

ML-TOOL-BENCH: TOOL-AUGMENTED PLANNING FOR ML TASKS

Anonymous authors

Paper under double-blind review

ABSTRACT

The development of autonomous machine learning (ML) agents capable of end-to-end data science workflows represents a significant frontier in artificial intelligence. These agents must orchestrate complex sequences of data analysis, feature engineering, model selection, and hyperparameter optimization, tasks that require sophisticated planning and iteration. While recent work on building ML agents has explored using large language models (LLMs) for direct code generation, tool-augmented approaches offer greater modularity and reliability. However, existing tool-use benchmarks focus primarily on task-specific tool selection or argument extraction for tool invocation, failing to evaluate the sophisticated planning capabilities required for ML Agents. In this work, we introduce a comprehensive benchmark for evaluating tool-augmented ML agents using a curated set of 61 specialized tools and 15 tabular ML challenges from Kaggle. Our benchmark goes beyond traditional tool-use evaluation by incorporating an in-memory named object management, allowing agents to flexibly name, save, and retrieve intermediate results throughout the workflows. We demonstrate that standard ReAct-style approaches struggle to generate valid tool sequences for complex ML pipelines, and that tree search methods with LLM-based evaluation underperform due to inconsistent state scoring. To address these limitations, we propose two simple approaches: 1) using shaped deterministic rewards with structured textual feedback, and 2) decomposing the original problem into a sequence of sub-tasks, which significantly improves trajectory validity and task performance. Using GPT-4o, our approach improves over ReAct by 16.52 percentile positions, taking the median across all Kaggle challenges. We believe our work provides a foundation for developing more capable tool-augmented planning ML agents.

1 INTRODUCTION

Autonomous agents capable of solving end-to-end machine learning (ML) tasks represent a critical frontier in artificial intelligence (Grosnit et al., 2024; Toledo et al., 2025; Yang et al., 2025; Chan et al., 2025). Such agents must be capable of doing: data preprocessing, feature engineering, model training, and hyperparameter tuning, while managing intermediate results and adapting their strategies based on the evolving context. Achieving this level of autonomy requires not only sophisticated planning, but also memory management and the capacity to coordinate multiple operations coherently. Large language models (LLMs) have recently been explored as the foundation for such agents (Grosnit et al., 2024; Chan et al., 2025; Huang et al., 2024). Early work has primarily focused on direct code generation, where the agent generates python code for completing a given ML task (Grosnit et al., 2024; Chan et al., 2025; Huang et al., 2024; Toledo et al., 2025). This paradigm has shown promise on competitive benchmarks inspired by Kaggle challenges, with some approaches achieving performance comparable to a Kaggle Master (Grosnit et al., 2024; Chan et al., 2025). Several benchmarks have also been proposed to evaluate the performance of LLMs on such tasks (Chan et al., 2025; Huang et al., 2024; Qiang et al., 2025; Jing et al., 2025; Zhang et al., 2025). However, any approach that relies on direct code generation is prone to key weaknesses: generated code is brittle (Abbassi et al., 2025; Liu et al., 2025), debugging typically requires multiple iterations, and reasoning is tightly coupled with execution (Liu et al., 2025; Chen et al., 2025).

An alternative paradigm equips LLMs with external tools, yielding tool-augmented agents that need to decide which tools to invoke and in what sequence, to solve the task. Tools offer modular, reusable

building blocks for data-science workflows: from preprocessing, to training, and evaluation. This design has proven effective in broader domains, including web navigation (Zhou et al., 2024b), operating systems (Bonatti et al., 2024), and code interpretation (Huang et al., 2024), yet its potential for ML workflows remains underexplored. Crucially, tool augmentation reformulates the problem as planning in a large action space: the agent must coordinate multi-step trajectories and retrieve and reuse intermediate artifacts (or results). Because the agent is restricted to a curated toolset, tool-augmented approaches decouple high-level reasoning from low-level code execution, improving modularity, reliability, and safety.

Existing benchmarks for tool use fall short on long-horizon planning. Most benchmarks and approaches evaluate whether agents can select the right tools and valid arguments. The Berkeley Function-Calling Leaderboard (BFCL) (Patil et al., 2025) measures single, parallel, and multiple function calling, and BFCL-v3 (Patil et al., 2025) extends this to multi-turn, multi-step settings. However, even BFCL-v3 emphasizes relatively shallow plans compared to ML workflows, which might require long-term planning, iterative refinement, and reuse of intermediate artifacts. Similarly ToolBench (Xu et al., 2023) provides a suite of diverse software tools, that span both single-step and multi-step action generation, but focuses on evaluating whether the LLM can correctly select tools and tool arguments.

