# MLE-Dojo: Interactive Environments for Empowering LLM Agents in Machine Learning Engineering

Rushi Qiang<sup>†</sup>; Yuchen Zhuang<sup>†</sup>; Yinghao Li<sup>†</sup>,
Dingu Sagar V K<sup>†</sup>, Rongzhi Zhang<sup>†</sup>, Changhao Li<sup>†</sup>, Ian Shu-Hei Wong<sup>†</sup>,
Sherry Yang<sup>§</sup>, Percy Liang<sup>§</sup>, Chao Zhang<sup>†</sup>, Bo Dai<sup>†</sup>

<sup>†</sup>Georgia Institute of Technology Stanford University

#### **Abstract**

We introduce MLE-Dojo, a Gym-style framework for systematically reinforcement learning, evaluating, and improving autonomous large language model (LLM) agents in iterative machine learning engineering (MLE) workflows. Unlike existing benchmarks that primarily rely on static datasets or single-attempt evaluations, MLE-Dojo provides an interactive environment enabling agents to iteratively experiment, debug, and refine solutions through structured feedback loops. Built upon 200+ real-world Kaggle challenges, MLE-Dojo covers diverse, open-ended MLE tasks carefully curated to reflect realistic engineering scenarios such as data processing, architecture search, hyperparameter tuning, and code debugging. Its fully executable environment supports comprehensive agent training via both supervised fine-tuning and reinforcement learning, facilitating iterative experimentation, realistic data sampling, and real-time outcome verification. Extensive evaluations of eight frontier LLMs reveal that while current models achieve meaningful iterative improvements, they still exhibit significant limitations in autonomously generating long-horizon solutions and efficiently resolving complex errors. Furthermore, MLE-Dojo's flexible and extensible architecture seamlessly integrates diverse data sources, tools, and evaluation protocols, uniquely enabling model-based agent tuning and promoting interoperability, scalability, and reproducibility. We open-source our framework and benchmarks to foster community-driven innovation towards next-generation MLE agents: https://github.com/MLE-Dojo/MLE-Dojo.

# 1 Introduction

Large language models (LLMs) have demonstrated remarkable capabilities across diverse coding tasks, including code generation, debugging, and refactoring [3, 14, 1, 23]. Despite these advances, machine learning engineering (MLE) tasks remain uniquely challenging due to their inherent complexity, specialized domain knowledge, and the extensive iterative experimentation required for developing and optimizing ML algorithms [17, 26, 2]. LLM-based agents hold significant promise for revolutionizing MLE by automating repetitive steps, generating boilerplate code, suggesting suitable algorithms [17, 26, 2].

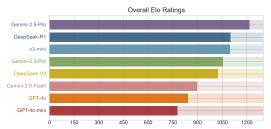


Figure 1: Benchmark evaluations of eight frontiers LLMs across 50 evaluation tasks in MLE-Dojo.

rithms, debugging implementations, and iteratively refining models. Ultimately, advanced LLM-based

<sup>\*</sup>Equal contribution.

Table 1: Summary of existing benchmarks of MLE agents with data resources and sample sizes.

Datasets	Interactive Gym	Executable Environment	_	Extensible Tasks	Flexible Scaffolds	Labiilar	CV	NLP	# Training	# Eval
AutoKaggle [24]	×	X	X	X	X	V	X	X	0	8
MLAgentBench [17]	X	×	×	X	X	<b>/</b>	1	/	0	13
DSBench [21]	X	×	X	X	X	<b>/</b>	/	/	0	74
DSEval [36]	X	<b>✓</b>	X	X	X	<b>/</b>	/	/	0	31
MLEBench [2]	X	<b>✓</b>	×	X	<b>/</b>	<b>/</b>	1	/	0	75
MLGym [26]	<b>✓</b>	<b>✓</b>	×	X	X	<b>✓</b>	1	/	0	13
MLE-Dojo	<b>✓</b>	<b>✓</b>	<b>V</b>	V	V	/	/	/	150	50

MLE agents could autonomously perform tasks such as literature search, hypothesis generation, experimental design, method implementation, result analysis, and even scientific dissemination. However, developing robust and autonomous MLE agents remains in its infancy, largely due to the absence of comprehensive benchmarks, interactive and executable training environments, and standardized evaluation frameworks capable of supporting rigorous, iterative experimentation. Several recent benchmarks have emerged to assess and facilitate progress in LLM-based coding agents [22, 34, 37, 17, 21, 26, 29]. However, most existing efforts focus mainly on isolated tasks such as data analysis or visualization, or competitions with narrow, non-interactive scenarios [22, 34, 37]. Such compartmentalized benchmarks fail to capture the inherent complexity and iterative nature of real-world MLE workflows, which require continuous experimentation, iterative debugging, structured feedback incorporation, and efficient resource management. Although broader benchmarks such as MLE-Bench [17] and DSBench [21] provide more diverse MLE tasks, they still lack interactive environments that support iterative experimentation and training paradigms such as fine-tuning or reinforcement learning, Moreover, in contrast to structured software engineering (SWE) tasks [20, 35], which only require a structured environment for code-based problem-solving and can be readily sourced from general web resources such as GitHub issues, typical MLE tasks necessitate systematically curated datasets and standardized training data, elements currently unavailable in existing MLE benchmarks or gym-style frameworks [2, 26]. This complexity not only increases storage and computational requirements, but also creates significant challenges when scaling to more comprehensive and diverse problem sets.

In this study, we introduce MLE-Dojo, an interactive, large-scale Gym-style environment explicitly designed for developing and benchmarking autonomous LLM agents on real-world MLE workflows. Built upon 200+ real-world Kaggle competitions across critical ML domains, including tabular data analysis, computer vision, natural language processing, and time series forecasting, etc.. MLE-Dojo provides carefully curated, executable tasks that closely replicate real-world engineering challenges. Among these, 150 tasks constitute the initial training dataset for MLE agents and are integrated into an interactive environment, enabling training trajectory sampling for both supervised fine-tuning and reinforcement learning. Moreover, MLE-Dojo incorporates pre-installed dependencies, standardized evaluation scripts, and real-time outcome verification, thus simplifying iterative experimentation and debugging processes. We conduct extensive empirical evaluations involving eight LLMs and multiple agent scaffolds. All evaluation results are publicly accessible via a continuously updated leaderboard, promoting transparent comparison, reproducibility, and collaborative research in this emerging field. Additionally, MLE-Do jo features a modular architecture that decouples agent capabilities from the underlying environment, facilitating seamless integration with diverse tools and data sources, enhancing interoperability, scalability, and supporting a robust ecosystem for developing generalizable MLE agents. Our main contributions are summarized as follows:

- Comprehensive Framework and Benchmark: We establish MLE-Dojo as a comprehensive and large-scale benchmark consisting of over 200 Kaggle MLE competitions, enabling systematic and rigorous evaluations of autonomous LLM agents.
- Interactive and Executable Environment: MLE-Dojo provides an interactive and fully executable Gym-style environment that facilitates iterative experimentation, including comprehensive training trajectories sampling for supervised fine-tuning and reinforcement learning.
- Advanced Functionalities and Scalability Supports: MLE-Dojo uniquely facilitates outcome verification, MLE agent tuning and seamless integration of diverse datasets and tools, significantly accelerating the development of robust, generalizable, and scalable MLE agents.

• Extensive Empirical Evaluation and Public Leaderboard: We conduct large-scale evaluations across multiple state-of-the-art LLMs and agent scaffolds, with results publicly available through an actively maintained long-term real-time leaderboard to foster community-driven innovation.

In the remainder of the paper, we discuss the related work in section 2 and thoroughly describe MLE-Dojo in section 3. We present the experimental setup and results of LLMs as MLE agents in section 4 and conclude the paper in section 5.

# 2 Related Works

**X-Agent.** LLM-based coding agents have significantly advanced automated software development, with impactful applications across two primary domains: SWE and MLE. SWE agents [33, 8, 32] excel at code generation, debugging, and refactoring for improved maintainability and performance. For example, SWE-agent [33] operates within a constrained agent-computer interface to facilitate file creation, repository navigation, and code testing; Magentic-One [8] expands the capability of SWE agents with web navigation features, further enhancing their practicality and usability. OpenHands [32] introduces an SWE agent with sandboxed environments to ensure safe command execution and verifiable web browsing, facilitating standardized benchmarking. MLE agents [15, 13, 24, 12, 16, 4, 19, 31], on the other hand, specialize in constructing, optimizing, and deploying comprehensive ML pipelines, including automatic architecture selection and hyperparameter tuning tailored to specific data characteristics. Early contributions focused on automating individual steps in ML engineering pipelines; for instance, CAAFE [15] leverages LLMs to iteratively generate semantically meaningful features for tabular data based on dataset description. More recent works have progressively expanded automation from single steps to encompass end-to-end ML workflows. For example, DS-Agent [13] integrates a case-based reasoning approach to retrieve, reuse, evaluate, and refine solutions based on historical experiences. AutoKaggle [24] develops iterative code execution, debugging, and comprehensive unit testing to ensure code correctness and logic consistency. Agent K v1.0 [12] provides an end-to-end autonomous data science agent designed to optimize workflows across diverse data science scenarios, while DataInterpreter [16] employs hierarchical graph modeling for complex problem decomposition and programmable node generation that refines and verifies subproblems via code generation. Advanced search and optimization approaches have further enhanced the capabilities of MLE agents. SELA [4] employs Monte Carlo Tree Search to enhance strategic planning during automated ML task-solving. AIDE [19] frames MLE as a code optimization problem, formulating trial-and-error as a tree search in the space of potential solutions. Here, we focus on evaluating and improving MLE agents.

**X-Bench.** Benchmarks for MLE agents not only provide foundational baselines but also inaugurate public leaderboards. As illustrated in Table 1, early benchmarks primarily targeted isolated actions, such as data visualization and simple data manipulation [22, 34, 37]. Recent benchmarks expand the scope of evaluations to encompass broader project-level tasks, often sourced from real-world Kaggle competitions. For example, AutoKaggle [24] assesses static LLM workflows across 8 tabular Kaggle competitions. MLAgentBench [17] incorporates 13 tasks from Kaggle and bespoke ML challenges, providing baseline solutions and evaluating how often agents can achieve at least 10% improvement over these baselines. However, such benchmarks often exhibit limitations in both task breadth and complexity. Benchmarks such as MLE-Bench [2] and DSBench [21] significantly broaden task diversity to 75 and 74 competitions, respectively, aiming to closely simulate real-world MLE scenarios. However, these platforms still lack robust interactive environments that support iterative experimentation, agent fine-tuning, and realistic training scenarios, critical components for effective development and evaluation of MLE agents. In contrast, MLE-Dojo combines comprehensive MLE task coverage with a fully interactive execution environment, offering unparalleled opportunities for developing, benchmarking, and refining advanced MLE agents.

**X-Gym (Dojo).** Interactive environments are crucial for evaluating and refining LLM-based agents. SWE-Gym [28] introduces an interactive environment specifically tailored for software engineering tasks, enabling extensive iterative training and validation. BrowserGym [7] facilitates the evaluation of web-based navigation and interaction tasks, capturing complex real-world knowledge work scenarios. Similarly, Collaborative-Gym [30] expands agent capabilities by enabling asynchronous, tripartite interaction among agents, humans, and task environments. The most recent and relevant work, ML-Gym [26], further integrates diverse ML research tasks into an interactive gym framework,

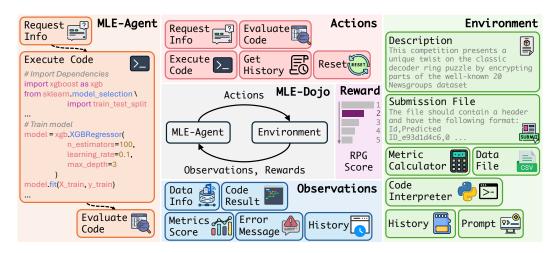


Figure 2: Overview of MLE-Dojo. The framework bridges MLE-Agents with MLE task environments through standardized interfaces for observation and action spaces.

providing an environment that supports reinforcement learning and iterative experimentation across open-ended research tasks. MLE-Dojo significantly extends both the quantity and complexity of tasks, offering 200+ competitions with an additional dedicated training set of 150 tasks. MLE-Dojo provides a concise and unified task format and interface to facilitate seamless task expansion, along with a flexible and user-friendly environment interface to support interactions between various agents and environments. It also offers meaningful and precise reward feedback, which enables effective trajectory sampling for supervised fine-tuning and reinforcement learning, thereby setting a new standard for comprehensive evaluation and training of MLE agents.

# 3 MLE-Dojo

We propose MLE-Dojo, a Gym-like framework that provides a standardized environment to evaluate MLE agents interacting with task-specific environments. As depicted in Figure 2, MLE-Dojo provides a unified interface for MLE agents, *i.e.*, LLM-based assistants that can write code to handle project-level ML tasks, such as completing data science competitions. Each task environment contains essential information, including datasets, evaluation metrics, analysis results, code execution outcomes, error messages, and interaction history, facilitating comprehensive agent-environment interactions. The agent usually takes actions to solve the task, such as requesting task information, writing and validating code, evaluating code, retrieving past history and so on.

From an agent's perspective, the environment is typically centered on a sampled MLE problem  $p \in \mathcal{P}$ , where  $\mathcal{P}$  is the task space. Each interaction between the environment and the agent can be formalized as a Partially Observable Markov Decision Process (POMDP), where at the time step t, the environment provides the current observation  $o_t \in \mathcal{O}$  and the reward  $r_t \in \mathcal{R}$ . Depending on the specific agent and environment design, the agent generates the next action  $a_{t+1} \in \mathcal{A}$  at step t+1 based on the history of prior interactions. Here,  $\mathcal{O}$ ,  $\mathcal{A}$ , and  $\mathcal{R}$  represent the observation space, action space, and reward space, respectively. The resulting interaction loop<sup>2</sup> is shown in Figure 3.

