. *All Hands AI Blog*, April 2025.
- [50] Hjalmar Wijk, Tao Lin, Joel Becker, Sami Jawhar, Neev Parikh, Thomas Broadley, Lawrence Chan, Michael Chen, Josh Clymer, Jai Dhyani, Elena Ericheva, Katharyn Garcia, Brian Goodrich, Nikola Jurkovic, Megan Kinniment, Aron Lajko, Seraphina Nix, Lucas Sato, William Saunders, Maksym Taran, Ben West, and Elizabeth Barnes. Re-bench: Evaluating frontier AI r&d capabilities of language model agents against human experts. CoRR, abs/2411.15114, 2024.
- [51] Zhen Xiang, Yi Zeng, Mintong Kang, Chejian Xu, Jiawei Zhang, Zhuowen Yuan, Zhaorun Chen, Chulin Xie, Fengqing Jiang, Minzhou Pan, et al. Clas 2024: The competition for llm and agent safety. In *NeurIPS 2024 Competition Track*, 2024.
- [52] Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree search. *arXiv preprint arXiv:2504.08066*, 2025.
- [53] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [54] Xingdi Yuan, Morgane M. Moss, Charbel El Feghali, Chinmay Singh, Darya Moldavskaya, Drew MacPhee, Lucas Caccia, Matheus Pereira, Minseon Kim, Alessandro Sordoni, and Marc-Alexandre Côté. debug-gym: A text-based environment for interactive debugging. *CoRR*, abs/2503.21557, 2025.

- [55] Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv* preprint arXiv:2504.03160, 2025.
- [56] Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *CoRR*, abs/2307.15043, 2023.

A Detailed Description of Research Competitions

A.1 LLM Merging [47] [Link]

Develop a novel and effective LLM merging method to improve performance on held out test set within the time constraints.

Description

Training high-performing large language models (LLMs) from scratch is a notoriously expensive and difficult task, costing hundreds of millions of dollars in compute alone. These pretrained LLMs, however, can cheaply and easily be adapted to new tasks via fine-tuning, leading to a proliferation of models that suit specific use cases. Recent work has shown that specialized fine-tuned models can be rapidly merged to combine capabilities and generalize to new skills.

The competition will provide the participants with a list of expert models that have already been trained on a task-specific dataset. The goal of this competition is to re-use the provided models to create a generalist model that can perform well on a wide variety of skills like reasoning, coding, maths, chat, and tool use. Along with these expert models, we have a set of hidden tasks that will be used to evaluate the submissions from participants.

A.2 Backdoor Trigger Recovery [51] [Link]

Backdoor Trigger Recovery for Code Generation Models

Description

Participants in this competition are tasked with developing algorithms to recover backdoor triggers embedded within large language models (LLMs) used for code generation. Each provided backdoored LLM contains multiple (trigger, target) pairs, where triggers are universal prompt injections designed to induce the generation of malicious code specified by the targets. In the development phase, participants receive a model finetuned with five known (trigger, target) pairs, while in the testing phase, the models include tens of secret (trigger, target) pairs related to various categories of harmful code generation. The objective is to predict the triggers corresponding to each provided target, adhering to a maximum token constraint of 10 tokens per trigger. Submissions will be evaluated using two metrics: recall, which measures the similarity between predicted and ground truth triggers, and the Reverse-Engineering Attack Success Rate (REASR), which assesses the effectiveness of the recovered triggers in eliciting the malicious code. Participants are provided with a starter dataset of 100 code generation queries and their correct outputs for method development and local evaluation, with additional data encouraged for enhancing method robustness. However, any attempts to access or guess the secret online evaluation dataset will be considered a rule violation.

A.3 Temporal Action Localisation [18] [Link]

Second Perception Test Challenge (ECCV 2024 Workshop) - Temporal Action Localisation Track

Description

The goal of this challenge is to develop methods that accurately **localize and classify actions** in untrimmed videos (up to 35 seconds long, 30 fps, max resolution 1080p) from a predefined set of classes.

Data

- **Training Data: Multimodal List**
- 1608 videos
- Includes both **action** and **sound** annotations
- Contains **video and audio features**
- **Validation Set**
- 401 videos, used to tune hyperparameters.

- **Test Set**
- Held-out set for final evaluation of your method's performance containing 5359 videos.

Output Format

For each video in test (or val), your model should output **all action segments**, with:

- 1. **Start timestamp**
- 2. **End timestamp**
- 3. **Predicted action class label**
- 4. **Confidence score**

Evaluation

- The main metric is Mean Average Precision (mAP), computed over your detected segments and averaged across:
- Different action classes
- IoU thresholds from 0.1 to 0.5 in increments of 0.1 (i.e., [0.1, 0.2, 0.3, 0.4, 0.5])
- You have separate splits for train, val, and test:
- Train on the training set.
- Use the validation set to tune, select models, etc.
- Evaluate final performance on the **test set**.

A.4 Rainfall Prediction [17] [Link]

Super-Resolution Rain Movie Prediction under Temporal Shifts

Description

The aim of the Weather4cast competition is to predict quantitatively future high resolution rainfall events from lower resolution satellite radiances. Ground-radar reflectivity measurements are used to calculate pan-European composite rainfall rates by the Operational Program for Exchange of Weather Radar Information (OPERA) radar network. While these are more precise, accurate, and of higher resolution than satellite data, they are expensive to obtain and not available in many parts of the world. We thus want to learn how to predict this high value rain rates from radiation measured by geostationary satellites operated by the European Organisation for the Exploitation of Meteorological Satellites (EUMETSAT).

Competition participants should predict the exact amount of rainfall for the next 8 hours in 32 time slots from an input sequence of 4 time slots of the preceding hour. The input sequence consists of four 11-band spectral satellite images. These 11 channels show slightly noisy satellite radiances covering so-called visible (VIS), water vapor (WV), and infrared (IR) bands. Each satellite image covers a 15 minute period and its pixels correspond to a spatial area of about 12km x 12km. The prediction output is a sequence of 32 images representing rain rates from ground-radar reflectivities. Output images also have a temporal resolution of 15 minutes but have higher spatial resolution, with each pixel corresponding to a spatial area of about 2km x 2km. So in addition to predicting the weather in the future, converting satellite inputs to ground-radar outputs, this adds a super-resolution task due to the coarser spatial resolution of the satellite data.

We provide training and validation data from one Eureopean region in 2019, and testing data from the same region in 2020, measuring a transfer learning performance under temporal shift. The task is to predict exact amount of rain events 4 hours into the future from a 1 hour sequence of satellite images. Rain rates computed from OPERA ground-radar reflectivities provide a ground truth.