In this work, we introduce ML-Tool-Bench, motivated by the lack of good benchmarks to assess planning approaches with tools in ML workflows. In particular, ML-Tool-Bench provides a benchmark to evaluate the planning capabilities of LLM agents on *tabular* Kaggle ML challenges. We introduce a curated suite of 61 tools sufficient to solve such tasks and assess performance across 15 Kaggle challenges spanning regression and classification.

We evaluate multiple agents using several different planning algorithms, on our benchmark. To enable agents to create, persist, and reuse intermediate artifacts, we adopt an in-memory, named-object management scheme: tools accept references to named objects, and agents can assign names to tool outputs. We refer to this as *scratchpad-augmented planning*: agents store and retrieve objects by name over multi-step trajectories, enabling tools to handle arbitrarily large or structured inputs, unlike prior benchmarks that restrict arguments to simple types (e.g., strings, integers, floats). We observe that simple methods like ReAct (Yao et al., 2023b) struggle to produce performant trajectories across our Kaggle benchmark. Monte Carlo Tree Search-based methods (Kocsis & Szepesvari, 2006; Silver et al., 2016) such as LATS (Zhou et al., 2024a), which rely on LLMs as value estimators, also underperform due to inconsistent trajectory scoring. In contrast, we propose two simple approaches: 1) combining shaped, deterministic rewards with textual feedback and 2) decomposing the original problem into a sequence of sub-tasks. These approaches outperform the baselines, yielding more performant tool trajectories. These results highlight the difficulty of autonomous ML planning and point toward tool-augmented systems that rely less on subjective LLM scoring as tool sets grow in size and complexity.

1. We introduce ML-Tool-Bench, a tool-augmented benchmark for end-to-end ML planning with 61 tools and 15 Kaggle challenges.
2. We formalize *scratchpad-augmented planning* via named-object management that supports arbitrarily large artifacts and reversible branching in search.
3. We propose *MCTS-Shaped*, an MCTS approach with shaped, deterministic rewards and targeted textual feedback, which improves trajectory validity and performance over ReAct and LATS.
4. We introduce *Hierarchical MCTS*, an approach that decomposes problems into sequenced sub-tasks, further improving validity and robustness. For GPT-4o, Hierarchical MCTS improves over LATS by 9.93 percentile positions on the leaderboard and over ReAct by 16.52 percentile positions (median across all competitions). For GPT-4.1-mini, it improves over MCTS-Shaped by 1.89 percentile positions, while both ReAct and LATS had a median percentile position of 0.

Together, these advances establish strong baselines for tool-augmented, end-to-end ML planning and reduce reliance on subjective LLM scoring.

2 RELATED WORK

Machine Learning Benchmarks for AI Agents: Most of the existing Data Science and ML benchmarks, provide the LLM agent access to write code that solves the task, and evaluate its performance. Chan et al. (2025) propose MLE-bench, a curated benchmark of 75 Kaggle challenges, that test real-

world ML engineering skills. They find that OpenAI’s o1-preview with the AI-Driven Exploration (AIDE) scaffolding (Jiang et al., 2025) achieves at least a level of Kaggle bronze medal in 16.9% of competitions in their benchmark. AIRA-dojo (Toledo et al., 2025) improves upon Chan et al. (2025), replacing AIDE (Jiang et al., 2025) with a different choice of operator set, to generate new candidate solutions, and using Monte Carlo Tree Search (MCTS) (Kocsis & Szepesvari, 2006) instead of greedy search, increasing the success rate of achieving a Kaggle medal from 39.6% to 47.7%. Huang et al. (2024) also propose a ML benchmark, called MLEAgentBench, containing a suite of 13 tasks, where the agent is allowed to perform actions like read/write files, execute code and inspect outputs. They construct a ReAct based agent (Yao et al., 2023b) (with Claude v3 Opus) and were able to build compelling ML models on MLEAgentBench with 37.5% average success rate. Qiang et al. (2025) propose MLE-Dojo, an interactive gym-style workflow for LLM agents in iterative ML engineering workflows, and build upon 200+ Kaggle challenges. To evaluate Data Science Agents, Jing et al. (2025) proposed a comprehensive benchmark that includes 466 data analysis tasks and 74 data modeling tasks, sourced from Eloquence and Kaggle competitions, and showed that state of the art LLMs and agents struggle on most tasks. Zhang et al. (2025) propose DataSciBench and demonstrate that closed source models (GPT, Claude etc.) outperform open source models on all metrics in their benchmark.