#### 3.1 Modular and User-Friendly Interface

The MLE-Dojo framework is designed with a strong emphasis on modularity, flexibility, and extensibility. Each component within the environment operates both independently and collaboratively and is fully decoupled, allowing seamless integration and extension through unified register mechanism. Specifically, we define the following core modules, each of which encapsulates a distinct aspect of the environment:

**Error:** It encodes a comprehensive hierarchy of error types, enabling fine-grained debugging and facilitating informative feedback from the environment to the user; **Interface:** It governs the execu-

<sup>&</sup>lt;sup>2</sup>For clarity and visual appeal, we present the core logic of the code in an intuitive manner. For detailed implementation, please refer to our code repository.

tion and interaction logic of native environment actions, serving as the backbone of agent-environment communication; **Feedback**: It translates interaction outcomes into structured, interpretable feedback to guide agent behavior; **Metric**: It defines a general metric base class, which can be subclassed to implement competition-specific evaluation metrics in a standardized and reusable manner.

Users can interact with the environment either by developing custom modules via the provided APIs or by directly using the pre-designed environment for seamless experimentation, requiring only a single call to env.step(action\_type, action\_args). The minimal and intuitive interaction logic sig-



Figure 3: Interaction loop in MLE-Dojo with theoretical model (left) and concrete Python API (right).

nificantly lowers the barrier for agent customization and development. We have also integrated agents of several scaffolds into MLE-Dojo, as detailed in the Appendix K. This design not only ensures clear separation of concerns, but also supports the easy incorporation of new functionalities, promoting both research reproducibility and rapid prototyping across diverse evaluation settings.

# 3.2 Extensible Task Space P

MLE-Dojo comprises more than 200 carefully selected tasks spanning diverse ML domains (Figure 11), including tabular data analysis, computer vision, natural language processing, etc. Each primary category further encompasses multiple specific task types—for instance, image classification within computer vision and sentiment analysis within natural language processing—yielding various distinct task types. The majority of these tasks represent complex, practi-

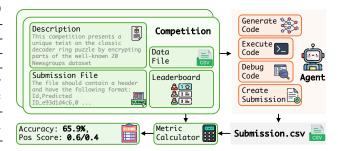


Figure 4: Overview of data structure in MLE-Dojo.

cal MLE challenges essential for developing robust LLM-based autonomous agents, sourced from real-world Kaggle<sup>3</sup> competitions.

We aggregate tasks from multiple sources, including 68 competitions from MLE-Bench [2] (excluding 7 tasks that are unavailable, excessively large, or tightly coupled with specific packages), 74 from DSBench [21], and 75 additionally downloaded through Kaggle's official API and carefully prepared. After removing duplicate entries across sources, we obtain a diverse collection of over 200 unique tasks. Following MLE-Bench [2], each task is standardized into a consistent data structure (Figure 4 and 10) to create a unified interface suitable for MLE-Dojo environment and LLM agents.

We partition MLE-Dojo into training and evaluation subsets with a 150:50 split. The evaluation set prioritizes tasks commonly referenced in prior benchmarks [2, 21], supplemented by strategically selected additional competitions to ensure domain diversity and representativeness. The larger training set provides extensive interactive experiences for LLM agents, facilitating robust agent training through both supervised fine-tuning and reinforcement learning.  $^4$ 

MLE-Dojo provides a standardized data interface and detailed instruction, simplify extending MLE-Dojo with new tasks (in a plug-and-play manner). A unified competition format is maintained for each task, including (1) *Basic information*: competition description and sample submission, (2) *Well-structured datasets*: well-split datasets through a general prepare.py script (most competitions do not originally provide labeled test data), and (3) *Leaderboards*: both public and private leaderboards assessed using competition-specific evaluation metrics. Following this standardized

<sup>&</sup>lt;sup>3</sup>Kaggle is a widely used online platform that hosts competitive ML tasks. Typically, a competition requires the submission of a prediction file, evaluated with a specific metric to compute the corresponding score.

<sup>&</sup>lt;sup>4</sup>We are continuing to extend our dataset to include more diverse, challenging and new tasks.

data format, users can easily incorporate their own competitions for testing purposes. Comprehensive format definitions and detailed guidance are available in Appendix C.

# 3.3 Observation Space $\mathcal{O}$

MLE-Dojo provides a rich observation space with five main components, enabling comprehensive, clear, and structured feedback during Gym-like interaction:

**Task Information.** For each competition, MLE-Dojo supplies comprehensive information required to solve tasks, including *Competition background*, *Goal description*, *Sample submission* and *Data structure*. This provides all necessary information a Machine Learning Engineer requires to solve the task, while avoiding the introduction of additional assistance.

**Code Execution Results.** Agents generate code as solutions for MLE tasks. MLE-Dojo provides detailed code execution results in a concise feedback format. Execution Results may include detailed error messages to facilitate the agent in identifying issues and performing debugging; they may also contain code execution outputs, which serve as an important source of information for the agent to deeply explore and analyze the data, gain insights of designed algorithm, etc.

**Evaluation Metric Scores.** For each task, MLE-Dojo implements well-defined evaluation functions, providing description-aligned criteria and metrics to validate submission format, compute specific metric scores and evaluate method performance. Performance is reported as both *raw scores* and *HumanRank score* (discussed in Section 3.5), offering absolute and relative performance insights.

**Interaction History.** In MLE-Dojo, MLE agents can access interaction histories in two ways: (1) *Conversation history*: Obtained from the agent side, including all LLM generations and environment observations, and (2) *Environment records*: Obtained from the environment side, objectively recording every agent action and environment observation. This dual approach accommodates different LLM agents and MLE task implementations, improving compatibility with both dialog-based LLMs and coding-focused agents. For reasoning models, conversation history effectively records extended reasoning procedures, while environment records better suit coding LLMs that may not support natural language interaction.

#### 3.4 Expandable Action Space A

The fundamental action space of MLE-Dojo consists of executable Python code. Agents can call request\_info to query task descriptions and dataset details, then generate Python code to solve specific MLE tasks using four core functions: validate\_code, execute\_code, get\_history and reset. Agents can utilize validate\_code to test code or capture useful information from code outputs in a light-weight manner. Through execute\_code, agent execute complete code including submission generation, metric evaluation and reward feedback. Agents must explicitly use "execute\_code" command to generate valid submission files and get scores; submissions are not automated. Additionally, agents can leverage the get\_history function to retrieve past interactions between the environment and themselves. MLE-Dojo also provides a user-defined action portal, offering an extensible action space for researchers to register additional actions. As long as new actions are registered through the portal and provided to the agent in context, the agent can learn to leverage these new actions to solve problems. We provide the details of our action space in Appendix E.

# 3.5 Reward Space and Environmental Feedback

MLE-Dojo introduces a reward mechanism specifically designed to reflect the quality of solutions. This mechanism utilizes quantitative metrics provided by each competition, such as accuracy, F1-score, etc. However, since competitions differ in evaluation metrics and score ranges, directly using absolute performance to evaluate LLM-based agents' generated code is not precise. To provide continuous and fine-grained feedback that rewards incremental improvements and innovative solutions, we propose using the relative position of the human leaderboard as reward instead of the coarse-grained medals used in existing benchmarks [2].

**HumanRank Score.** We calculate the relative position score of the current submission on the leaderboard of human competitors. Suppose that the submission ranks at position p among a total of N submissions on the leaderboard. Then, the position score is computed as:  $s = 1 - \frac{p}{N}$ .

Table 2: Main experiments of LLMs as MLE Agents on MLE tasks in MLE-Dojo.

Task Categories $(\rightarrow)$		MLE-Lite			Tabular			NLP			CV	
Models $(\downarrow)$	AUP↑	H-Rank $(\%)^{\uparrow}$	Elo <sup>↑</sup>	<b>AUP</b> <sup>↑</sup>	H-Rank $(\%)^{\uparrow}$	Elo <sup>↑</sup>	<b>AUP</b> <sup>↑</sup>	H-Rank $(\%)^{\uparrow}$	Elo <sup>↑</sup>	AUP↑	H-Rank $(\%)^{\uparrow}$	Elo <sup>↑</sup>
gpt-4o-mini [18]	1.492	21.21	753	0.724	15.37	839	0.837	13.14	758	1.172	10.17	754
gpt-4o [18]	1.448	27.85	830	0.691	18.97	861	0.842	29.97	903	1.418	18.76	761
o3-mini [27]	1.895	56.48	1108	0.739	32.65	1019	0.992	37.46	1056	1.892	35.02	1207
DeepSeek-v3 [25]	1.825	44.26	1004	0.727	37.85	1015	0.977	28.41	1028	1.784	26.75	1067
DeepSeek-r1 [6]	1.852	58.43	1137	0.678	38.13	1053	0.988	28.48	1103	1.844	34.26	1083
Gemini-2.0-Flash [9]	1.696	33.50	847	0.689	30.36	923	0.884	28.39	860	1.607	20.35	978
Gemini-2.0-Pro [10]	1.796	48.61	1064	0.787	37.46	1139	0.970	30.93	1028	1.651	23.07	973
Gemini-2.5-Pro [11]	1.919	61.95	1257	0.798	42.64	1150	0.998	38.45	1266	1.915	42.83	1177

The HumanRank score indicates the percentage of human competitors the agent surpasses on the leaderboard for a given competition. A higher score reflects stronger performance relative to human participants. To prevent bias between public and private leaderboards, we compute the relative scores on each leaderboard independently and then use their average as the final reward.

We adopt the HumanRank Score as the reward in our environment. It is fully aligned with the original performance metrics-achieving a higher original score is strictly positively correlated with obtaining a higher HumanRank Score, regardless of the specific metric used. Moreover, HumanRank is a normalized score within the [0, 1] range, which resolves the issue of varying score magnitudes across different tasks and enables it to serve as a unified and informative reward.

# 4 LLMs as MLE Agents in MLE-Dojo

#### 4.1 Experiment Setups

**Backbone LLMs.** We consider different backbones to test the effectiveness of MLE-Dojo in the evaluation and improvement of LLMs as MLE Agents. MLE Agent leverages native actions and interacts with the environment through a straightforward logic. See implementation and prompt details in Appendix I and Appendix K.1. Specifically, we consider gpt-4o-mini (2024-07-18) [18], gpt-4o (2024-11-20) [18], o3-mini (2025-01-31) [27] from OpenAI, Gemini-2.0-Flash [9], Gemini-2.0-Pro (exp) [10], and Gemini-2.5-Pro (exp-3-25) [11] from Google, and DeepSeek-v3 (2025-03-24) [25] and DeepSeek-r1 [6] from DeepSeek as evaluation backbone LLMs. For non-reasoning models, we set temperature=0.0 and top-p=1.0 to ensure reproducible evaluations. We take the best performance of two runs per task per model.

**Evaluation Metrics.** To ensure a comprehensive evaluation, we consider *Area Under the Performance Profile (AUP)* [26], *HumanRank Score (H-Rank, %)*, and *Elo* ranking [5] together as metrics. Additional implementation details are available in appendix **F**.

**Environment Configurations.** Prompt, history, time, and memory configurations are designed to rigorously evaluate LLMs' capabilities in long-context handling, instruction-following, reasoning, and coding under resource-constrained conditions, closely mirroring realistic Machine Learning Engineering scenarios. We do not restrict the number of submission attempts to enable continuous improvement. See Appendix G for detailed configurations.

### 4.2 Main Results

Table 2 presents a comprehensive evaluation of eight LLMs as MLE Agents across 50 evaluation tasks. Reasoning models such as o3-mini, DeepSeek-r1, and Gemini-2.5-Pro consistently achieve high rankings across all metrics, demonstrating strong adaptability and effectiveness as MLE Agents. Additionally, Figure 6 further illustrates the Performance Profiles along with the corresponding AUP curves. Models like Gemini-2.0-Pro exhibit balanced performance profiles, achieving moderate but consistently reliable results across various tasks. This comprehensive comparative perspective highlights the strengths and limitations of each model, offering practical insights into their suitability for different MLE scenarios. Specifically, the HumanRank Score provides an absolute measure of performance incorporating human benchmarks, the Elo Score clarifies competitive relationships

through pairwise analyses, and the Performance Profiles with AUP scores assess robustness and consistency across performance variations. Together, these evaluation approaches in MLE-Dojo form a robust, multifaceted perspective on the capabilities and limitations of LLMs as MLE agents, enhancing the overall interpretation and reliability of our experimental conclusions.

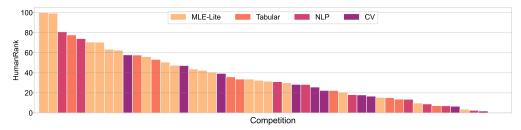


Figure 5: Task difficulty sorted by average HumanRank score.

Furthermore, we define the difficulty level of different tasks with the average performance of different models in comparison with the human leaderboard. Figure 5 illustrates the average performance distribution across 8 frontier models on the tasks. As shown in the figure, CV tasks are the most challenging-none of them have an average HumanRank score above 60, and more than half fall below 30. For MLE-Lite tasks, the average HumanRank scores mostly exceed 30. Difficulty levels of tasks in other domains are more evenly distributed.

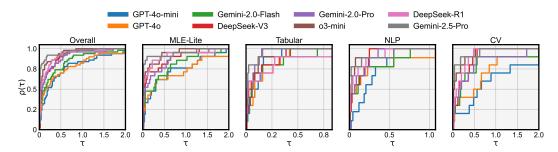


Figure 6: Performance profiles and corresponding AUP curves for evaluating the robustness of LLMs across four ML tasks. The x-axis represents the performance ratio threshold  $\tau$ , while the y-axis indicates the fraction of tasks for which a model achieves performance within a factor  $\tau$  of the best-performing model.