A.5 Machine Unlearning [48] [Link]

Machine Unlearning Challenge

One-sentence summary

Develop efficient algorithms for "machine unlearning" such that, after forgetting certain training data, the resulting model closely matches one that was never trained on that data in the first place.

Description

We focus on **machine unlearning**, i.e., "removing the influence" of a subset of the training data (the *forget set*) from a trained model, so that the resulting model behaves similarly to one trained *without* that subset. This is especially relevant for privacy regulations (e.g., "right to be forgotten"), where individuals can request removal of their data from a model.

Goal

Our goal is to compare the strengths and weaknesses of different unlearning methods under a *shared* and *standardized* evaluation. Participants receive:

- 1. A **pre-trained** model (trained on facial images, CASIA-SURF, to predict age group in test phase, CIFAR-10 in dev phase).
- 2. A **forget set** (data samples to remove) and a **retain set** (the rest of training data).
- 3. A hidden **test set** for final scoring.
- **Output**: An unlearned model that should:
- **Erase** the forget set's influence to match the behavior of a retrained model that never saw those forget samples.
- **Retain** good accuracy on the remaining data and on the test set.
- **Finish** within provided compute/runtime constraints.

Data & Evaluation

- **Dataset**: CASIA-SURF, containing facial images labeled by age group (10 classes) in test phase, CIFAR-10 in dev phase.
- **Pretrained model**: A classifier trained for 30 epochs on the entire dataset.
- **Forgetting**: Must "remove" any trace of the forget set.
- **Utility**: Must stay accurate on the retain data and a hidden test set.
- **Metrics**:
- 1. **Forgetting quality** compares unlearned model θ_u to a model retrained from scratch θ_r without the forget set.
- 2. **Utility** checks retain/test accuracy relative to θ_r .
- 3. **Efficiency** run under time constraints (< 8h on provided compute).

The challenge uses an *online* evaluation on Kaggle. Each submitted unlearning method will be run multiple times against multiple "original" and "retrained-from-scratch" checkpoints, producing a final score that balances forgetting quality and model utility.

A.6 Next Product Recommendation [26] [Link]

This task focuses on next product recommendation by predicting the most likely product a customer will engage with based on session data and product attributes, using test data from English, German, and Japanese locales.

Description

For each session, the participant should predict 100 product IDs (ASINs) that are most likely to be engaged with. The product IDs should be stored in a list and are listed in decreasing order of confidence, with the most confident prediction at index 0 and least confident prediction at index 99. Evaluation is performed using mean reciprocal rank where the rank in your list of the ground truth next item is being assessed. For each session, you will be provided with the locale of the user and a list of products already viewed in that session. A separate file has metadata about each product.

Table 3: Comparison between MLRC-BENCH and existing work on automated scientific discovery in machine learning with LLM agents. "~" means that some but not all of the tasks in that benchmark require the indicated capability. "Compute Constraints" indicates whether the solution code must adhere to specified runtime and GPU memory limitations.

	Problem Identification		Experiment Design	Code Implementation	Evaluation Method	Evaluation Object	Compute Constraints	Continual Updates
AI Scientist [36]	✓	✓	✓	✓	LLM & Human Judge	Paper		
Can LLMs Generate Novel Research Ideas? [44]	✓	✓	✓		Human Judge	Idea Proposal		
DiscoPOP [35]		✓		✓	Performance -Based	Function-Level Code		
MLAgentBench [20]		~		✓	Performance -Based	File-Level Code		
MLE-Bench [7]		~		✓	Performance -Based	File-Level Code		
MLGym-Bench [39]		~		✓	Performance -Based	File-Level Code		
RE-Bench [50]		✓		✓	Performance -Based	Repository -Level Code	✓	
MLRC-BENCH (Ours)		✓		✓	Performance -Based	Repository -Level Code	✓	✓

A.7 Cross-Domain Meta Learning [6] [Link]

The competition focuses on cross-domain meta-learning for few-shot image classification, challenging participants to develop scalable and robust models that can quickly adapt to diverse tasks with varying numbers of classes ("ways") and training examples per class ("shots") across domains like healthcare, ecology, and manufacturing.

Description

Goal and Data

This competition challenges participants to develop meta-learning models that adapt quickly to few-shot classification tasks across ten diverse domains (e.g., healthcare, ecology, manufacturing). Drawing on the newly expanded Meta Album meta-dataset (10 image datasets unified at 128×128 resolution), the final evaluation tasks vary in "ways" (2–20 classes) and "shots" (1–20 training examples per class). By combining such heterogeneous tasks, the challenge highlights the importance of scalability, robustness to domain shifts, and flexible generalization in the "any-way any-shot" meta-learning setting. 5 datasets will be used for training and 5 will be used for testing.

Participants develop a 'MetaLearner' whose 'meta_fit' function returns a 'Learner' whose 'fit' function returns a 'Predictor' with a 'predict' function.

Evaluation and Metric

Submissions are evaluated with blind testing on ten representative datasets. Each task includes a support set (training) and a query set (testing), and the competition's primary metric is a random-guess normalized balanced accuracy. First, a balanced classification accuracy (bac) is computed by averaging per-class accuracies (i.e., macro-average recall). Then, to account for varying numbers of classes (ways), the bac is normalized by the expected performance of random guessing. This ensures a fair comparison across tasks with different ways/shots configurations and highlights each model's true ability to learn effectively from limited examples in multiple domains.

B Additional Discussion on Related Work

Scientific discovery in machine learning typically includes four main stages: **Problem Identification**, where gaps in existing methods are recognized; **Method Proposal**, which introduces a new approach to address the issue; **Experiment Design**, involving the selection of datasets, baselines, and metrics for evaluation; and **Code Implementation**, where the method is realized through executable code. While prior work [36, 44] covers Problem Identification and Experiment Design, evaluation could be subjective based on idea proposal or final paper. Instead, MLRC-BENCH focuses on the critical

stages of proposing and implementing novel methods, enabling objective performance assessment. The differences between our benchmark and existing work on automating ML research workflow with LLM agents are presented in Table 3.

C Additional Analysis

C.1 Case Study

We present two case studies below to illustrate failed solutions implemented by LLM agents. Please see Appendix G for the concrete code implemented by AI agents.