Learning in Tool augmented LLMs: Solving ML challenges solely through the invocation of a fixed set of tools, in the correct sequential order, remains relatively unexplored. Approaches such as ARTIST (Singh et al., 2025), ReTool (Feng et al., 2025), StepTool (Yu et al., 2024), ToRL (Li et al., 2025), and ToolPlanner (Wu et al., 2024) couple reasoning and tool use for LLMs, using Reinforcement Learning to learn robust strategies for tool use. Recently, methods to fine-tune LLMs on responses containing tool usage have also been proposed (Schick et al., 2023; Qin et al., 2023; Gou et al., 2023; Patil et al., 2023).

Alternately, tree search methods (Yao et al., 2023a; Hao et al., 2023; Zhou et al., 2024a; Zhuang et al., 2023) have also been used to generate valid tool use trajectories. Zhuang et al. (2023) employs A* search, Hao et al. (2023) adopts Monte Carlo Tree Search (MCTS) and uses LLM as the world model, Zhou et al. (2024a) uses MCTS with value functions obtained from an LLM and self-reflection, and Yao et al. (2023a) explores Breadth-First Search (BFS) and Depth-First Search (DFS). However, these methods either depend on heuristic cost functions or leverage LLM feedback as a value function, and they are primarily applied to problems with relatively shallow depth. LATS (Zhou et al., 2024a) and Toolchain* (Zhuang et al., 2023) are the only approaches that explore planning with tools while the others restrict themselves to reasoning or toy domains. Feng et al. (2024) propose TS-LLM, an AlphaZero-inspired tree-search framework for LLMs that integrates a learned value function to guide decoding. The trajectories generated from tree search can further be used to fine-tune and improve the LLM, and TS-LLM has been shown to scale to tree depths of up to 64. Another approach, ReST-MCTS (Zhang et al., 2024), adopts a similar strategy to TS-LLM; however, in this case the per-step rewards are inferred directly from MCTS, whereas TS-LLM infers them using TD- λ (Sutton, 1988).

Tool Benchmarks: Benchmarks for LLM tool use largely emphasize correct tool selection and argument specification rather than extended planning. ToolBench (Xu et al., 2023) covers diverse software tools for single- and multi-step tasks but underplays long-horizon coordination. The Berkeley Function Calling Leaderboard (BFCL) (Patil et al., 2025) evaluates single, parallel, and multi-step calls, though plans remain shallow. τ -Bench (Yao et al., 2024) focuses on human-agent interaction under domain rules, highlighting alignment and information gathering more than proactive planning.

3 ML-TOOL-BENCH

Each task in ML-Tool-Bench, can be formalized as a Markov Decision Process (MDP) $(\mathcal{S}, \mathcal{A}, \mathcal{T}, R)$ (Puterman, 2014). The *state space* \mathcal{S} consists of the entire interaction history: all AI, Human, and Tool messages together with artifacts such as dataframes and ML models. Whenever a tool is executed, its observations (e.g., outputs, errors, logs) are appended to the history and folded into the state, so that the state maintains an up-to-date record of both conversational and artifact changes. The initial state s_0 comprises the Kaggle challenge description along with the dataset.

The *action space* is defined as: $\mathcal{A} = (\mathcal{A}_{\text{tool}} \cup \{\emptyset\}) \times (\mathcal{A}_{\text{reason}} \cup \{\emptyset\})$, where $\mathcal{A}_{\text{tool}}$ denotes the set of all tool invocations together with their full parameterizations (not just tool identity, but also

argument values and hyperparameters). This makes the benchmark challenging, since the effective size of $\mathcal{A}_{\text{tool}}$ can be very large rather than a small, discrete set. The set $\mathcal{A}_{\text{reason}}$ is the space of free-form reasoning steps, which we model as natural-language strings. The null element \emptyset denotes “no action” in that component, allowing tool-only, reason-only, both, or neither at a step. Reasoning actions organize information, plan future steps, and inject prior knowledge; tool actions modify or analyze data and train/evaluate models, thereby updating the state’s artifacts.

The *transition function* $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ maps a state–action pair (s, a) to the next state by appending the messages generated by the agent’s action, appending tool messages (i.e., the observations produced), and updating artifacts accordingly.

The *reward function* R evaluates progress and can be instantiated in several ways: (i) an outcome reward granted upon successful challenge completion; (ii) a shaped reward providing intermediate credit for measurable progress; or (iii) an LLM-based evaluation of the current state, using the LLM as a judge (Zheng et al., 2023) and absolving us from providing the reward function.