#### 4.3 Cost Analysis

Figure 7 illustrates the cost-performance relationship across different LLMs and task categories. Reasoning models (*e.g.*, DeepSeek-r1) typically incur higher costs due to their premium pricing structures and longer solution outputs. Even reasoning models with comparatively lower pricing, such as o3-mini, tend to produce longer outputs due to more complex reasoning processes. Notably, tasks involving computer vision and deep neural network training pipelines consistently generate longer codes compared to classical ML tasks (*e.g.*, tabular analysis) executed on CPUs. While cost generally correlates with solution complexity and token consumption, some models, such as DeepSeek-r1, achieve competitive performance with significantly fewer tokens, highlighting potential cost-efficiency opportunities. We further investigate both history length and solution length of different models in Appendix H.2.

# 4.4 Step-wise Performance Dynamics

Figure 8 presents step-wise performance dynamics across different models. Among reasoning models, o3-mini consistently achieves high performance within the initial steps (typically within the first five) and maintains stable scores in subsequent steps. Conversely, DeepSeek-r1 and Gemini-2.5-Pro exhibit gradual improvements, achieving comparable or superior performance in intermediate to later steps. Non-reasoning models occasionally outperform reasoning models at early or intermediate steps but generally show limited improvement as steps progress, resulting in lower final scores.

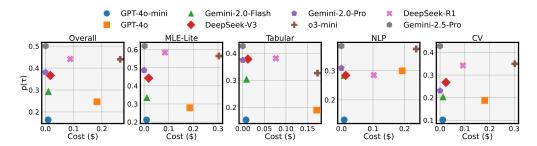


Figure 7: Relationship between average computational cost and performance across evaluated LLMs and task categories. Each point represents the average cost per task, with specific attention given to reasoning vs. non-reasoning model cost dynamics. Note that Gemini-2.0-Pro and Gemini-2.5-Pro are excluded from cost analysis (only for performance reference) due to current free usage.

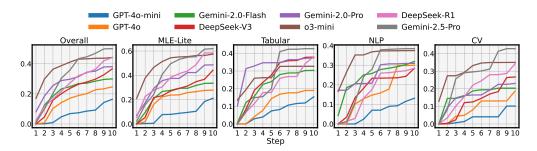


Figure 8: Step-wise HumanRank performance comparisons between reasoning and non-reasoning models, considering only code execution and validation steps. Information requests are excluded.

#### 4.5 Error Analysis

Figure 9 shows the average failure rates across tasks, highlighting validation, execution, and overall errors. Gemini-2.5-Pro maintains the lowest overall failure rate, aligning with its consistently high performance. In contrast, DeepSeek-r1, despite achieving strong performance, experiences relatively high failure rates in both execution and validation categories. gpt-40

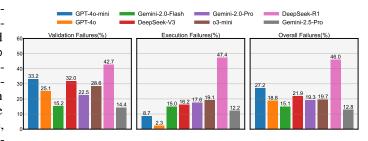


Figure 9: Average failure rates across tasks.

and gpt-4o-mini, due to their conservative strategies, achieve low execution failure rates but face comparatively high validation failures. Conversely, Gemini-2.0-Flash successfully balances aggressive execution with one of the lowest overall failure rates, second only to Gemini-2.5-Pro. We further investigate excution error types in Appendix H.3.

#### 5 Conclusion

We introduced MLE-Dojo, a Gym-style framework designed for training, evaluating, and benchmarking autonomous MLE agents. By providing an interactive and realistic environment built upon a large-scale collection of real-world Kaggle competitions, MLE-Dojo enables systematic experimentation, rigorous outcome verification, and iterative feedback crucial for advancing LLM-driven MLE workflows. Our extensive empirical evaluations established foundational baselines and highlighted both capabilities and critical limitations of current state-of-the-art models and agent architectures. We publicly release our framework, benchmarks, and leaderboard to encourage transparent comparison, reproducible research, and community-driven progress toward next-generation of fully autonomous MLE agents. Future work includes expanding MLE-Dojo to incorporate domain-specific deep research and support for multi-agent collaborative scenarios.

# Acknowledgments and Disclosure of Funding

We thank the anonymous reviewers and area chairs for their valuable feedback. This work was supported in part by ONR (N000142512173), NSF ECCS (2401391), and NSF IIS (2403240).

#### References

- [1] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- [2] J. S. Chan, N. Chowdhury, O. Jaffe, J. Aung, D. Sherburn, E. Mays, G. Starace, K. Liu, L. Maksin, T. Patwardhan, et al. Mle-bench: Evaluating machine learning agents on machine learning engineering. *arXiv* preprint arXiv:2410.07095, 2024.
- [3] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. D. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, et al. Evaluating large language models trained on code. *arXiv* preprint arXiv:2107.03374, 2021.
- [4] Y. Chi, Y. Lin, S. Hong, D. Pan, Y. Fei, G. Mei, B. Liu, T. Pang, J. Kwok, C. Zhang, et al. Sela: Tree-search enhanced llm agents for automated machine learning. *arXiv* preprint *arXiv*:2410.17238, 2024.
- [5] W.-L. Chiang, L. Zheng, Y. Sheng, A. N. Angelopoulos, T. Li, D. Li, B. Zhu, H. Zhang, M. Jordan, J. E. Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. In *Forty-first International Conference on Machine Learning*, 2024.
- [6] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.
- [7] A. Drouin, M. Gasse, M. Caccia, I. H. Laradji, M. D. Verme, T. Marty, D. Vazquez, N. Chapados, and A. Lacoste. Workarena: How capable are web agents at solving common knowledge work tasks? In *Forty-first International Conference on Machine Learning*, 2024.
- [8] A. Fourney, G. Bansal, H. Mozannar, C. Tan, E. Salinas, F. Niedtner, G. Proebsting, G. Bassman, J. Gerrits, J. Alber, et al. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468*, 2024.
- [9] Google. Gemini 2.0 flash: Our powerful workhorse model with low latency and enhanced performance, built to power agentic experiences. *Google DeepMind Blog*, 2025.
- [10] Google. Gemini 2.0 pro: Our best model yet for coding performance and complex prompts. *Google DeepMind Blog*, 2025.
- [11] Google. Gemini 2.5: Our most intelligent ai model. *Google Blog*, 2025.
- [12] A. Grosnit, A. Maraval, J. Doran, G. Paolo, A. Thomas, R. S. H. N. Beevi, J. Gonzalez, K. Khandelwal, I. Iacobacci, A. Benechehab, et al. Large language models orchestrating structured reasoning achieve kaggle grandmaster level. *arXiv preprint arXiv:2411.03562*, 2024.
- [13] S. Guo, C. Deng, Y. Wen, H. Chen, Y. Chang, and J. Wang. DS-agent: Automated data science by empowering large language models with case-based reasoning. In *Forty-first International Conference on Machine Learning*, 2024.
- [14] D. Hendrycks, S. Basart, S. Kadavath, M. Mazeika, A. Arora, E. Guo, C. Burns, S. Puranik, H. He, D. Song, and J. Steinhardt. Measuring coding challenge competence with APPS. In Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2), 2021.
- [15] N. Hollmann, S. Müller, and F. Hutter. Large language models for automated data science: Introducing caafe for context-aware automated feature engineering. *Advances in Neural Information Processing Systems*, 36, 2024.

- [16] S. Hong, Y. Lin, B. Liu, B. Liu, B. Wu, C. Zhang, C. Wei, D. Li, J. Chen, J. Zhang, et al. Data interpreter: An Ilm agent for data science. arXiv preprint arXiv:2402.18679, 2024.
- [17] Q. Huang, J. Vora, P. Liang, and J. Leskovec. MLAgentbench: Evaluating language agents on machine learning experimentation. In *Forty-first International Conference on Machine Learning*, 2024.
- [18] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [19] Z. Jiang, D. Schmidt, D. Srikanth, D. Xu, I. Kaplan, D. Jacenko, and Y. Wu. Aide: Ai-driven exploration in the space of code. *arXiv preprint arXiv:2502.13138*, 2025.
- [20] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- [21] L. Jing, Z. Huang, X. Wang, W. Yao, W. Yu, K. Ma, H. Zhang, X. Du, and D. Yu. DSBench: How far are data science agents from becoming data science experts? In *The Thirteenth International Conference on Learning Representations*, 2025.
- [22] Y. Lai, C. Li, Y. Wang, T. Zhang, R. Zhong, L. Zettlemoyer, W.-t. Yih, D. Fried, S. Wang, and T. Yu. Ds-1000: A natural and reliable benchmark for data science code generation. In *International Conference on Machine Learning*, pages 18319–18345. PMLR, 2023.
- [23] Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- [24] Z. Li, Q. Zang, D. Ma, J. Guo, T. Zheng, M. Liu, X. Niu, Y. Wang, J. Yang, J. Liu, et al. Autokaggle: A multi-agent framework for autonomous data science competitions. *arXiv* preprint arXiv:2410.20424, 2024.
- [25] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, et al. Deepseek-v3 technical report. arXiv preprint arXiv:2412.19437, 2024.
- [26] D. Nathani, L. Madaan, N. Roberts, N. Bashlykov, A. Menon, V. Moens, A. Budhiraja, D. Magka, V. Vorotilov, G. Chaurasia, D. Hupkes, R. S. Cabral, T. Shavrina, J. Foerster, Y. Bachrach, W. Y. Wang, and R. Raileanu. Mlgym: A new framework and benchmark for advancing ai research agents, 2025.
- [27] OpenAI. Openai o3-mini: Pushing the frontier of cost-effective reasoning. OpenAI Blog, 2025.
- [28] J. Pan, X. Wang, G. Neubig, N. Jaitly, H. Ji, A. Suhr, and Y. Zhang. Training software engineering agents and verifiers with swe-gym. *arXiv preprint arXiv:2412.21139*, 2024.
- [29] S. Schmidgall and M. Moor. Agentrxiv: Towards collaborative autonomous research. *arXiv* preprint arXiv:2503.18102, 2025.
- [30] Y. Shao, V. Samuel, Y. Jiang, J. Yang, and D. Yang. Collaborative gym: A framework for enabling and evaluating human-agent collaboration. *arXiv* preprint arXiv:2412.15701, 2024.
- [31] P. Trirat, W. Jeong, and S. J. Hwang. Automl-agent: A multi-agent llm framework for full-pipeline automl. *arXiv preprint arXiv:2410.02958*, 2024.
- [32] X. Wang, B. Li, Y. Song, F. F. Xu, X. Tang, M. Zhuge, J. Pan, Y. Song, B. Li, J. Singh, et al. Openhands: An open platform for ai software developers as generalist agents. *arXiv* preprint *arXiv*:2407.16741, 2024.
- [33] J. Yang, C. E. Jimenez, A. Wettig, K. Lieret, S. Yao, K. R. Narasimhan, and O. Press. Swe-agent: Agent-computer interfaces enable automated software engineering. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

- [34] P. Yin, W.-D. Li, K. Xiao, A. Rao, Y. Wen, K. Shi, J. Howland, P. Bailey, M. Catasta, H. Michalewski, et al. Natural language to code generation in interactive data science notebooks. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 126–173, 2023.
- [35] D. Zan, Z. Huang, W. Liu, H. Chen, L. Zhang, S. Xin, L. Chen, Q. Liu, X. Zhong, A. Li, et al. Multi-swe-bench: A multilingual benchmark for issue resolving. *arXiv preprint arXiv:2504.02605*, 2025.
- [36] Y. Zhang, Q. Jiang, X. Han, N. Chen, Y. Yang, and K. Ren. Benchmarking data science agents. *arXiv preprint arXiv:2402.17168*, 2024.
- [37] Y. Zhang, Y. Pan, Y. Wang, and J. Cai. Pybench: Evaluating llm agent on various real-world coding tasks. *arXiv preprint arXiv:2407.16732*, 2024.

# **A** Limitations and Broader Impacts

#### A.1 Limitations

**Resource Limitations.** While MLE-Dojo provides a rich, interactive environment to train and evaluate LLM agents, it also imposes computational and storage demands. Executing full MLE pipelines-including data preprocessing, model training, and debugging-across hundreds of competition tasks requires potential API credits, CPU/GPU resources and large-scale disk storage. Although we provide competitions of varying difficulty levels and data sizes—organized in ascending order of size for user clarity—comprehensively training or evaluating MLE agents still demands substantial computational and storage resources.

**Data Privacy and Licensing.** MLE-Dojo is constructed atop real-world Kaggle competitions, each of which may be governed by different data usage licenses and privacy policies. Although the framework does not redistribute any proprietary data directly, it facilitates automated downloading and usage of datasets via the official Kaggle API. Users are solely responsible for reviewing, understanding, and adhering to the license agreements associated with each dataset. To assist in compliance, we provide direct links to the license terms and competition rules for all included tasks in https://github.com/MLE-Dojo/MLE-Dojo/blob/main/prepare/licenses.json.

# A.2 Broader Impacts

**Potential Positive Societal Impacts.** MLE-Dojo aims to advance the development of intelligent, autonomous agents capable of supporting and accelerating the machine learning engineering process. By enabling automation of repetitive and error-prone engineering tasks, the framework may help reduce the barrier to entry for ML practitioners, democratize access to model development tools, and empower domain experts (e.g., in healthcare, education, and scientific research) to more effectively apply machine learning solutions without requiring deep expertise in engineering workflows. Furthermore, the open-source nature of MLE-Dojo fosters reproducibility, transparency, and community-led innovation in the design and evaluation of LLM agents for real-world MLE tasks.

Potential Negative Societal Impacts. Agents developed using MLE-Dojo could inadvertently learn and amplify biases from training datasets (such as the Kaggle challenges used), leading to unfair or discriminatory model outcomes if deployed without careful auditing. The autonomous generation of complex MLE code by these agents may also introduce security vulnerabilities if not subjected to rigorous review. Additionally, the powerful capabilities for automating model development, if not responsibly governed, could be misused for creating harmful or privacy-invasive applications, or could lower the barrier for actors to develop such systems. Mitigation strategies should include incorporating fairness metrics, robustness checks, and security evaluations within the MLE-Dojo framework itself, promoting ethical guidelines and best practices for agent development and deployment, ensuring mechanisms for robust human oversight, and fostering research into privacy-preserving and explainable agent behaviors.