LLM Merging Challenge: The objective is to develop a novel and effective algorithm to merge several (in our case, two) expert models into a single model that demonstrates improved performance on a held-out test set within the given time constraints. The MLAB agent (o3-mini) proposed a median aggregation of parameters, which slightly underperforms the baseline of a mean aggregation. We hypothesize that the median, while robust to extreme outliers, typically exhibits higher statistical variability when merging multiple parameter sets, especially with fewer models.

Machine Unlearning Challenge: The goal is to develop efficient algorithms that enable a model to "forget" specific training data, such that the resulting model closely resembles one that was never trained on that data in the first place. The MLAB agent (claude-3-5-sonnet-v2) proposed a Gradient Ascent Unlearning Algorithm, a two-phase approach combining gradient ascent for forgetting and fine-tuning for retaining knowledge. Specifically, the algorithm first performs gradient ascent on the forget set to maximize loss (achieving unlearning) and then fine-tunes the model on the retain set to restore the desired knowledge. While this approach sounds promising in theory, it scored significantly lower than the baseline. We hypothesize that by separating the gradient ascent on the forget set and the fine-tuning on the retain set into two distinct phases, the model may not effectively balance these two conflicting objectives. In contrast, a joint optimization approach—where both objectives are optimized at each gradient update—might better balance the processes of "forgetting" and "retaining" knowledge.

C.2 Cost-Effectiveness Analysis

In Figure 6, we analyze the agents' success rates in a cost-controlled setting, motivated by recent work [28] emphasizing the importance of jointly optimizing both performance and cost in agent design. Llama 3.1 405b Instruct offers the most favorable trade-off, achieving higher success rates than GPT-40 and Claude 3.5 Sonnet at a significantly lower cost.

Although incorporating an ideation phase before implementation improves overall performance compared to the implementation-only MLAB setting, it incurs additional costs due to the generation of research ideas. Nevertheless, we believe the performance gain will increasingly justify the added cost as base models continue to grow stronger, particularly for complex research problems where strategic high-level planning leads to substantial gains in final outcomes.

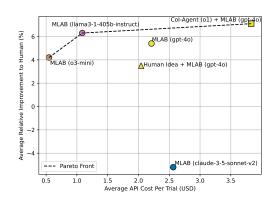


Figure 6: We perform a cost-effectiveness analysis of various setups. On the x-axis, we plot API cost, where lower is better, and on the y-axis, we show reletive improvement to human (Section 3.3), where higher is better.

¹⁰We exclude the gemini-exp-1206 model from this figure because it was experimental and its pricing was unavailable at the time of writing.

¹¹We estimate the API cost for Llama models based on Amazon Bedrock service pricing.

Table 4: For each research competition and agent, we report the test-phase *best percentage improvement in the performance metric over the baseline among 8 trails* provided in the starter code. Additionally, we present the improvements achieved by the top human participants at the time of competition under the same setup. Best performing agent in each task is highlighted in **bold**. Agents can only achieve marginal performance gains compared to human experts, and in many cases, the agents' solutions even degrade baseline performance.

Agent	temporal -action-loc	llm -merging			rainfall -pred		backdoor -trigger	Avg
MLAB (gemini-exp-1206)	-1.3	3.4	-3.2	0.6	91.4	3.5	80.4	25.0
MLAB (llama3-1-405b-instruct)	1.5	-0.7	-14.9	0.0	66.7	3.8	71.7	18.3
MLAB (o3-mini)	0.9	-0.7	-14.9	0.6	53.3	2.2	38.8	11.5
MLAB (claude-3-5-sonnet-v2)	2.2	3.4	-14.9	12.3	31.0	-58.6	247.9	31.9
MLAB (gpt-4o)	0.9	1.4	-14.9	2.6	100.8	-11.1	64.5	20.6
w/ Human Idea	1.5	-0.7	-14.9	8.9	26.1	4.2	54.5	11.4
w/ CoI-Agent Idea (o1)	1.0	-0.7	-14.9	0.6	83.6	7.3	24.9	14.5
Top Human in Competition	284.6	68.2	304.5	412.6	212.0	61.9	621.3	280.7

D Additional Results

This section presents additional results that complement the findings reported in the main paper and appendix.

- Table 4 reports the absolute improvement over the baseline, supplementing the success rate results shown in Table 2 (Section 4.2).
- Figure 7 displays the correlation heatmap between objective and subjective metrics when LLM-as-a-Judge is applied without code as input, complementing Figure 4 (Section 4.4).
- Figure 8 shows inference-time scaling results broken down by task, complementing the aggregate results in Figure 5 (Appendix 4.5).
- Figures 9 and 10 provide a task-level analysis of the implementation process, extending the results in Figure 2 (Appendix 4.3).

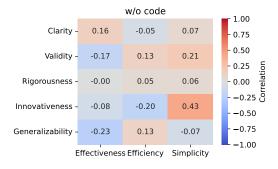


Figure 7: Correlation heatmap between objective and subjective metrics when LLM-as-a-Judge is done without code as input. The "with code" version is shown in Figure 4.

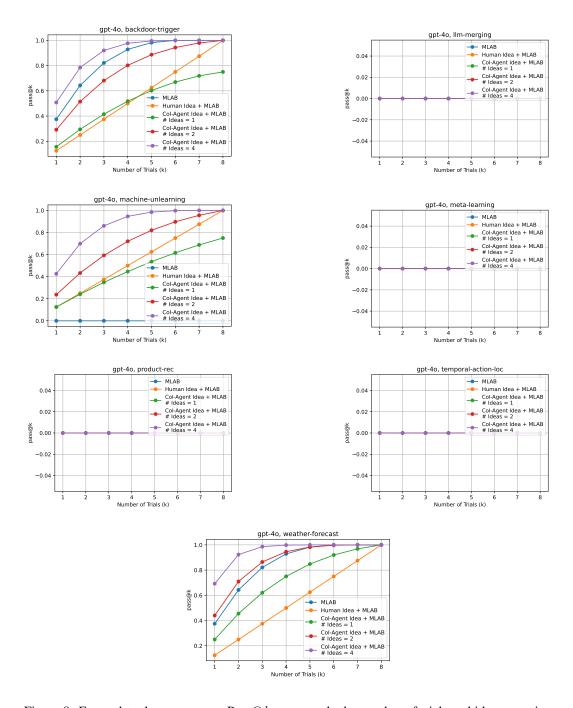


Figure 8: For each task, we measure Pass@k as we scale the number of trials and ideas, running MLAB for eight trials per idea. Pass@k is the probability that at least one of k trials converges to a successful implementation, defined as the agent closes at least 5% of the gap between baseline and top human participant scores.