#### A.3 Ethical Statements

**Data Usage and Consent.** MLE-Dojo does not host or redistribute any proprietary datasets. All tasks are built upon publicly available Kaggle competitions, and datasets are accessed via the official Kaggle API in compliance with their respective licenses. We ensure that no personal or sensitive information is retained within the framework. Users are explicitly informed that it is their responsibility to read and adhere to the terms of use and data handling policies provided by each competition. Where applicable, competitions selected for inclusion in the benchmark are reviewed to ensure they do not contain personally identifiable or ethically sensitive data.

Fair Use and Responsible Deployment. MLE-Dojo is intended solely for academic research, educational purposes, and the responsible development of AI systems. We discourage the use of this framework for automating competition participation without appropriate attribution or for circumventing fair use policies. We also emphasize the importance of transparent model development practices when using MLE-Dojo, particularly in scenarios where LLM agents are intended for downstream deployment or integration into critical systems.

Bias and Representational Fairness. While MLE-Dojo includes a broad spectrum of tasks across multiple ML domains, the selection of tasks is constrained by the availability of open competitions and may not fully represent all application domains or demographic contexts. Moreover, since LLM agents can inherit and amplify biases present in training data or model pretraining corpora, it is crucial that users critically evaluate outputs, particularly in sensitive contexts (e.g., healthcare, hiring, or education). Future versions of MLE-Dojo will explore the inclusion of bias detection and mitigation tools to support ethical agent development.

**Open-Source Commitment.** We commit to maintaining MLE-Dojo as an open-source project under a permissive license, enabling transparent inspection, community-driven development, and reproducibility. We welcome contributions that enhance the ethical robustness of the framework, including the addition of fairness metrics, safety guardrails, and improved documentation on responsible usage.

#### **B** Disclaimer

The dataset is made available exclusively for educational and academic research purposes, with the intention of advancing scholarly investigations in relevant domains. Users of the dataset must adhere to the following terms and conditions:

- **Data Source and Accuracy**: Although reasonable efforts have been undertaken to curate and organize the dataset, the providers do not warrant the accuracy, completeness, or currency of the information. Users are strongly encouraged to independently verify the data and bear full responsibility for any analyses, interpretations, or conclusions derived from it.
- Usage Restrictions: This dataset is strictly limited to non-commercial use. Any commercial exploitation, product development, or profit-oriented application based on this dataset requires prior explicit written authorization from the dataset providers.
- **Privacy and Legal Compliance**: Users must ensure that their usage of the dataset complies with all relevant legal frameworks, particularly those concerning data privacy, protection, and security. The dataset providers shall not be held accountable for any legal liabilities or consequences resulting from improper or unauthorized usage.
- Non-Infringement of Rights: The dataset includes pre-processed content derived from external sources and is distributed solely for non-commercial research purposes. The dataset providers do not assert ownership over the original data and expressly acknowledge the rights of the original creators. It is the user's responsibility to ensure that their use of the dataset does not violate any copyright or intellectual property laws.
- **Disclaimer of Liability**: The dataset is provided "as is" without any express or implied warranties. The dataset providers shall not be liable for any direct, indirect, incidental, or consequential damages arising from the use of the dataset, including but not limited to financial loss, data misinterpretation, or third-party claims.

# C Unified Data Structure

Users can flexibly incorporate new tasks originating from diverse sources, which can be standardized into a unified competition format with minimal overhead. These personalized competitions can then be seamlessly integrated into our environment for evaluation, benchmarking and training.

To ensure clarity, modularity, and scalability, each competition is organized under a unified and minimalistic directory structure. As shown in Figure 10, the root directory is composed of three major subdirectories: data/, utils/, and an optional info/ folder. The data/ directory is further divided into private/ and public/ subfolders, corresponding respectively to hidden and accessible phases of the competition. Each contains standardized files such as test\_answer.csv, leaderboard.csv, and input-output formats (e.g., sample\_submission.csv, description.txt), thereby facilitating consistent evaluation procedures. The utils/ folder encapsulates task-specific scripts for data preparation and metric evaluation (i.e., prepare.py, metric.py). Lastly, the optional info/ directory provides supplementary metadata such as web links, original descriptions, and data schema. This design adheres to a uniform organizational paradigm that simplifies integration, supports automation, and improves the transparency of competition configurations.

```
competition/
   data/
      private/
        test_answer.csv
        private_leaderboard.csv (Optional)
        public_leaderboard.csv (Optional)
      public/
        (train/test data)
        description.txt
        sample_submission.csv
        data_structure.txt (Optional)
   utils/
     _prepare.py
     _metric.py
   info (Optional)/
     description.txt
     raw_description.txt
     web_url.txt
     data_structure.txt
```

Figure 10: Standardized directory structure for each competition.

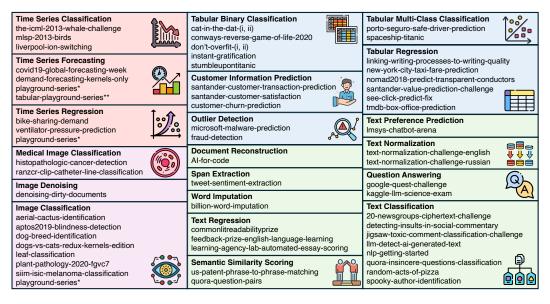


Figure 11: Overview of task diversity in MLE-Dojo, highlighting representative examples from four major domains; time series, computer vision, tabular data, and natural language processing.

# D Data Details

To establish a unified and flexible framework capable of accommodating diverse and emerging tasks, we develop customized data preparation scripts for each competition. These scripts systematically structure datasets, generate representative example submissions, and enable efficient local testing and evaluation. Given that most Kaggle competitions lack publicly available labels for test datasets, our scripts carefully split the original training data into new, clearly defined training and test subsets. We further implement tailored evaluation metrics specific to each competition, derived from a common base metric class. We begin with an initial set of around 600 Kaggle competitions. We've already manually reviewed them to confirm their suitability based on clarity of descriptions, data availability, and relevance of evaluation metrics, while excluding tasks with excessively large or unwieldy datasets. Our final collection balances challenge complexity with computational feasibility, yielding a practical,

scalable, and user-friendly benchmark of 200 competitions as our first release. We plan progressive releases of processed datasets to support ongoing benchmarking efforts.

Table 3 presents all competitions included in our first release, sorted in ascending order by data size. The associated tags provide information on each task's metric, category, modality, or domain. Users are encouraged to flexibly select tasks based on specific needs.

Table 3: Kaggle Competition Data Overview

<b>Competition Name</b>	Size	Tags
playground-series-s3e12 titanic	25.56 kB 93.08 kB	Beginner, Tabular, Binary Classification, Custom Metric Binary Classification, Tabular, Beginner, Categorization Accuracy
playground-series-s3e5	225.6 kB	Beginner, Tabular, Cohen Kappa Score
playground-series-s3e13	284.43 kB	Beginner, Tabular, Multiclass Classification, Health Conditions, MAP{K}
mercedes-benz-greener- manufacturing	351.45 kB	Automobiles and Vehicles, Regression, Tabular, $R^2$ score (coefficient of determination)
icr-identify-age-related-conditions	356.88 kB	Tabular, Binary Classification, Health, Weighted Multi- class Loss
kaggle-llm-science-exam	364.21 kB	Physics, NLP, MAP{K}
playground-series-s3e22	386.6 kB	Beginner, Tabular, Multiclass Classification, Animals, Health, F1 Score
playground-series-s3e3	455.63 kB	Tabular, Beginner, Binary Classification, Business, Area Under Receiver Operating Characteristic Curve
playground-series-s3e9	485.23 kB	Beginner, Tabular, Regression, Root Mean Squared Error
kobe-bryant-shot-selection	708.3 kB	Basketball, Binary Classification, Tabular, Log Loss
tabular-playground-series- jul-2021	826.96 kB	Tabular, Pollution, Time Series Analysis, Mean Columnwise Root Mean Squared Logarithmic Error
bike-sharing-demand	1.12 MB	Cycling, Tabular, Time Series Analysis, Root Mean Squared Logarithmic Error
home-data-for-ml-course	1.15 MB	Mean Absolute Error
spaceship-titanic	1.24 MB	Beginner, Tabular, Binary Classification, Categorization Accuracy
covid19-global-forecasting-week-2	1.26 MB	Coronavirus, Tabular, Mean Columnwise Root Mean Squared Logarithmic Error
covid19-global-forecasting-week-3	1.41 MB	Tabular, Coronavirus, Mean Columnwise Root Mean Squared Logarithmic Error
nlp-getting-started	1.43 MB	Text, Binary Classification, NLP, Custom Metric
covid19-global-forecasting- week-1	1.63 MB	Coronavirus, Tabular, Mean Columnwise Root Mean Squared Logarithmic Error
tabular-playground-series- jan-2022	1.73 MB	Tabular, Time Series Analysis, SMAPE
playground-series-s3e2	1.88 MB	Binary Classification, Tabular, Health Conditions, Beginner, Area Under Receiver Operating Characteristic Curve
spooky-author- identification	1.9 MB	Multiclass Classification, Literature, Linguistics, Multiclass Loss
covid19-global-forecasting- week-4	1.95 MB	Tabular, Coronavirus, Mean Columnwise Root Mean Squared Logarithmic Error
liberty-mutual-group- property-inspection-	2 MB	Housing, Normalized Gini Index
prediction		
playground-series-s3e25	2.11 MB	Beginner, Tabular, Regression, Earth Science, Median Absolute Error
us-patent-phrase-to-phrase- matching	2.14 MB	NLP, Text, PearsonCorrelationCoefficient
movie-review-sentiment- analysis-kernels-only	2.44 MB	Text, Multiclass Classification, Categorization Accuracy
playground-series-s3e6	2.7 MB	Beginner, Tabular, Housing, Root Mean Squared Error
commonlitreadabilityprize	2.93 MB	Text, Regression, Root Mean Squared Error
playground-series-s3e14	2.99 MB	Beginner, Tabular, Regression, Mean Absolute Error
detecting-insults-in-social- commentary	3.02 MB	Area Under Receiver Operating Characteristic Curve

Table 3: Kaggle Competition Data Overview (continued)

<b>Competition Name</b>	Size	Tags
walmart-recruiting-store- sales-forecasting	3.22 MB	Time Series Analysis, Weighted Mean Absolute Error
prudential-life-insurance-	3.4 MB	Tabular, QuadraticWeightedKappa
assessment	2.96 MD	Tout Internat Custom Matric
tweet-sentiment-extraction playground-series-s3e7	3.86 MB 3.86 MB	Text, Internet, Custom Metric Beginner, Tabular, Binary Classification, Area Under
		Receiver Operating Characteristic Curve
llm-detect-ai-generated- text	4.43 MB	Education, Primary and Secondary Schools, Binary Classification, Text Generation, Roc Auc Score
unimelb	4.53 MB	Area Under Receiver Operating Characteristic Curve
playground-series-s4e2	4.59 MB	Beginner, Time Series Analysis, Tabular, Multiclass Classification, Accuracy Score
playground-series-s4e3	5.49 MB	Beginner, Tabular, Multiclass Classification, Binary Classification, Manufacturing, Mean Columnwise Area Under Receiver Operating Characteristic Curve
tabular-playground-series- sep-2022	5.73 MB	Tabular, SMAPE
nomad2018-predict- transparent-conductors	6.24 MB	Chemistry, Mean Columnwise Root Mean Squared Logarithmic Error
amazon-employee-access-challenge	6.39 MB	Area Under Receiver Operating Characteristic Curve
playground-series-s3e1	6.48 MB	Regression, Tabular, Housing, Beginner, Root Mean Squared Error
tabular-playground-series- aug-2022	7.21 MB	Tabular, Binary Classification, Area Under Receiver Operating Characteristic Curve
playground-series-s3e18	7.63 MB	Beginner, Tabular, Binary Classification, Multilabel Classification, Roc Auc Score
poker-rule-induction	8.3 MB	Multiclass Classification, Card Games, Tabular, Categorization Accuracy
playground-series-s4e4	8.4 MB	Beginner, Tabular, Regression, Mean Squared Log Error
allstate-purchase- prediction-challenge	8.97 MB	Categorization Accuracy
playground-series-s3e16	9.05 MB	Beginner, Tabular, Regression, Animals, Mean Absolute Error
feedback-prize-english- language-learning	9.3 MB	NLP, Education, Primary and Secondary Schools, Custom Metric
dont-call-me-turkey	11.04 MB	Binary Classification, Tabular, Animals, Area Under Receiver Operating Characteristic Curve
tabular-playground-series- apr-2021	12.66 MB	Beginner, Tabular, Binary Classification, Categorization Accuracy
playground-series-s3e19	12.79 MB	Beginner, Tabular, Time Series Analysis, SMAPE
playground-series-s3e17	12.86 MB	Beginner, Tabular, Roc Auc Score
GiveMeSomeCredit	14.47 MB	Area Under Receiver Operating Characteristic Curve
google-quest-challenge	14.85 MB	Text, NLP, Mean Columnwise Spearman's r (rank correlation coefficient)
forest-cover-type-kernels- only	15.38 MB	Tabular, Forestry, Categorization Accuracy
playground-series-s4e6	16.2 MB	Beginner, Tabular, Education, Accuracy Score
novozymes-enzyme- stability-prediction	16.39 MB	Chemistry, SpearmanR
integer-sequence-learning	17.91 MB	Tabular, Categorization Accuracy
random-acts-of-pizza	17.97 MB	Binary Classification, Text, Internet, Area Under Receiver Operating Characteristic Curve
playground-series-s3e23	18.63 MB	Beginner, Tabular, Binary Classification, Roc Auc Score
demand-forecasting- kernels-only	18.7 MB	Tabular, SMAPE
playground-series-s3e8	18.96 MB	Beginner, Tabular, Regression, Root Mean Squared Error
playground-series-s3e10	20.94 MB	Beginner, Tabular, Log Loss
playground-series-s4e1	21.65 MB	Beginner, Tabular, Binary Classification, Banking, Roc Auc Score

Table 3: Kaggle Competition Data Overview (continued)

Competition Name	Size	Tags
afsis-soil-properties	21.7 MB	Mean Columnwise Root Mean Squared Error
sberbank-russian-housing-	22.71 MB	Banking, Housing, Regression, Tabular, Root Mean
market		Squared Logarithmic Error
playground-series-s3e24	22.79 MB	Beginner, Tabular, Binary Classification, Health, Roc Aud Score
aerial-cactus-identification	25.4 MB	Earth and Nature, Image, Plants, Area Under Receiver Operating Characteristic Curve
crowdflower-weather- twitter	25.4 MB	Root Mean Squared Error
predicting-red-hat-	26.74 MB	Tabular, Business, Area Under Receiver Operating Charac
business-value	20.74 MD	teristic Curve
DontGetKicked	29.52 MB	Gini Index
battlefin-s-big-data-	30.23 MB	Mean Absolute Error
combine-forecasting- challenge	30.23 MD	Mean Absolute Error
bioresponse	31.06 MB	Log Loss
tabular-playground-series-	31.4 MB	Tabular, Time Series Analysis, Cities and Urban Areas,
mar-2022 chaii-hindi-and-tamil-	31.4 MB	Regression, Mean Absolute Error Text, Languages, Custom Metric
question-answering		
nbme-score-clinical- patient-notes	35.73 MB	Text, Medicine, Education, NLP, Custom Metric
leaf-classification	36.05 MB	Image, Multiclass Classification, Multiclass Loss
learning-agency-lab-	36.2 MB	Education, NLP, Primary and Secondary Schools, Cohen
automated-essay-scoring-2	30.2 MD	Kappa Score
20-newsgroups-ciphertext-	36.97 MB	Multiclass Classification, Text, F-Score (Macro)
challenge	30.97 WID	ividiticiass Classification, Text, 1-Score (iviacio)
dont-overfit-ii	38.6 MB	Tabular, Binary Classification, Area Under Receiver Operating Characteristic Compa
rossmann-store-sales	39.85 MB	ating Characteristic Curve Tabular, Time Series Analysis, Root Mean Square Percent age Error
sf-crime	45.24 MB	Multiclass Classification, Tabular, Crime, Multiclass Loss
imaterialist-challenge-	48.59 MB	MeanBestErrorAtK
furniture-2018	.0.0 1.12	Treating to the first treating to the first treating to the first treating
playground-series-s3e11	49.83 MB	Beginner, Tabular, Regression, Mean Squared Log Error
word2vec-nlp-tutorial	54.37 MB	Text, Binary Classification, Movies and TV Shows, Area Under Receiver Operating Characteristic Curve
jigsaw-toxic-comment-	55.18 MB	Text, Mean Columnwise Area Under Receiver Operating
classification-challenge	55 55 MD	Characteristic Curve
see-click-predict-fix	55.75 MB	Root Mean Squared Logarithmic Error
denoising-dirty-documents	58.81 MB	Image, Root Mean Squared Error
homesite-quote-conversion	65.13 MB	Binary Classification, Tabular, Area Under Receiver Operating Characteristic Curve
airbnb-recruiting-new-user- bookings	67.85 MB	Hotels and Accommodations, Recommender Systems, Tabular, NDCG{K}
cat-in-the-dat	67.96 MB	Tabular, Binary Classification, Area Under Receiver Opeating Characteristic Curve
tmdb-box-office-prediction	70.24 MB	Tabular, Movies and TV Shows, Root Mean Squared Logarithmic Error
home-depot-product- search-relevance	72.93 MB	Tabular, Root Mean Squared Error
facial-keypoints-detection	80.86 MB	Image, Root Mean Squared Error
nyc-taxi-trip-duration	89.91 MB	Tabular, Regression, Root Mean Squared Logarithmic Error
text-normalization- challenge-english-language	95.51 MB	Linguistics, Languages, Text, Categorization Accuracy
bnp-paribas-cardif-claims- management	103.75 MB	Banking, Tabular, Binary Classification, Log Loss
playground-series-s4e5	104.17 MB	Beginner, Tabular, Logistic Regression, R2 Score

Table 3: Kaggle Competition Data Overview (continued)

<b>Competition Name</b>	Size	Tags
wsdm-cup-multilingual- chatbot-arena	113.98 MB	Languages, Text Conversation, Accuracy Score
santander-customer-	119.04 MB	Tabular, Binary Classification, Banking, Area Under Re-
satisfaction playground-series-s3e20	119.66 MB	ceiver Operating Characteristic Curve Beginner, Tabular, Time Series Analysis, Root Mean Squared Error
text-normalization- challenge-russian-language	125.87 MB	Text, Linguistics, Languages, Categorization Accuracy
tabular-playground-series- mar-2021	131.65 MB	Tabular, Logistic Regression, Binary Classification, Area Under Receiver Operating Characteristic Curve
allstate-claims-severity	142.87 MB	Regression, Tabular, Mean Absolute Error
tabular-playground-series- jan-2021	144.43 MB	Tabular, Regression, Root Mean Squared Error
cat-in-the-dat-ii	145.84 MB	Binary Classification, Area Under Receiver Operating Characteristic Curve
liverpool-ion-switching	146.08 MB	Biology, F-Score (Macro)
tabular-playground-series-feb-2021	154.66 MB	Regression, Tabular, Root Mean Squared Error
yelp-recsys-2013	179.46 MB	Root Mean Squared Error
lmsys-chatbot-arena	184.19 MB	Languages, Text Conversation, Log Loss
llm-classification-	184.19 MB	Languages, Text Conversation, Log Loss
finetuning stumbleupon	196.18 MB	Text, Tabular, Internet, Area Under Receiver Operating Characteristic Curve
playground-series-s3e4	198.09 MB	Tabular, Binary Classification, Beginner, Area Under Receiver Operating Characteristic Curve
the-winton-stock-market- challenge	213.13 MB	Finance, Tabular, Weighted Mean Absolute Error
conways-reverse-game-of- life-2020	251.11 MB	Simulations, Board Games, Custom Metric
facebook-recruiting-iv- human-or-bot	260.68 MB	Binary Classification, Tabular, Internet, Area Under Receiver Operating Characteristic Curve
the-icml-2013-whale-challenge-right-whale-	293.14 MB	Area Under Receiver Operating Characteristic Curve
redux porto-seguro-safe-driver- prediction	300.58 MB	Tabular, Binary Classification, Normalized Gini Index
statoil-iceberg-classifier- challenge	302.1 MB	Binary Classification, Weather and Climate, Image, Log Loss
tabular-playground-series- aug-2021	337.38 MB	Tabular, Regression, Banking, Root Mean Squared Error
flavours-of-physics-kernels- only	436.14 MB	Custom Metric
tgs-salt-identification- challenge	483.07 MB	Geology, Image, Custom Metric
linking-writing-processes- to-writing-quality	485.71 MB	Education, NLP, Primary and Secondary Schools, Mean Squared Error
grupo-bimbo-inventory- demand	502.44 MB	Tabular, Food, Root Mean Squared Logarithmic Error
quora-question-pairs	523.24 MB	Text, Tabular, Linguistics, Internet, Log Loss
tabular-playground-series- may-2022	574.22 MB	Tabular, Binary Classification, Area Under Receiver Operating Characteristic Curve
bigquery-geotab- intersection-congestion	577.97 MB	Tabular, Regression, Cities and Urban Areas, Geospatial Analysis, Root Mean Squared Error
mlsp-2013-birds	585.1 MB	Area Under Receiver Operating Characteristic Curve
santander-customer-	606.35 MB	Banking, Tabular, Binary Classification, Area Under Re-
transaction-prediction	(11 (( 15)	ceiver Operating Characteristic Curve
loan-default-prediction global-wheat-detection	611.66 MB 643.57 MB	Mean Absolute Error Image, Plants, Custom Metric
giobai-wiicat-detection	UIVI   C.C+U	image, i failts, Custom Metric

Table 3: Kaggle Competition Data Overview (continued)

Competition Name	Size	Tags
tabular-playground-series-dec-2021	693.17 MB	Tabular, Multiclass Classification, Categorization Accuracy
ventilator-pressure- prediction	698.79 MB	Tabular, Medicine, Biology, Mean Absolute Error
whale-categorization- playground	726.74 MB	Image, Animals, MAP{K}
dog-breed-identification	750.43 MB	Multiclass Classification, Animals, Image, Multiclass Los
plant-pathology-2020-	823.79 MB	Image, Agriculture, Mean Columnwise Area Under Re-
fgvc7 dogs-vs-cats-redux-kernels- edition	854.51 MB	ceiver Operating Characteristic Curve Image, Animals, Binary Classification, Log Loss
benchmark-bond-trade- price-challenge	910.54 MB	Weighted Mean Absolute Error
pubg-finish-placement- prediction	965.66 MB	Video Games, Tabular, Mean Absolute Error
instant-gratification	972.6 MB	Tabular, Binary Classification, Area Under Receiver Ope ating Characteristic Curve
tabular-playground-series- nov-2021	1.04 GB	Tabular, Binary Classification, Area Under Receiver Ope ating Characteristic Curve
petfinder-pawpularity- score	1.04 GB	Image, Root Mean Squared Error
santander-value-prediction- challenge	1.08 GB	Banking, Finance, Root Mean Squared Logarithmic Error
champs-scalar-coupling	1.22 GB	Chemistry, Tabular, Regression, Custom Metric
avazu-ctr-prediction	1.28 GB	Log Loss
tabular-playground-series- sep-2021	1.37 GB	Tabular, Binary Classification, Area Under Receiver Operating Characteristic Curve
billion-word-imputation	1.7 GB	Text, Linguistics, Levenshtein Mean
plant-seedlings- classification	1.81 GB	Plants, Image, Multiclass Classification, Custom Metric
tabular-playground-series- feb-2022	1.87 GB	Tabular, Multiclass Classification, Categorization Accuracy
shopee-product-matching	1.92 GB	Image, Text, Retail and Shopping, Custom Metric
AI4Code	2.16 GB	NLP, Text, Computer Science, Custom Metric
the-nature-conservancy- fisheries-monitoring	2.27 GB	Image, Multiclass Classification, Multiclass Loss
jigsaw-unintended-bias-in- toxicity-classification	2.38 GB	Text, NLP, Custom Metric
uw-madison-gi-tract- image-segmentation	2.47 GB	Image, Medicine, Dice3DHausdorff
ashrae-energy-prediction	2.61 GB	Tabular, Energy, Root Mean Squared Logarithmic Error
stanford-covid-vaccine	2.68 GB	Biology, Biotechnology, Coronavirus, Public Health, Cu tom Metric
acquire-valued-shoppers- challenge	3.07 GB	Area Under Receiver Operating Characteristic Curve
facebook-recruiting-iii- keyword-extraction	3.19 GB	Custom Metric
invasive-species- monitoring	3.35 GB	Image, Plants, Area Under Receiver Operating Characteritic Curve
tabular-playground-series- oct-2021	3.49 GB	Tabular, Binary Classification, Area Under Receiver Operating Characteristic Curve
kuzushiji-recognition	4.51 GB	Image, Multiclass Classification, History, Japan, Custom Metric
nfl-player-contact- detection	5.01 GB	Health, Football, Video, Tabular, Matthews correlation coefficient
bengaliai-cv19	5.18 GB	Image, Multiclass Classification, Weighted Categorization Accuracy
new-york-city-taxi-fare- prediction	5.7 GB	Regression, Tabular, Root Mean Squared Error

Table 3: Kaggle Competition Data Overview (continued)

Competition Name	Size	Tags
cassava-leaf-disease-	6.19 GB	Image, Multiclass Classification, Plants, Categorization
classification		Accuracy
quora-insincere-questions-	6.56 GB	Text, Binary Classification, Custom Metric
classification		
histopathologic-cancer-	7.76 GB	Cancer, Medicine, Research, Area Under Receiver Operat-
detection		ing Characteristic Curve
microsoft-malware-	8.47 GB	Area Under Receiver Operating Characteristic Curve
prediction		
bms-molecular-translation	8.87 GB	Chemistry, Image, Levenshtein Mean
aptos2019-blindness-	10.22 GB	Image, Multiclass Classification, Medicine, Healthcare,
detection		QuadraticWeightedKappa
ranzcr-clip-catheter-line-	13.13 GB	Image, Multilabel Classification, Mean Columnwise Area
classification		Under Receiver Operating Characteristic Curve
plant-pathology-2021-	16.1 GB	Image, Plants, Custom Metric
fgvc8		
smartphone-decimeter-	22.9 GB	Geospatial Analysis, Signal Processing, Research, Tabular,
2022		Mobile and Wireless, Custom Metric
multi-modal-gesture-	23.23 GB	Custom Metric
recognition		
osic-pulmonary-fibrosis-	23.99 GB	Image, Healthcare, Laplace Log Likelihood
progression		
freesound-audio-tagging-	26.15 GB	Audio, Weighted Label Ranking Average Precision
2019		
hms-harmful-brain-activity-	26.4 GB	Tabular, Research, Signal Processing, Healthcare, Kull-
classification		back Leibler Divergence
hotel-id-2021-fgvc8	26.65 GB	Image, Public Safety, MAP{K}
imet-2020-fgvc7	29.46 GB	F-Score Beta (Micro)
predict-volcanic-eruptions-	31.25 GB	Signal Processing, Geology, Physics, Mean Absolute Error
ingv-oe		
airbus-ship-detection	31.41 GB	Image, IntersectionOverUnionObjectSegmentationBeta
alaska2-image-steganalysis	32.27 GB	Custom Metric
hubmap-kidney-	32.97 GB	Image, Health, Biology, Dice
segmentation		D 1 0 D 11 101 1 27 D/TT)
h-and-m-personalized-	34.56 GB	Recommender Systems, Retail and Shopping, MAP{K}
fashion-recommendations	24 05 GD	T 1440 B
draper-satellite-image-	36.07 GB	Image, MASpearmanR
chronology	25 02 GD	T. G. C. T.
vesuvius-challenge-ink-	37.02 GB	Image Segmentation, History, Image Text Recognition,
detection	46 60 GD	DiceFBeta
iwildcam-2019-fgvc6	46.68 GB	Image, Multiclass Classification, F-Score (Macro)
herbarium-2020-fgvc7	63.16 GB	Image, Plants, F-Score (Macro)
cdiscount-image-	78.12 GB	Multiclass Classification, Categorization Accuracy
classification-challenge	07.76 CD	M D C AZZ
inaturalist-2019-fgvc6	87.76 GB	MeanBestErrorAtK
icecube-neutrinos-in-deep-	117.18 GB	Tabular, Astronomy, Physics, MeanAngularError
ice	110 50 CD	Multiplace Classification Di-1 C-titi- A
iwildcam-2020-fgvc7	118.59 GB	Multiclass Classification, Biology, Categorization Accu-
aiim aavid10 d-tti	120 51 CD	racy Image Multilahal Classification Custom Matric
siim-covid19-detection	128.51 GB	Image, Multilabel Classification, Custom Metric
rsna-miccai-brain-tumor-	136.85 GB	Image, Binary Classification, Healthcare, Area Under
radiogenomic-classification	156 02 CP	Receiver Operating Characteristic Curve
seti-breakthrough-listen	156.02 GB	Astronomy, Signal Processing, Science and Technology,
harbarium 2021 f0	161 0 CP	Area Under Receiver Operating Characteristic Curve
herbarium-2021-fgvc8	161.9 GB	Image, Plants, F-Score (Macro)
herbarium-2022-fgvc9	163.17 GB	Plants, Image, F-Score (Macro)
vinbigdata-chest-xray-	205.96 GB	Image, Healthcare, Custom Metric
abnormalities-detection		