Figure 9: For each task, we track the percentages of changes of performance, runtime, and lines of code compared to baseline across iterative refinement of implementations within a trial of LLM-based MLAB agent.

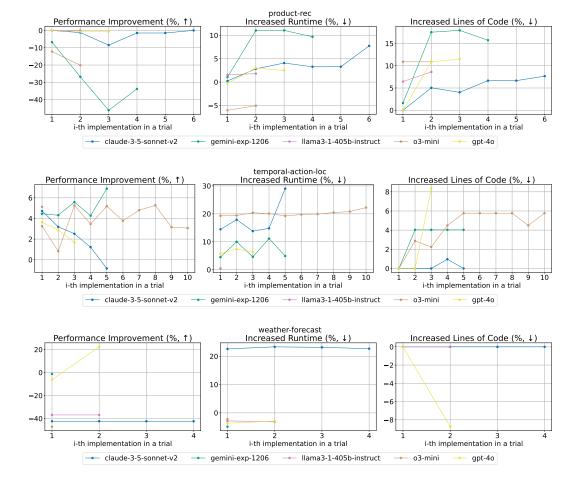


Figure 10: (Cont'd) For each task, we track the percentages of changes of performance, runtime, and lines of code compared to baseline across iterative refinement of implementations within a trial of LLM-based MLAB agent.

E Prompts for LLM-as-a-Judge

Prompt for LLM Explainer in Figure 1

Analyze the following Python code with comments and generate a high-level idea proposal summarizing:

- 1. The main goal or purpose of the method or algorithm implemented.
- 2. The general approach or methodology used to achieve the goal.
- 3. Any core assumptions or requirements underlying the implementation.

Focus on providing a conceptual overview rather than implementation details.

Code: {code}

Provide the summary as an idea proposal, avoiding references to the code itself. Focus on describing the approach and methodology as a standalone concept.

Prompt for LLM Judge in Figure 1

You are an AI assistant whose primary goal is to assess the quality and soundness of scientific methods across diverse dimensions, in order to aid researchers in refining their methods based on your evaluations and feedback, thereby enhancing the impact and reach of their work.

You are going to evaluate a scientific method for its {metric} in addressing a research problem, focusing on how well it is described in a clear, precise, and understandable manner that allows for replication and comprehension of the approach.

As part of your evaluation, you can refer to the research problem, which will help in understanding the context of the proposed method for a more comprehensive assessment.

Research problem: {researchProblem}

Now, proceed with your {metric} evaluation approach that should be systematic:

- Start by thoroughly reading the proposed method and its rationale, keeping in mind the context provided by the research problem, and existing studies mentioned above.
- Next, generate a review and feedback that should be constructive, helpful, and concise, focusing on the {metric} of the method.
- Finally, provide a score on a 5-point Likert scale, with 1 being the lowest, please ensuring a discerning and critical evaluation to avoid a tendency towards uniformly high ratings (4-5) unless fully justified:

The criteria for {metric} evaluation: {criteria}

I am going to provide the proposed method with its code implementation, as follows:

Proposed method: {Method} Code implementation:

{code}

After your evaluation of the above content, please respond **only** with a valid JSON object in the following format: { "Review": "Your review here", "Feedback": "Your feedback here", "Rating": "Your rating here" }

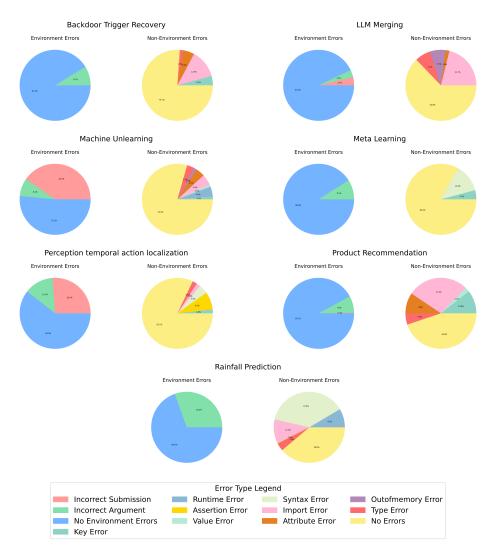


Figure 11: Distribution of environment and non-environment errors across different MLRC-Bench tasks. Each task is represented with two pie charts: one for errors related to the environment (e.g., submission issues, argument mismatches) and another for non-environment errors (e.g., runtime failures, memory issues) .

F Agent Trace Analysis

In this section, we analyze the agent traces for *Gemini-exp-1206* across different tasks. We collect trajectories across 7 tasks with 8 runs for each task, resulting in a total of 56 trajectories.

F.1 Error Type Categorization

We categorize *Gemini-exp-1206*'s actions on MLRC-Bench tasks into two main types:

- Non-execute Steps: Steps where the action "Execute Script" was not invoked.
- Execute Steps: Steps where the action "Execute Script" was invoked.

Non-execute errors are further classified into incorrect argument and incorrect submissions. Execute errors include key errors, value errors, type errors, assertion errors, runtime errors, attribute errors, out-of-memory errors, import errors, and syntax errors.

We briefly explain the errors encountered:

- EnvError: Occurs when submissions do not match the leaderboard records, files are missing, or arguments are passed incorrectly.
- **KeyError**: Results from passing incorrect argument names or not registering methods.
- ValueError: Triggered by invalid parameters, such as an improper learning rate or an empty parameter list.
- TypeError: Occurs from unexpected keyword arguments.
- AssertionError: Occurs when conditions such as shape compatibility or divisibility are not met.
- RuntimeError: Typically related to tensor shape issues.
- AttributeError: Happens when a required attribute is missing.
- OutofmemoryError: Indicates a CUDA out-of-memory condition.
- **ImportError**: Occurs when a module cannot be imported.
- SyntaxError: Triggered by syntax issues, such as a missing comma.

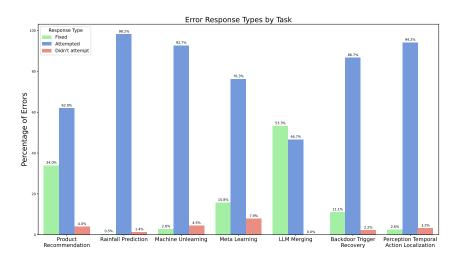


Figure 12: Error response distribution across tasks. For each task, errors are classified as *Fixed* (fully resolved), *Attempted* (partially addressed), or *Didn't attempt* (unresolved). These labels were assigned by GPT-4o-mini after evaluating each error along with all its subsequent steps (action, reflection, thought, and observation).