Table 4 presents 50 competitions used for evaluation in our benchmark, encompassing four categories of tasks: MLE-Lite, Tabular, NLP, and CV.

Table 4: Competitions for Evaluation

Competition Name	Size	Tags
	MI	E-Lite Competitions
spooky-author-	1.9 MB	Multiclass Classification, Literature, Linguistics, Multi-
identification		class Loss
detecting-insults-in-social- commentary	3.02 MB	Area Under Receiver Operating Characteristic Curve
nomad2018-predict-	6.24 MB	Chemistry, Mean Columnwise Root Mean Squared Loga-
transparent-conductors		rithmic Error
random-acts-of-pizza	17.97 MB	Binary Classification, Text, Internet, Area Under Receiver Operating Characteristic Curve
aerial-cactus-identification	25.4 MB	Earth and Nature, Image, Plants, Area Under Receiver Operating Characteristic Curve
leaf-classification	36.05 MB	Image, Multiclass Classification, Multiclass Loss
jigsaw-toxic-comment-	55.18 MB	Text, Mean Columnwise Area Under Receiver Operating
classification-challenge denoising-dirty-documents	58.81 MB	Characteristic Curve Image, Root Mean Squared Error
text-normalization-	95.51 MB	Linguistics, Languages, Text, Categorization Accuracy
challenge-english-language	, , , , , , , , , , , , , , , , , , , ,	
text-normalization-	125.87 MB	Text, Linguistics, Languages, Categorization Accuracy
challenge-russian-language the-icml-2013-whale-	293.14 MB	Area Under Receiver Operating Characteristic Curve
challenge-right-whale-	293.14 MID	Area Onder Receiver Operating Characteristic Curve
redux		
tabular-playground-series-	574.22 MB	Tabular, Binary Classification, Area Under Receiver Oper-
may-2022 mlsp-2013-birds	585.1 MB	ating Characteristic Curve Area Under Receiver Operating Characteristic Curve
tabular-playground-series-	693.17 MB	Tabular, Multiclass Classification, Categorization Accu-
dec-2021		racy
dog-breed-identification plant-pathology-2020-	750.43 MB 823.79 MB	Multiclass Classification, Animals, Image, Multiclass Loss Image, Agriculture, Mean Columnwise Area Under Re-
fgvc7	623.79 MID	ceiver Operating Characteristic Curve
dogs-vs-cats-redux-kernels-	854.51 MB	Image, Animals, Binary Classification, Log Loss
edition	5.7.CD	D 'TIID (M C IF
new-york-city-taxi-fare- prediction	5.7 GB	Regression, Tabular, Root Mean Squared Error
histopathologic-cancer-	7.76 GB	Cancer, Medicine, Research, Area Under Receiver Operat-
detection	40.00.00	ing Characteristic Curve
aptos2019-blindness- detection	10.22 GB	Image, Multiclass Classification, Medicine, Healthcare, QuadraticWeightedKappa
ranzcr-clip-catheter-line-	13.13 GB	Image, Multilabel Classification, Mean Columnwise Area
classification		Under Receiver Operating Characteristic Curve
	Ì	NLP Competitions
kaggle-llm-science-exam	364.21 kB	Physics, NLP, MAP@K
llm-detect-ai-generated-	4.43 MB	Education, Primary and Secondary Schools, Binary Classifaction, Taxt Conception, Page Avg Second
text 20-newsgroups-ciphertext-	36.97 MB	fication, Text Generation, Roc Auc Score Multiclass Classification, Text, F-Score (Macro)
challenge	30.77 IVID	Francisco Classification, Text, 1-Score (Macro)
lmsys-chatbot-arena	184.19 MB	Languages, Text Conversation, Log Loss
stumbleupon	196.18 MB	Text, Tabular, Internet, Area Under Receiver Operating
linking-writing-processes-	485.71 MB	Characteristic Curve Education, NLP, Primary and Secondary Schools, Mean
to-writing-quality	103.71 1111	Squared Error
quora-question-pairs	523.24 MB	Text, Tabular, Linguistics, Internet, Log Loss
AI4Code quora-insincere-questions-	2.16 GB 6.56 GB	NLP, Text, Computer Science, Custom Metric Text, Binary Classification, Custom Metric
classification	0.JU UD	Teat, Billary Classification, Custom Metric
		CV Competitions
facial-keypoints-detection	80.86 MB	Image, Root Mean Squared Error
negponie detection	50.00 1,115	Continued on next nac

<b>Competition Name</b>	Size	Tags
whale-categorization- playground	726.74 MB	Image, Animals, MAP@K
petfinder-pawpularity-	1.04 GB	Image, Root Mean Squared Error
score bengaliai-cv19	5.18 GB	Image, Multiclass Classification, Weighted Categorization Accuracy
cassava-leaf-disease- classification	6.19 GB	Image, Multiclass Classification, Plants, Categorization Accuracy
bms-molecular-translation	8.87 GB	Chemistry, Image, Levenshtein Mean
imet-2020-fgvc7	29.46 GB	F-Score Beta (Micro)
airbus-ship-detection alaska2-image-steganalysis	31.41 GB 32.27 GB	Image, IntersectionOverUnionObjectSegmentationBeta Custom Metric
draper-satellite-image- chronology	36.07 GB	Image, MASpearmanR
	Ta	abular Competitions
demand-forecasting- kernels-only	18.7 MB	Tabular, SMAPE
dont-overfit-ii	38.6 MB	Tabular, Binary Classification, Area Under Receiver Operating Characteristic Curve
santander-customer- satisfaction	119.04 MB	Tabular, Binary Classification, Banking, Area Under Receiver Operating Characteristic Curve
liverpool-ion-switching	146.08 MB	Biology, F-Score (Macro)
conways-reverse-game-of- life-2020	251.11 MB	Simulations, Board Games, Custom Metric
porto-seguro-safe-driver- prediction	300.58 MB	Tabular, Binary Classification, Normalized Gini Index
santander-customer- transaction-prediction	606.35 MB	Banking, Tabular, Binary Classification, Area Under Receiver Operating Characteristic Curve
ventilator-pressure- prediction	698.79 MB	Tabular, Medicine, Biology, Mean Absolute Error
instant-gratification	972.6 MB	Tabular, Binary Classification, Area Under Receiver Operating Characteristic Curve
santander-value-prediction- challenge	1.08 GB	Banking, Finance, Root Mean Squared Logarithmic Error

We ranked the difficulty of competitions based on the average HumanRank scores achieved by all models across each competition. Table D presents all competitions for evaluation ranked by difficulty from easiest to hardest. This difficulty ranking provides a partial indication of which competitions MLE Agents perform relatively better in comparison to human competitors. However, due to potential biases introduced by factors such as dataset splits, the ranking should be interpreted as a reference rather than a definitive measure.

Table 5: Competition Difficulty

Competition	Avg. HumanRank	Category
tabular-playground-series-dec-2021	1.000000	MLE-Lite
detecting-insults-in-social-commentary	0.995000	MLE-Lite
llm-detect-ai-generated-text	0.808599	NLP
santander-customer-satisfaction	0.777327	Tabular
20-newsgroups-ciphertext-challenge	0.741197	NLP
plant-pathology-2020-fgvc7	0.706198	MLE-Lite
dogs-vs-cats-redux-kernels-edition	0.704814	MLE-Lite
histopathologic-cancer-detection	0.635389	MLE-Lite
aerial-cactus-identification	0.625717	MLE-Lite
airbus-ship-detection	0.578445	CV
demand-forecasting-kernels-only	0.576934	Tabular
nomad2018-predict-transparent-conductors	0.560720	MLE-Lite
dont-overfit-ii	0.532964	Tabular
random-acts-of-pizza	0.506494	MLE-Lite

Competition	Avg. HumanRank	Category
aptos2019-blindness-detection	0.474492	MLE-Lite
draper-satellite-image-chronology	0.472431	CV
jigsaw-toxic-comment-classification-challenge	0.437596	MLE-Lite
the-icml-2013-whale-challenge-right-whale-	0.423450	MLE-Lite
redux		
spooky-author-identification	0.405318	MLE-Lite
facial-keypoints-detection	0.393214	CV
porto-seguro-safe-driver-prediction	0.359001	Tabular
santander-value-prediction-challenge	0.337203	Tabular
tabular-playground-series-may-2022	0.336664	MLE-Lite
leaf-classification	0.324687	MLE-Lite
text-normalization-challenge-english-language	0.315865	MLE-Lite
lmsys-chatbot-arena	0.311389	NLP
dog-breed-identification	0.300586	MLE-Lite
petfinder-pawpularity-score	0.283909	CV
quora-question-pairs	0.283346	NLP
alaska2-image-steganalysis	0.257705	CV
cassava-leaf-disease-classification	0.223718	CV
conways-reverse-game-of-life-2020	0.223404	Tabular
denoising-dirty-documents	0.204969	MLE-Lite
stumbleupon	0.182091	NLP
imet-2020-fgvc7	0.179036	CV
whale-categorization-playground	0.166983	CV
mlsp-2013-birds	0.154272	MLE-Lite
liverpool-ion-switching	0.152072	Tabular
ventilator-pressure-prediction	0.136828	Tabular
kaggle-llm-science-exam	0.135041	NLP
text-normalization-challenge-russian-language	0.098380	MLE-Lite
linking-writing-processes-to-writing-quality	0.089086	NLP
instant-gratification	0.071988	Tabular
quora-insincere-questions-classification	0.070380	NLP
bms-molecular-translation	0.066147	CV
ranzer-clip-catheter-line-classification	0.036199	MLE-Lite
AI4Code	0.025269	NLP
bengaliai-cv19	0.018395	CV
new-york-city-taxi-fare-prediction	0.000843	MLE-Lite
santander-customer-transaction-prediction	0.000000	Tabular

# **E** Action Space

**request\_info.** When an LLM agent requires necessary information to solve an MLE problem, it can call the request\_info function to access task descriptions, sample submissions, data directories, output directories, and data structures. The response includes all necessary competition information without further analysis or information pruning.

validate\_code. The validate\_code action performs basic syntax checking and runtime validation. It provides a compilation trial of the code and returns detailed syntax or execution error messages. It also serves as a critical interface for running data analysis code and extracting deeper insights from the resulting outputs. This action is more lightweight than execute\_code without submission requirements, which may serve as a debugging or information printing tool.

**execute\_code.** Unlike validate\_code, the execute\_code action performs complete code execution, submission verification and evaluation. Only through this action can the generated submission of the code be validated and evaluated using the corresponding metric. Each use of the execute\_code action corresponds to a full competition submission. Restriction on invocations of this action allow for different and flexible methodological and experimental design needs. Only through the invocation of this action can the generated submission be validated and evaluated using the corresponding metric.

**get\_history.** Since memory is a crucial component in agent design, MLE-Dojo provides the get\_history function to access past experiences and enable learning from memory. This flexible design accommodates different agent architectures and scaffolding approaches.

**reset.** MLE-Dojo offers a reset function to reset the entire environment and restart from scratch. This action serves as the default mechanism designed to accommodate a wide range of use cases.

# **F** Implementation Details

#### F.1 Evaluation Metrics Details

We include additional details of evaluation metrics design as follows:

**AUP Score.** We use the AUP score to systematically evaluate and compare multiple methods across diverse tasks. A performance profile captures the proportion of tasks where a given method performs within a certain factor of the best-performing method. Lower performance ratios indicate better results, though we invert this ratio for metrics like accuracy or  $R^2$ , where higher values represent better performance. If a method fails to produce a valid solution (infeasible), we assign it a penalty relative to the worst feasible performance. Integrating these performance profiles provides the AUP score, offering a single robust measure of each method's overall effectiveness across the benchmark.

We adopt the *performance profile* and *AUP score* from ML-Gym [26] and to compare the effectiveness of different methods across a set of benchmark tasks. We make a few minor modifications to better align with our framework. For each backbone model  $m \in \mathcal{M}$  and task  $t \in \mathcal{T}$ , the performance ratio is defined as

$$r_{t,m} = \frac{\ell_{t,m}}{\min\limits_{m' \in \mathcal{M}} \ell_{t,m'}},\tag{1}$$

where  $\ell_{t,m}$  denotes the performance metric of backbone model m on task t, and  $\mathcal{M}$  is the set of all evaluated methods. This formulation assumes that lower metric values correspond to better performance.

The performance profile curve of backbone model m is then defined as the cumulative distribution of the log-scaled performance ratio:

$$\rho_m(\tau) = \frac{1}{|\mathcal{T}|} \left| \left\{ t \in \mathcal{T} \mid \log_{10} r_{t,m} \le \tau \right\} \right|. \tag{2}$$

This function quantifies the proportion of tasks for which backbone model m performs within a  $\tau$ -log distance of the best backbone model on that task.

To aggregate performance into a single scalar score, we compute the Area Under the Profile (AUP) curve:

$$AUP_m = \int_1^{\tau_{\text{max}}} \rho_m(\tau) d\tau, \tag{3}$$

where  $\tau_{\text{max}}$  is the smallest value such that  $\rho_m(\tau_{\text{max}}) = 1$  for all  $m \in \mathcal{M}$ .

For metrics where higher values indicate better performance (e.g., Accuracy, R<sup>2</sup>), we invert the ratio computation:

$$r_{t,m} = \frac{\max_{m' \in \mathcal{M}} \ell_{t,m'}}{\ell_{t,m}}.$$
 (4)

If a backbone model fails to obtain a valid score for a task t, its ratio is defined relative to the lowest valid score across all backbone models  $m_{\rm bottom}$  as:

$$r_{t,m} = (1 + \varepsilon) \cdot r_{t,m_{\text{bottom}}}, \quad \text{with } \varepsilon = 1.0.$$
 (5)

This is a relatively reasonable design choice that appropriately penalizes the performance of backbone models that fail to produce valid scores. Additionally, we impose an upper bound of 100 on the value of  $r_{t,m}$  to mitigate the risk of biased results caused by extremely large ratios.

**HumanRank Score.** The HumanRank score (Eq. F.1) measures the relative ranking of a submission within the competition leaderboard. Submissions receive higher scores if they achieve a better

rank among all participants. Suppose that the submission ranks at position p among a total of N submissions on the leaderboard. Then, the position score is computed as:  $s=1-\frac{p}{N}$ . To prevent bias between public and private leaderboards, we compute the relative scores on each leaderboard independently and then use their average as the final score.

**Elo ranking.** We adopt the Elo ranking calculation algorithm from Chatbot Arena [5]. We utilize Elo ratings to perform systematic pairwise comparisons of methods across competitions. Two complementary Elo calculation methods ensure stability and reliability: an online linear update with a low K-factor, providing stable incremental ratings, and a Bradley-Terry model using logistic regression, directly fitting ratings to all comparisons for robust estimates. Additionally, we apply bootstrapping techniques to compute confidence intervals, offering clear measures of uncertainty and confidence in the resulting ratings.

Each competition can be viewed as a battle, in which two backbone models are pitted against each other. Their respective performance scores on the competition determine the outcome-win, loss, or tie-between the two. Based on these pairwise outcomes, we calculate their Elo rankings. In our setting, we assume that the large language models (LLMs) under evaluation are static, allowing us to estimate a stable skill rating using the Bradley-Terry model. Specifically, we reformulate the pairwise win-loss outcomes between models as a logistic regression problem. Let  $\mathcal{M}$  denote the set of all models with cardinality  $|\mathcal{M}|=p$ , and let  $r_i$  denote the (logit-transformed) latent skill rating of model i. For every observed battle between model i and model j, we construct a pair of training samples  $(\mathbf{x}_{ij},y_{ij})$  and  $(\mathbf{x}_{ii},y_{ji})$  such that:

$$\mathbf{x}_{ij} = \log B \cdot (\mathbf{e}_i - \mathbf{e}_j), \quad y_{ij} = 1, \quad \mathbf{x}_{ji} = \log B \cdot (\mathbf{e}_j - \mathbf{e}_i), \quad y_{ji} = 0,$$

where B is the base of the logistic function (default B=10), and  $\mathbf{e}_i$  is a p-dimensional one-hot vector indicating model i. The number of duplicate entries per pair is determined by the total number of battles (wins, losses, and ties), with each entry weighted accordingly. The logistic regression is then solved as

$$\hat{\mathbf{r}} = \arg\min_{\mathbf{r}} \sum_{k=1}^{n} w_k \cdot \log \left( 1 + \exp(-y_k \cdot \mathbf{x}_k^{\top} \mathbf{r}) \right),$$

where  $w_k$  is the sample weight for the k-th comparison. The final Elo score for each model i is computed as  $s_i = S \cdot \hat{r}_i + R_0$ , where S is a scaling factor (default S = 400), and  $R_0$  is the initial Elo offset (default  $R_0 = 1000$ ). To estimate confidence intervals, we apply a non-parametric bootstrap procedure. Given the original battle dataset  $\mathcal{D}$ , we resample with replacement to generate  $\mathcal{D}_1, \ldots, \mathcal{D}_R$ , and re-estimate Elo scores as  $\mathbf{s}^{(1)}, \ldots, \mathbf{s}^{(R)}$ . The final Elo score for each model is reported as the pointwise median across bootstrap samples. All results in our study are reported using R = 100 bootstrap rounds with fixed random seed.

#### **G** Experiment Configurations.

The total number of steps is set to 15, with agents having full access to their interaction histories. A unified, concise prompt with clear instructions is provided, without additional extraneous information. The maximum runtime per session is 12 hours, and GPU memory is limited to 32 GB. The maximum input token length is set at 50,000, while each output round is capped at 8,192 tokens. Experiment settings of prompt, history, time, and memory configurations are designed to rigorously evaluate LLMs' capabilities in long-context handling, instruction-following, reasoning, and coding under resource-constrained conditions, closely mirroring realistic Machine Learning Engineering scenarios. We do not restrict the number of submission attempts to enable continuous improvement. To generate valid submission files and scores, agents must explicitly use the "execute\_code" command; submissions are not automated. See Appendix for detailed configurations. We pre-install commonly used Python packages, though agents may install additional packages within their generated code as needed.

#### **H** Additional Experimental Analysis

# H.1 Additional Analysis of Evaluation Metrics

We identify several key observations supported by analytical reasoning regarding the evaluation metrics utilized: HumanRank Score, Elo Score, and Performance Profiles with the AUP score.

HumanRank Scores Reflect Comprehensive Absolute Performance. The HumanRank Score (%) converts raw, task-specific metrics into relative rankings against human performance, yielding scores ranging from 0 (lowest) to 1 (highest). A high HumanRank score thus indicates superior performance relative to human benchmarks across tasks of diverse difficulty. This uniform normalization allows for comprehensive comparison and mitigates inconsistencies from heterogeneous raw scores. Our main experimental results illustrate that stronger models (e.g., o3-mini, DeepSeek-r1, and Gemini-2.5-Pro) consistently exhibit profiles enclosing those of less capable models, confirming their overall superior performance. Nonetheless, by equally weighting tasks irrespective of complexity, the averaging process may mask task-specific strengths or weaknesses and obscure the direct competitive relationships between individual models.

Elo Scores Highlight Pairwise Relative Performance. The Elo Score provides rankings based on explicit win-loss relationships determined through pairwise comparisons on each task. Utilizing the standard Elo rating system, these scores intuitively represent the relative ordering of models based on their competitive outcomes. Elo scores generally reflect a pattern where superior models encompass weaker ones, yet variations may arise. Notably, models with comparatively lower absolute average performance scores may still achieve higher Elo scores in certain task categories due to advantageous win-loss records. Elo scores thus effectively reveal relative model strengths and weaknesses but do not quantify the absolute magnitude of performance gaps.

**Performance Profiles and AUP Scores Demonstrate Robustness.** Performance Profiles and AUP scores accumulate performance ratios, which are relative values of raw scores, over varying thresholds. These metrics illustrate model robustness and consistency across multiple evaluation conditions. While these integrated scores offer insightful perspectives into relative performance dynamics, disparities in raw score scales across tasks can lead to comparability challenges. Despite potential biases due to varying raw scales, these metrics successfully capture nuanced relative performances, reflecting the models' adaptability across different performance thresholds.

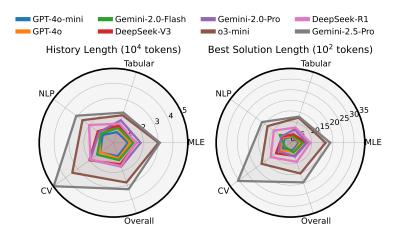


Figure 12: Total chat history length (left) and best solution length (right).

# H.2 History Length and Solution Length

Figure 12 shows both the total chat history length, including all interaction prompts, actions, and generated codes, and the length of the best solution generated by each model. The total chat history length closely aligns with the best solution length, where both metrics positively correlate with overall model performance. Reasoning models typically generate notably longer solutions compared to non-reasoning models. Moreover, stronger-performing models frequently produce more extended solutions, which often correspond to higher performance scores. Although increased solution length does not inherently ensure superior outcomes, it generally indicates a model's capability to explore more intricate and sophisticated solution strategies, a hallmark predominantly observed in more capable reasoning models.

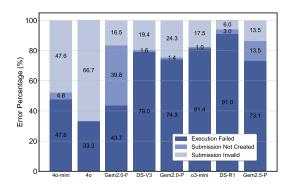


Figure 13: Execution error types.

# H.3 Error Type

**Execution Error Types.** To further analyze the nature of execution errors, we categorized failures into "Execution Failed", "Submission Not Created", and "Submission Invalid" (Figure 13). Stronger models, generating longer, more complex code, are more susceptible to "Execution Failed" errors. However, successful executions from these models frequently yield valid submissions and high performance. Models employing conservative validation strategies (e.g., gpt-4o-mini, gpt-4o) significantly reduce "Execution Failed" errors but continue to encounter issues related to submission creation and validity, indicating limitations in the validation approach.

# I Prompt Details

We design our prompts to be concise while containing all the information necessary for the LLM to perform the task. In doing so, we avoid providing extra assistance beyond what is required, while also ensuring that no essential information is omitted. This careful balance allows our benchmark to serve as a fair and comprehensive evaluation of LLMs as MLE Agents.

Below, we present the prompts used for the MLE Agents, which consist of four components. The **System Instruction** is an overall directive provided at the beginning of the task, detailing the task setting, the available actions, the expected output format, and an overview of the environment. During task execution, **Error** and **Reflection** as dynamic prompts are supplied to guide the LLM toward the next step, based on the current state: a concise prompt is used in the event of an error, while an informative prompt is used upon success, both incorporating feedback from the environment. In cases where the LLM produces outputs in an incorrect format, **Parse Error** prompt is triggered to prompt the LLM to correct its output. The detailed prompts are as follows:

```
System Instruction
You are a top-ranked Kaggle grandmaster with extensive competition experience.
Your objective is to solve a Kaggle competition, with the goal of maximizing
the Position Score (Your rank in the leaderboard) in limited steps.
You must use Machine Learning/Deep Learning/Computer Vision/NLP/etc.
methods to solve the problem, the score of random guess or without any ML/DL/CV/NLP
methods will be cancelled finally.
You are likely to train models according to specific competition requirements.
You have access to a GPU and several CPUs for training DL/ML models.
Use cuda and PyTorch for faster training whenever needed.
You have a total of {num_actions} actions available.
You have a total of {time_left} seconds, including code execution time.
You have access to exactly three actions with params, and receive corresponding
feedback after each action:
1. request_info: Retrieve specific competition information
    · params: info_type (str), must be one of: "overview", "sample_submission",
   "data_structure", "data_path", "output_path"
```

```
- feedback: information you requested
2. validate_code: Test (partial) code execution for debugging purposes or
   "print" information in the output
   - params: code (str)
   - feedback: execution result (success or failure), error message if failed,
   code output if success
3. execute_code: Run completed code, generate submission and get evaluation
   - params: code (str)
   - feedback: execution result, submission status, evaluation score
Code requirements:
- Request all information needed first
- Read all data files from data_dir
- Save all submissions to output_dir, should match test_data length
- Don't add, delete, or modify any files in data_dir
- Use "print" or other functions to output information in the feedback
- No plotting or visualization is allowed
- Refer to Sample Submission for the output format
- Code should be self-contained and not rely on any variables or state outside
- Code for submission should be completely runnable, otherwise it will be
 considered as failed
- Optimize your Model/Parameters/Data Processing/Algorithm for continuous improvement
Only if "execute_code" action taken, code successfully executed and valid submission
generated, you'll be able to get a Position Score (Your rank in the leaderboard) for
this competition.
Response format requirements, strictly one of the following:
    "action": "request_info",
    "params": {{
        "info_type": "<info_type>"
}}
or
{{
    "action": "validate_code",
    "params": {{
        "code": "<code>"
    }}
}}
or
{{
    "action": "execute_code",
    "params": {{
        "code": "<code>"
}}
- Must be valid JSON format
- No additional text or formatting allowed
```

```
Execution failed, details below:

#### Error Start ####
{observation}
#### Error End ####

You still have {num_actions} actions available.
You still have {time_left} seconds left.

Output your next action strictly following Response format requirements.
```

```
The results of your previous action:

#### Results Start ####
{observation}

#### Results End ####

You still have {num_actions} actions available.
You still have {time_left} seconds left.
Optimize your Model/Parameters/Data Processing/Algorithm for continuous improvement.

Output your next action strictly following Response format requirements.
```

```
Parse Error
The response can't be parsed as valid JSON.
Fix the error following the response format requirements.
First, only fix the JSON format error, don't change the contents of the response.
Then make sure the <code> is in completely runnable format for "python3 str(<code>)".
Only focus on the format, don't change the meaningful contents.
### Response Start ###
{response}
### Response End ###
### Error Start ###
{error}
### Error End ###
Response format requirements, strictly one of the following:
{{
    "action": "request_info",
    "params": {{
        "info_type": "<info_type>"
    }}
}}
{{
    "action": "validate_code",
    }}
or
{{
    "action": "execute_code",
    "params": {{
        "code": "<code>"
    }}
}}
- Must be valid JSON format
- No additional text or formatting allowed
```

# J Code Analysis

For the same competition, different models exhibit a rich diversity in the solutions or code they generate; for different competitions, even the same model is capable of producing specific and targeted solutions tailored to each task. We present the following analyses to illustrate how different models formulate their solutions when interacting with various competitions in MLE-Dojo.