Figure 11 shows that *Gemini-exp-1206* successfully completes a considerable number of steps without errors, yet its performance varies noticeably across tasks. In particular, while Meta Learning displays

relatively few issues, Rainfall Prediction exhibits a higher frequency of "hallucination" based errors such as incorrect argument handling, non-existent file references, and invalid parameter choices. This discrepancy indicates that certain tasks present greater challenges for the model, likely due to more complex or less familiar contexts.

Within the **Execute Steps**, the most frequent error types are import, value, and type errors, reflecting a tendency to reference nonexistent modules, pass invalid parameters, or supply arguments of the wrong data type. On the **Non-execute Steps** side, incorrect arguments remain a recurring challenge, showing another case where the agent seems to be "hallucinating" the argument names.

Taken together, these findings highlight the generally robust completion of tasks, but also highlight the need to refine the agent's internal checks to reduce parameter mismatches and submission errors. Strengthening agent self-verification strategies could help mitigate hallucinations and further align its outputs with the intended specifications of each task.

F.2 Error Response Distribution

Figure 12 presents an overview of how errors are handled across the seven tasks, highlighting the proportion of errors that were fully resolved (*Fixed*), partially addressed (*Attempted*), or left unaddressed (*Did not attempt*). These groupings were derived by passing each error along with all its next steps— containing its *action*, *reflection*, *thought* and *observation*— to GPT-4o-mini, and the error was then labeled based on whether it was successfully resolved, partially addressed, or not addressed at all.

From Figure 12, we observe notable variations in error-handling effectiveness across tasks. Specifically, *LLM Merging* demonstrates the highest proportion of fully resolved (*Fixed*) errors, indicating more effective resolution strategies, whereas *Rainfall Prediction*, *Backdoor Trigger Recovery*, *Machine Unlearning*, and *Perception Temporal Action Localization* predominantly exhibit errors that are only partially addressed (*Attempted*). Meanwhile, *Meta Learning* has the largest share of errors categorized as *Did not attempt*. These distinctions highlight task-specific differences in error management.

We also observe a consistently high percentage of errors categorized as *Attempted* across nearly all tasks, indicating that the agent often struggles to fully resolve errors. This broadly suggests challenges in the agent's comprehension or planning capabilities when addressing complex errors, potentially pointing to difficulties in fully interpreting the underlying problem or effectively formulating corrective actions. Additionally, the notable variability in fully resolved (*Fixed*) and unaddressed (*Did not attempt*) errors across tasks implies that certain tasks inherently pose greater cognitive complexity or ambiguity, further exacerbating the agent's difficulty in error resolution. The prompts used for this annotation are shown in Appendix I.

F.3 Error Solve Rate

To further expand on this analysis of errors, we also show which error types are more effectively resolved and highlight their associated complexity during task execution. Errors which were categorised as *Fixed* are treated as solved while errors belonging to the other two categories are treated as unsolved. Using the prompt in Appendix I, we also had GPT-40-mini return the step at which the error was fixed for those that were categorised as fixed.

Figure 13 provides insights into the solve rates across different error types, revealing variability in the agent's efficiency in resolving specific errors. Among the error types, *Out of Memory Error* achieved the highest solve rate, suggesting that these errors are relatively straightforward for the agent to diagnose and address. In contrast, *Syntax Errors* and *Environment Errors* demonstrated lower resolution rates, while *Value Error*, *Runtime Error*, and *Assertion Error* were never fixed, highlighting their inherent complexity or ambiguity.

Additionally, the average number of steps taken to resolve errors further underscores these differences. Notably, *Out of Memory Errors* required the highest average number of steps, indicating that, although these errors are ultimately resolved at a high rate, their resolution involves a complex, multi-step process. Conversely, when *Syntax Errors* and *Environment Errors* are fixed, they tend to be resolved more quickly, suggesting that these issues, while more challenging to fix overall, can be diagnosed and corrected with fewer steps when addressed successfully.

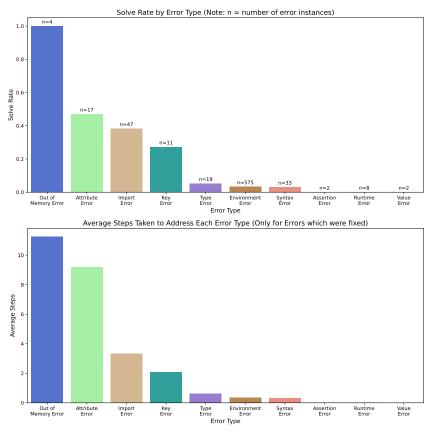


Figure 13: Solve rate and average steps taken for resolving various error types. The top chart shows the proportion of errors successfully resolved (*Solve Rate*), annotated with the total number of instances per error type. The bottom chart illustrates the average number of steps required to achieve resolution, only errors which were fixed were used to calculate average steps.

F.4 Per-Step Action Distribution

Figure 14 presents how frequently each action (*List Files, Understand File, Edit Script (AI), Execute Script, Copy File, Undo Edit Script, Inspect Script Lines*, and *Final Answer*) is used over the maximum allowed 50 steps. This breakdown helps us observe when the agent transitions from environment exploration to iterative code refinement and debugging. In particular, *Rainfall Prediction* was not used for this analysis, as it was run for 100 steps.

Early steps are dominated by environment-inspection actions, particularly *List Files* and *Understand File*, which give the agent context about available files and their contents. As the trajectory progresses, the agent increasingly relies on *Edit Script (AI)* and *Execute Script* for iterative code modifications and testing, while *Inspect Script Lines* helps to target debugging. *Undo Edit Script* is used far less frequently, suggesting that the agent rarely reverts to a previous state. This pattern highlights an iterative development approach, but also indicates that the agent may underutilize rollback strategies when encountering errors. Although *Final Answer* typically signals the end, some runs exhibit early submission, indicating missed opportunities for further refinements.

F.5 Per-Step Stage Distribution

In this section, we analyze the per-step stage distribution, categorizing the steps into seven stages based on GPT-40 annotations: Understanding & Exploration, Baseline Assessment, Problem Analysis & Idea Generation, Implementation, Debugging & Error Handling, Experimental Refinement, and Final Evaluation & Submission. Each

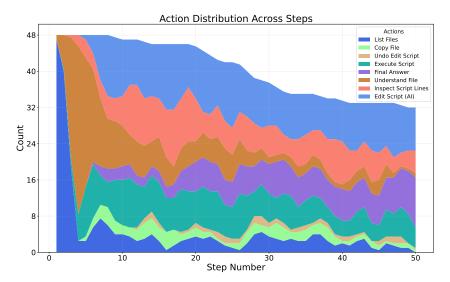


Figure 14: Distribution of Actions Taken Across Steps. This visualization depicts how frequently different types of actions were taken at each step by the agent.

step in the agent's trajectory—comprising its *Reflection*, *Thought*, *Action Input*, and *Action*—was labeled by GPT-40, which matched the step content to the most relevant stage criteria.