#### J.1 Solution diversity of models

When faced with the same competition, different models often produce diverse solutions and code implementations, leading to great diversity of solutions. This diversity arises from variations in the choice of Python packages, the selection of computational devices (CPU or GPU), data preprocessing strategies, and model architectures.

Code 1 is a code snippet extracted from one of the step actions by **Gemini-2.5-Pro** on the 20-newsgroups-ciphertext-challenge competition. In this step, the model utilizes *sklearn* as main python package, CPU as device, employs classical feature engineering techniques, and applies a Logistic Regression model to address the task. Code 2 is a code snippet extracted from one of the step actions by **DeepSeek-r1** on the same competition. In contrast to previous approach, it prioritizes the use of GPU devices, adopts PyTorch as the core framework, and designs and trains the model in a more deep learning-oriented manner. This highlights the diverse actions and solution strategies adopted by different models when addressing the same problem.

```
import pandas as pd
import os
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import Pipeline
train_df = pd.read_csv(os.path.join(data_dir, 'train.csv'))
test_df = pd.read_csv(os.path.join(data_dir, 'test.csv'))
test_df = pd.read_csv(os.path.join(data_dir,
sample_submission = pd.read_csv(os.path.join(data_dir, 'sample_submission.csv'))
# Define features (X) and target (y)
X_train = train_df['ciphertext']
y_train = train_df['target']
X_test = test_df['ciphertext']
# Increased max_features and C again for the final attempt
pipeline = Pipeline([
  .
('<mark>tfidf</mark>', TfidfVectorizer(analyzer='<mark>char'</mark>, ngram_range=(2, 6), max_features=2000)),
  ('clf', OneVsRestClassifier(LogisticRegression(solver='liblinear', random_state=42,
       C=70.0))
# Train the model
pipeline.fit(X_train, y_train)
predictions = pipeline.predict(X_test)
submission_df = pd.DataFrame({'Id': test_df['Id'], 'Predicted': predictions})
```

Code Listing 1: Code snippet of Gemini-2.5-Pro on 20-newsgroups-ciphertext-challenge

```
import pandas as pd
import torch
import torch.nn as nn
from torch.utils.data import Dataset, DataLoader
from sklearn.preprocessing import LabelEncoder
from pathlib import Path
import math
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
class CharDataset(Dataset):
   def __init__(self, texts, labels=None):
       self.texts = texts
       self.labels = labels
   def __len__(self):
       return len(self.texts)
   def __getitem__(self, idx):
       text = self.texts[idx][:max_length].ljust(max_length, ' ')
       encoded = [char_to_idx.get(c, 0) for c in text]
       tensor = torch.tensor(encoded, dtype=torch.long)
       if self.labels is not None:
          label = torch.tensor(self.labels[idx], dtype=torch.long)
          return tensor, label
          return tensor
```

```
class CipherTransformer(nn.Module):
   def __init__(self, vocab_size=len(chars)+1, d_model=512, nhead=16, num_layers=8,
        num_classes=20):
       super().__init__()
       self.embed = nn.Embedding(vocab_size, d_model)
       self.pos_encoder = PositionalEncoding(d_model)
       encoder_layer = nn.TransformerEncoderLayer(
           d_model=d_model,
           nhead=nhead,
          dim_feedforward=2048,
           dropout=0.1,
           batch_first=True
       self.transformer = nn.TransformerEncoder(encoder_layer, num_layers)
       self.classifier = nn.Sequential(
          nn.Linear(d_model, 1024),
          nn.ReLU(),
          nn.Dropout(0.5),
          nn.Linear(1024, num_classes)
   def forward(self, x):
       x = self.embed(x)
       x = self.pos\_encoder(x)
       x = self.transformer(x)
       x = x.mean(dim=1)
       return self.classifier(x)
```

Code Listing 2: Code snippet of DeepSeek-R1 on 20-newsgroups-ciphertext-challenge

#### J.2 Solution diversity of competitions

A single model may adopt highly tailored solutions across competitions of varying nature. A comparison between Code 1 and Code 3 demonstrates that **Gemini-2.5-Pro** utilizes markedly distinct approaches in response to the specific demands of each task. This highlights the diverse actions and solution strategies adopted addressing different problems. MLE-Dojo offers a broad and flexible space that enables various models or agent scaffolds to fully leverage their respective capabilities without limitation.

```
import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import Dataset, DataLoader
from PIL import Image
import torchvision.transforms as transforms
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
class CactusDataset(Dataset):
   def __init__(self, df, img_dir, transform=None, is_test=False, tta_transform=None):
    self.df = df
       self.img_dir = img_dir
       self.transform = transform # Base transform (ToTensor, Normalize)
       self.tta_transform = tta_transform # Augmentation transform for TTA
       self.is_test = is_test
       self.image_files = df['id'].values
       if not self.is_test:
           self.labels = df['has_cactus'].values
   def __len__(self):
       return len(self.df)
   def __getitem__(self, idx):
       img_name = os.path.join(self.img_dir, self.image_files[idx])
       image = Image.open(img_name).convert('RGB')
       if self.is_test and self.tta_transform:
```

```
# Apply TTA transforms
           images = [self.transform(self.tta_transform(image, step)) for step in
                range(TTA_STEPS)]
           image_stack = torch.stack(images)
           return image_stack, self.image_files[idx]
       elif self.transform: # Standard transform for training or non-TTA test
           image = self.transform(image)
       if self.is test:
           return image, self.image_files[idx]
           label = torch.tensor(float(self.labels[idx]))
           return image, label
class SimpleCNN(nn.Module):
   def __init__(self):
       super(SimpleCNN, self).__init__()
self.conv1 = nn.Conv2d(3, 32, kernel_size=3, padding=1)
       self.relu1 = nn.ReLU()
       self.pool1 = nn.MaxPool2d(kernel_size=2, stride=2)
       self.conv2 = nn.Conv2d(32, 64, kernel_size=3, padding=1)
       self.relu2 = nn.ReLU()
       self.pool2 = nn.MaxPool2d(kernel_size=2, stride=2)
       self.conv3 = nn.Conv2d(64, 128, kernel_size=3, padding=1)
       self.relu3 = nn.ReLU()
       self.pool3 = nn.MaxPool2d(kernel_size=2, stride=2)
       self.fc1 = nn.Linear(128 * 4 * 4, 512)
       self.relu4 = nn.ReLU()
       self.dropout = nn.Dropout(0.5)
       self.fc2 = nn.Linear(512, 1)
   def forward(self, x):
       x = self.pool1(self.relu1(self.conv1(x)))
       x = self.pool2(self.relu2(self.conv2(x)))
x = self.pool3(self.relu3(self.conv3(x)))
       x = x.view(x.size(0), -1)
       x = self.relu4(self.fc1(x))
       x = self.dropout(x)
       x = self.fc2(x)
       return x
model = SimpleCNN().to(device)
```

Code Listing 3: Code snippet of Gemini-2.5-Pro on aerial-cactus-identification

#### J.3 Statistics of Method Diversity

We further investigate the distribution of model types implemented by different Large Language Models (LLMs) in Figure 14. To facilitate this analysis, the implemented models were systematically categorized based on their architectural complexity and type, utilizing the following mapping for brevity:

- Adv-NN: Advanced Neural Networks, encompassing contemporary and complex architectures such as Transformers (e.g., BERT, RoBERTa), EfficientNets, and LSTMs.
- Cls-NN: Classic Neural Networks, referring to well-established architectures like Convolutional Neural Networks (CNNs), ResNet variants, and DenseNets.
- **SD-NN**: Self-defined Neural Networks, indicating implementations involving customdesigned network structures or significant modifications to standard architectures, often specified within the solution code itself.
- Sim-LN: Simple Linear Models, comprising fundamental algorithms like Logistic Regression, Linear Regression, Support Vector Machines (with linear kernels), and Naive Bayes classifiers.
- Tree: Tree-based Ensemble Models, including algorithms such as Random Forest, Gradient Boosting Machines (GBM), XGBoost, LightGBM, and CatBoost.

• N/A: Not Applicable, denoting instances where the LLM did not generate a specific modeling solution or where the generated code did not contain an identifiable primary model.

Analysis of the model distribution across the evaluated LLMs and competition tasks reveals several pertinent observations. Firstly, a considerable degree of heterogeneity exists in the model choices preferred by different LLMs for identical tasks. Secondly, the selection of model type strongly correlates with the nature of the competition data; **Tree** and **Sim-LN** models are predominantly employed for tabular datasets, reflecting their established effectiveness in such domains, whereas **Cls-NN**, **Adv-NN**, and **SD-NN** are the favored choices for image and natural language processing tasks. Thirdly, certain LLMs exhibit discernible tendencies: for example, some models appear more frequently to propose **SD-NN** solutions, potentially indicating a higher propensity for generating novel or customized architectures, while others might more commonly default to simpler baselines (**Sim-LN**, **Tree**) or yield **N/A** results. The frequent occurrence of **N/A** across various LLM-task combinations also suggests variability in the LLMs' capabilities to successfully generate relevant and complete modeling code for diverse problem specifications. This distribution underscores the varying strategies and potential biases inherent in different LLMs when approaching machine learning problem-solving via code generation.

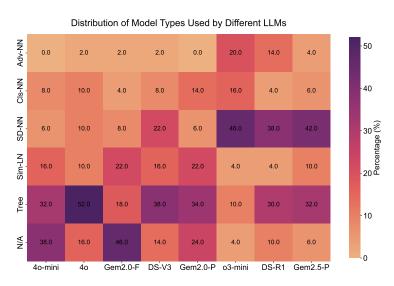


Figure 14: Statistics of method implementation across 50 evaluation tasks of eight frontiers LLM.

# K Agent Scaffolds

#### K.1 MLE Agent

In MLE-Dojo, the agent scaffold for benchmarking is referred to as the MLE Agent. It leverages the environment's natively supported action space and adopts a minimalistic yet effective design logic. LLMs are provided with a comprehensive and detailed initial instruction at the beginning of each task. This instruction includes the task description, output format, objectives, available actions, and resource constraints. Subsequently, the agent interacts with the environment through concise prompts that guide the LLM to take new actions based on feedback, in order to iteratively improve its performance. The MLE Agent maintains a fixed-length interaction history window and utilizes structured, formatted outputs as both actions and their corresponding contents. The MLE Agent serves as a lightweight yet representative scaffold that supplies essential guidance without offering additional hints or assistance. It enables systematic benchmarking and evaluation of an LLM's overall capabilities in the MLE-Dojo environment, including context understanding, analytical reasoning, instruction following, and code generation.

#### K.2 AIDE

AIDE [19] adopts a problem-solving paradigm inspired by how human data scientists approach challenges-through iterative refinement based on performance feedback. Central to its methodology is a technique termed *Solution Space Tree Search*, which enables systematic exploration and optimization over the space of candidate solutions. This framework comprises three core components: (1) a *Solution Generator*, which proposes new candidate solutions either by drafting from scratch or by modifying existing ones (e.g., fixing bugs or introducing improvements); (2) an *Evaluator*, which runs each candidate and quantitatively assesses its performance against the task objective; and (3) a *Solution Selector*, which identifies the most promising solution to seed the next iteration. Through repeated application of this feedback-driven cycle, AIDE efficiently navigates the solution space and converges towards optimal or near-optimal solutions. This iterative, adaptive process combines algorithmic rigor with human-like creativity, enabling AIDE to solve complex data science problems with remarkable effectiveness.

**Implementation.** We seamlessly integrate AIDE as the agent scaffold within MLE-Dojo, enabling its interaction with the environment through a single line of core code. The original score feedback mechanism in AIDE can be effortlessly replaced with either the HumanRank score or the actual raw competition score provided by MLE-Dojo. The interaction process can be fully executed by a single call to agent.step(exec\_callback=exec\_callback), ensuring compatibility without altering AIDE's core logic. Specifically, we make the following modifications: (1) We revise the prompting strategy to emphasize the requirement of generating a correctly formatted submission.csv file at the designated location. (2) The feedback mechanism and the criterion for selecting candidate nodes are modified to rely on the true raw score and HumanRank score provided by MLE-Dojo, instead of using the model's own validation performance.

# **NeurIPS Paper Checklist**

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Both our abstract and the end of the introduction clearly present our contributions, and our environment, code, and experimental results are well aligned with our claims.

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
  contributions made in the paper and important assumptions and limitations. A No or
  NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the limitation in Section A.1 in Appendix.

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

#### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include theoretical results.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We open-source all the code for data-preparation and running experiments; we also clearly present details of experiment settings in Section 4.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

# 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We open-souce all the code and present the link in the abstract.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
  to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

#### 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide all the experimental settings/details in Section 4.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

# 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We benchmark 8 models on 50 tasks and take the best performance of 3 trials on each task.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

#### 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We include the Experiments compute resources in the Experiment Configurations part in Section 4.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

#### 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We discuss both potential positive societal impacts and negative societal impacts of the work in Appendix.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: We discuss the data usage in Appendix, regarding the licenses of Kaggle data.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All the license and terms are mentioned for each kaggle task. And all the code we utilized or adapted is properly credited.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: All new assets introduced in the paper are well documented together in the Github repository mentioned in the abstract.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.