Figure 16 visualizes the distribution of these seven stages over the course of the maximum allowed 50 steps. In particular, *Rainfall Prediction* was not used for this analysis, as it was ran for 100 steps.

The early steps are predominantly labeled **Understanding & Exploration**, reflecting the initial focus of the agent on examining files, reviewing the environment, and clarifying task requirements. A smaller portion of these early steps is allocated to **Baseline Assessment**, where the agent measures the performance of the unmodified solution to establish a reference point.

As the agent progresses, the distribution shifts noticeably toward **Implementation**, reflecting a transition from initial passive information gathering to active code modifications. Notably, the agent dedicates very few steps to **Problem Analysis & Idea Generation**, suggesting a rapid move from conceptual planning to execution. This change is often accompanied by a surge in **Debugging & Error Handling** steps, as newly introduced modifications lead to runtime or logical errors that must be diagnosed and fixed. The close interplay between **Implementation** and **Debugging & Error Handling** underscores the iterative nature of the agent's development process.

Interestingly, it should be noted that the agent continues to spend a substantial number of steps in the **Understanding & Exploration** stage. This ongoing emphasis highlights the inherent complexity and cognitive demands of repository-level tasks, which often require extensive file navigation and conceptual understanding.

Toward the latter steps, a subset of runs proceeds to **Experimental Refinement**, engaging in repeated re-runs, parameter tuning, and exploring alternative strategies to optimize performance. However, in many cases, the agent transitions relatively quickly to **Final Evaluation & Submission**. This early move towards final submission implies potential underuse of iterative enhancement cycles, indicating an area for improvement in the agent's approach. The prompts used for this annotation are shown in Appendix H.

F.6 Stage Timelines

Using the stage annotation from the previous section, we now extend our analysis by visualizing stage timelines for each task and run. Figure 15 depict the duration the model spends in each stage, ranging from **Understanding & Exploration** to **Final Evaluation & Submission**, with block widths proportional to the time allocated. The overall run durations are also displayed, providing context for the stage-wise time distribution. Notably, *Rainfall Prediction* was not used for this analysis, as it was ran for 10 hours.

Stage Timelines Across Multiple Tasks



Figure 15: Combined stage timelines across multiple tasks. Each timeline represents an individual run, with block widths proportional to the time spent in each stage. The total duration of each run is shown on the right.

Table 5: Highest Capability Levels Across Experimental Runs for Evaluated Agents. This table reports the highest capability level, as defined by L1–L8 metric, achieved by each agent over eight runs across seven distinct tasks. Each run is assigned a numeric score corresponding to its level (e.g., L6 = 6, L5 = 5, and so on).

Agent	temporal -action-loc	llm -merging	product -rec	rainfall -pred	meta -learning	machine -unlearning	backdoor -trigger
MLAB (gemini-exp-1206)	4	5	5	6	4	6	6
MLAB (llama3-1-405b-instruct)	5	3	5	6	3	6	6
MLAB (o3-mini)	5	3	5	6	3	5	6
MLAB (claude-3-5-sonnet-v2)	5	5	5	6	3	4	6
MLAB (gpt-4o)	5	5	5	6	3	4	6
Human Idea + MLAB (gpt-4o)	5	3	5	6	3	6	6
CoI-Agent (o1) Idea + MLAB (gpt-4o)	5	3	5	6	3	6	5

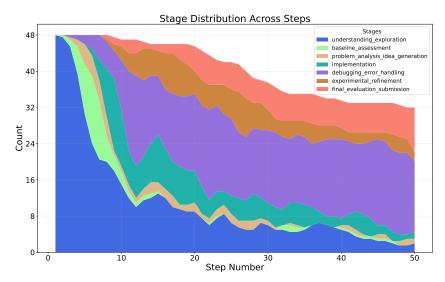


Figure 16: Stage distribution across each step, annotated using GPT-40 and grouped into seven distinct stages to illustrate shifts in task focus and activity over the course of all tasks.

F.7 Capability Level

We categorize each experimental run into one of eight capability levels (L1 to L8) based on its performance relative to both a baseline and the top human solution. Definitions of each level are described below:

- L1: No Valid Output. The agent fails to generate any valid evaluation outputs on either the development or test set, indicating a complete inability to produce usable predictions.
- L2: Test Submission Failure. The agent processes the development set but fails to produce a valid submission on the test set, meaning that while some processing occurs, the pipeline does not yield a final result.
- L3: Unimproved but Valid. The agent produces valid predictions for both the development and test sets yet remains below the baseline performance throughout the run.
- L4: Overfitting. The agent outperforms the baseline on the development set but falls short on the test set, suggesting that the model may have overfitted to the development data.
- L5: Baseline-Comparable. The agent's test performance surpasses the baseline but remains under 5% of the margin by which the top human solution exceeds the baseline. In this range, performance is very close to the baseline level.
- **L6: Notable Gains.** The agent's test performance exceeds the baseline by an improvement margin between 5% and 50% of the gap between the baseline and the top human solution. In practical terms, this level is our "success" scenario because it indicates the agent has closed a meaningful portion of the gap above the baseline.
- L7: Near Human-Level. The agent captures between 50% and 100% of the improvement margin from the baseline to the top human solution, demonstrating that the performance is approaching that of the best human score.
- L8: Superhuman. The agent exceeds top human performance not only by delivering superior quantitative results, but also by demonstrating the creative ability to generate novel ideas and implement them effectively. This level signifies that the agent can innovate beyond established human benchmarks.

This metric places an agent's performance on a tiered scale, relative to both the baseline and the top human solution, ensuring that any level of improvement (or lack thereof) is still meaningfully captured, even when the agent falls short of surpassing the baseline. Table 5 shows that rainfall-pred and backdoor-trigger are relatively easier tasks in our benchmark as the agent can achieve a

meaningful improvement over the baseline (L6), though still far behind the human. The other tasks appear to be very difficult for all agents, as they cannot achieve capability levels greater than L5.

G Agent Code Samples

Here we show the two examples of solutions generated by LLM agents mentioned in Section C.1.

```
# High-Level Overview:
   # Purpose: This script merges multiple HuggingFace model checkpoints
      into a single base model by computing the median
   #
              of each corresponding parameter across different
      checkpoints.
   # Methodology:
4
       1. Load multiple HuggingFace model checkpoints along with their
      configurations.
   #
      2. For every parameter in the models, stack the corresponding
      tensors along a new dimension and compute the median value.
          The median is computed in float precision and then cast back
      to the original data type.
      3. Load the base model and its tokenizer.
      4. Update the base model's parameters with the aggregated median
      values and set the model to evaluation mode.
  #
10
11
   # Key Steps:
       - Load checkpoints and configurations.
12
       - Iterate over each parameter, compute the median across
      checkpoints with careful type conversion.
       - Load the base model and tokenizer.
14
       - Update the model state and prepare it for inference.
15
16
17
   import torch
  from methods.BaseMethod import BaseMethod
19
  from peft import get_peft_model, set_peft_model_state_dict
20
21
   class MedianMethod(BaseMethod):
22
       MedianMethod performs the merging of multiple checkpoint models
23
          by computing the median of each parameter.
24
       This class extends the BaseMethod to load HuggingFace
25
          checkpoints, compute a robust median-aggregated state,
       and update the base model accordingly.
26
28
       def __init__(self, name):
29
30
           Initialize the MedianMethod instance.
31
32
           Parameters:
33
               name (str): The identifier for this method instance.
34
35
           Returns:
36
37
               None
38
           # Call the parent BaseMethod's initialization method.
39
           super().__init__(name)
40
41
       def run(self):
42
43
           Execute the merging pipeline and load the updated base model.
44
45
46
           Detailed Steps:
               1. Load HuggingFace model checkpoints and configurations.
```

```
- Uses a helper function to populate
48
                      'self.loaded_models' with state dictionaries from
                      different checkpoints.
               2. Merge the checkpoints by iterating over each parameter
49
                   kev:
                  - For each parameter, retrieve the corresponding
                      tensor from every loaded model.
                  - Detach the tensor from the computation graph and
51
                      move it to CPU.
                  - Stack these tensors along a new dimension (dim=0) to
52
                      form a single tensor.
                  - Convert the stacked tensor to float for precise
53
                      median computation,
                    compute the median along the new dimension, and cast
54
                        the result back to the original data type.
               3. Load the base model's architecture and its tokenizer
55
                   via helper functions.
               4. Update the base model's parameters with the merged
56
                   state dictionary and set it to evaluation mode.
           Returns:
               torch.nn.Module: The updated base model, now containing
59
                   the median-aggregated parameters.
61
           # Step 1: Load HuggingFace model checkpoints and
62
              configurations.
           # This helper function populates self.loaded_models with
63
              state dictionaries from different checkpoints.
           super()._load_huggingface_models_and_configs()
64
65
           # Step 2: Merge checkpoints by computing the median of each
66
              parameter across all loaded models.
           # Retrieve all model state dictionaries as a list.
67
           all_models = list(self.loaded_models.values())
68
69
70
           # Assume all models share the same architecture; extract
               parameter names from the first model.
           all_parameter_names = all_models[0].keys()
72
           # Iterate over each parameter name.
73
           for parameter_name in all_parameter_names:
               # Retrieve the parameter tensor for the current parameter
75
                   from each model,
               # detaching from its computation graph and moving it to
76
                   CPU to ensure consistency in merging.
               param_list = [model[parameter_name].detach().cpu() for
                   model in all_models]
78
               # Stack the tensors along a new dimension (dim=0) to
                   create a single tensor.
               stacked_params = torch.stack(param_list, dim=0)
80
81
               # Compute the median across the new dimension.
82
               # The tensor is first cast to float for precision during
                   the median computation,
               # then the median result is cast back to the original
84
                   data type.
               median_value = torch.median(stacked_params.float(),
                   dim=0) [0].to(stacked_params.dtype)
86
               # Save the computed median tensor in the merged_model
87
                   dictionary for later use.
               self.merged_model[parameter_name] = median_value
```

```
# Step 3: Load the base model's architecture and its
90
                tokenizer.
            # These helper functions initialize the base model and
91
                configure its tokenizer.
            self._load_base_model()
92
            self._load_tokenizer()
93
94
            # Step 4: Update the base model with the merged parameters.
95
            # Load the merged state dictionary into the base model.
96
            self.base_model.load_state_dict(self.merged_model)
97
98
            # Set the base model to evaluation mode to disable
99
                training-specific layers like dropout.
            self.base_model.eval()
100
101
102
            # Return the updated base model ready for inference.
            return self.base_model
103
```

Listing 1: Median Merging Solution by MLAB (o3-mini) for the LLM Merging Challenge

```
Gradient Ascent Unlearning Algorithm
   -----
   Purpose: Selectively unlearn specific training samples while
      retaining knowledge of others
   Methodology: Two-phase approach combining gradient ascent and
      fine-tuning
   Key Steps:
6
   1. Phase 1: Gradient ascent on forget set to maximize loss
      (unlearning)
   2. Phase 2: Fine-tuning on retain set to restore desired knowledge
8
9
10
   from copy import deepcopy
11
   import torch
   from torch import nn, optim
13
   from methods.BaseMethod import BaseMethod
14
15
   DEVICE = 'cuda' if torch.cuda.is_available() else 'cpu'
16
   class GradientAscentUnlearning(BaseMethod):
18
       def __init__(self, name):
19
           """Initialize the unlearning method
20
           Args:
23
               name: Name identifier for the method
24
25
           super().__init__(name)
26
27
       def get_name(self):
           """Return the name of the unlearning method
28
29
           Returns:
30
               String identifier for the method
31
           return "gradient_ascent_unlearning"
33
34
       def run(self, net, retain_loader, forget_loader, val_loader):
35
           """Implement two-phase unlearning using gradient ascent and
36
              fine-tuning
37
           Args:
38
               net: The model to be unlearned
39
40
               retain_loader: DataLoader for retained training data
               forget_loader: DataLoader for data to be forgotten
```

```
val_loader: DataLoader for validation data
42
43
44
            Returns:
                The unlearned model
45
46
            criterion = nn.CrossEntropyLoss()
48
            # Phase 1: Gradient Ascent on forget set
49
            optimizer_forget = optim.SGD(net.parameters(), lr=0.0001,
50
                                         momentum=0.9, weight_decay=5e-4)
51
52
            for epoch in range(2): # 2 epochs for forgetting
53
                net.train()
54
55
                for batch_idx, sample in enumerate(forget_loader):
                     # Handle different data formats (dict vs tuple)
56
57
                    if isinstance(sample, dict):
                         inputs = sample["image"]
58
                         targets = sample["age_group"]
59
                    else:
60
                         inputs, targets = sample
                    inputs, targets = inputs.to(DEVICE),
62
                        targets.to(DEVICE)
63
                    optimizer_forget.zero_grad()
64
65
                    outputs = net(inputs)
66
                    loss = criterion(outputs, targets)
                    # Multiply gradients by -1 for gradient ascent
67
                    loss.backward()
68
69
                    for param in net.parameters():
                         param.grad = -param.grad
70
                    optimizer_forget.step()
            # Phase 2: Fine-tune on retain set
73
            optimizer_retain = optim.SGD(net.parameters(), lr=0.01,
74
                                         momentum=0.9, weight_decay=5e-4)
            scheduler =
76
                torch.optim.lr_scheduler.CosineAnnealingLR(optimizer_retain
                T_{max}=1
            # 5 epochs for fine-tuning
78
            for epoch in range (5):
                net.train()
                for batch_idx, sample in enumerate(retain_loader):
81
                    # Handle different data formats (dict vs tuple)
82
                    if isinstance(sample, dict):
    inputs = sample["image"]
83
84
                         targets = sample["age_group"]
85
86
                    else:
                         inputs, targets = sample
87
                    inputs, targets = inputs.to(DEVICE),
88
                        targets.to(DEVICE)
20
                    optimizer_retain.zero_grad()
90
                    outputs = net(inputs)
91
92
                    loss = criterion(outputs, targets)
93
                    loss.backward()
                    optimizer_retain.step()
94
95
                scheduler.step()
97
            net.eval()
            return net
```

Listing 2: Gradient Ascent Unlearning Solution by MLAB (claude-3-5-sonnet-v2) for the Machine Unlearning Challenge

H Prompt for Stage Annotation

Prompt for LLM Stage Annotator

You are a researcher. You are given the following trace of an AI agent working on ML research challenges:

{output_json_str}

Your task is to analyze every step in the trace and assign a stage to each step. Use the following 7 stages. For each stage, use the reasoning guidelines provided to decide if a step belongs to that stage.

1. Understanding & Exploration:

- Description: Investigate the problem statement, explore the codebase, review data files, and understand evaluation metrics. This stage is about gathering context and building a solid grasp of the task and environment
- Reasoning Guideline: Assign a step to this stage if it focuses on examining available resources, reading documentation or files, exploring the code structure, or otherwise building an initial understanding of the project.

2. Baseline Assessment:

- Description: Evaluate the unmodified baseline solution's performance to collect performance metrics and establish a reference benchmark.
- Reasoning Guideline: Assign a step to this stage if it focuses on measuring the performance of the original, unaltered solution, collecting data for baseline comparison, and ensuring the initial performance level is documented. Do not assign a step to this stage if it executes the solution after changes have been made.

3. Problem Analysis & Idea Generation:

- Description: Analyze the baseline results to identify shortcomings and brainstorm potential improvements or alternative strategies.
- Reasoning Guideline: Assign a step to this stage if it is centered on evaluating baseline outcomes, identifying issues, or generating ideas and strategies for potential improvements.

4. Implementation:

- Description: Develop and integrate the proposed modifications into the codebase by editing, extending, or refactoring the existing solution.
- Reasoning Guideline: Assign a step to this stage if it involves writing new code, modifying existing code, or integrating changes aimed at improving the solution.

5. Debugging & Error Handling:

- Description: Identify, isolate, and fix any errors or unexpected behaviors introduced during implementation to ensure the solution runs reliably.
- Reasoning Guideline: Assign a step to this stage if it is focused on diagnosing problems, investigating error messages, or making corrections to ensure proper functionality.

6. Experimental Refinement:

- Description: Re-run experiments on an already implemented solution and iteratively test various configurations, tune parameters, and compare alternative approaches to upgrade performance.
- Reasoning Guideline: Assign a step to this stage if it involves re-executing or adjusting an implemented solution, making upgrades and modifications to improve performance after the initial implementation has been established.

7. Final Evaluation & Submission:

- Description: Conduct a comprehensive evaluation of the refined solution against benchmarks and prepare the solution for final submission.
- Reasoning Guideline: Assign a step to this stage if it involves performing a final, thorough evaluation of the solution's performance, verifying that all improvements meet the required criteria, and preparing for submission.

Your response must be a JSON object where the keys are the step numbers (as strings) and the values are the corresponding stage numbers (from 1 to 7) that best describe the agent's activity at that step.

IMPORTANT: When assigning a stage, review the steps before and after each step to understand the broader context.

IMPORTANT: The original trace has {original_step_count} steps. Your response MUST contain exactly {original_step_count} keys, numbered from "1" to "{original_step_count}".

Example output format:

```
"1": 1,
"2": 1,
"3": 4,
"4": 6,
"5": 7,
"6": 7,
```

I Prompt for Error Response Annotation

Prompt for LLM Error Response Annotator

Below is a detailed chain-of-thought from an agent after encountering an error message:

```
{error_step}
```

Based on the provided debugging steps, classify the agent response regarding the error as follows:

- 1 -> Fixed the error: The agent identified the issue and implemented a solution that resolved the error.
- 2 -> Tried to fix the error but didn't: The agent attempted to address the error but the fix was not successful.
- 3 -> Didn't even try to fix the error and just went off doing something else: The agent did not directly attempt to resolve the error but instead focused on other tasks unrelated to fixing it.

If the error was fixed (status -> 1), also identify which step number was the error fixed at.

Return a JSON with two fields:

- Status: The number (1, 2, or 3) corresponding to the classification
- FixedAtStep: The step number where the error was fixed (only if Status is 1, otherwise null)

J Impact Statement

The ability of AI agents to perform high-quality ML research could accelerate breakthroughs in critical domains such as healthcare, climate modeling, and AI safety. However, agents capable of autonomously generating novel methods at scale may also amplify risks if their outputs outpace human understanding or oversight. While current agents underperform human researchers, MLRC-BENCH highlights the need for ongoing monitoring of their capabilities to ensure alignment with ethical standards and societal goals. By releasing this benchmark, we aim to enhance transparency and encourage the development of safer, more reliable AI research agents. We caution against deploying such systems without robust safeguards and urge the community to prioritize evaluations that balance innovation with accountability.