3.1 SCRATCHPAD

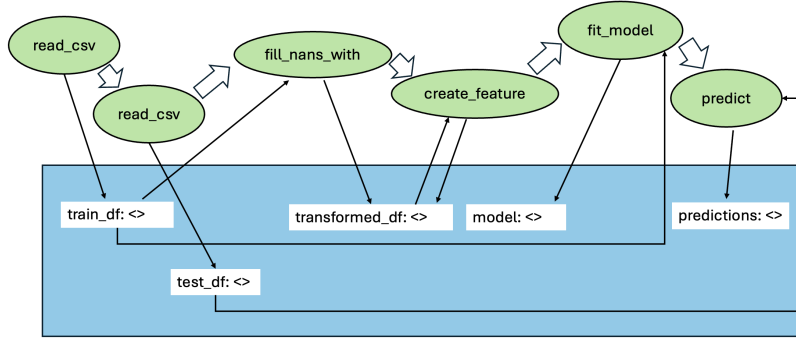


Figure 1: An illustration of our named-object management scheme. Green circles denote tool calls; the blue rectangle denotes the scratchpad (a key–value store). Each tool can read any named object from the scratchpad and write outputs back to it, depending on their read-write behavior. Arrows into a tool indicate inputs; arrows from a tool to the scratchpad indicate outputs. `read_csv` is a set tool; `fill_nans_with`, `fit_model`, and `predict` are get-set tools; `create_feature` is an override tool. There are two `read_csv` tool calls in the figure, one for train data and one for test.

Solving an ML challenge often involves storing large dataframes, models, and other complex artifacts as they cannot be directly passed as tool inputs by an LLM. A naive workaround is to maintain a single dataframe and model object that the agent incrementally modifies via tool calls. However, a single erroneous call can corrupt these objects, forcing a restart of the trajectory, and the agent becomes inflexible to create and reuse intermediate variables.

To address this, we adopt an in-memory, named-object management scheme: an agent assigns names to tool outputs, and tools accept references to named objects as inputs. Thus, agents can pass complex objects to tools by specifying the name under which the object is stored in the scratchpad. An illustration of this approach is presented in Figure 1. Implementing this requires modifying tools to operate on named references rather than raw objects; we describe these changes next.

3.2 TOOLS

We grant the agent access to a curated suite of 61 tools spanning data loading, data cleaning, feature engineering, and modeling. These tools are designed to be reasonably sufficient for solving tabular regression and classification tasks. Agent performance depends on the available toolset: in principle, a very large collection would maximize flexibility, but it results in an increased action space and complicates planning. We therefore adopt a fixed, compact tool set that trades some flexibility for a more tractable planning, while remaining adequate to solve the Kaggle challenges considered. For

modeling, we restrict to tree-based learners: Random Forest, XGBoost, LightGBM, and CatBoost, and linear/logistic regression, in light of the strong performance of tree-based methods on tabular Kaggle challenges (Grinsztajn et al., 2022). For more information on tools and how arbitrary user defined tools are modified to operate on named references rather than objects, refer to Appendix E

3.3 KAGGLE CHALLENGES

We select 15 tabular Kaggle ML challenges for ML-Tool-Bench: eight classification (binary and multiclass) and seven regression. These tasks are chosen so that they are solvable with our tool set. Several datasets are large (e.g., New York City Taxi Fare Prediction is ~ 2.5 GB), so we randomly sample 10,000 data points from each competition’s training set to keep planning computationally tractable. Because Kaggle test labels are hidden, we create an internal evaluation split by reserving 20% of the sampled training data as a test set with ground-truth labels. We evaluate using each competition’s official metric and report agent performance as the corresponding public-leaderboard percentile. Our evaluation metric is chosen to accommodate a collection of regression and multi-class classification tasks. Note that Kaggle leaderboards are computed on a test set, the labels to which we do not have access to; our reported results are computed on our held-out test split. For more information on the Kaggle challenges, refer to Appendix C

4 APPROACHES

4.1 REACT

ReAct (Yao et al., 2023b) is a prompting framework that interleaves natural-language reasoning (*Thought*) with tool interaction (*Action*) and the subsequent *Observation* from the environment due to tool calling. ReAct augments the agent’s action space to include the space of language, to account for thoughts or reasoning traces that do not affect the environment. Thoughts compose useful information from the current context and update the context to support future reasoning or actions. By explicitly exposing intermediate chain-of-thought alongside tool calls, ReAct enables agents to plan, invoke tools, and revise plans based on feedback. However, ReAct is unidirectional and can neglect potential alternative continuations from certain states, leading to locally optimal solutions (Zhuang et al., 2023; Zhou et al., 2024a).

4.2 MONTE CARLO TREE SEARCH (MCTS)

MCTS (Kocsis & Szepesvari, 2006) is a search algorithm that has achieved remarkable success in challenging domains such as Go (Silver et al., 2016) and Atari (Ye et al., 2021). MCTS builds a search tree where nodes correspond to states and edges correspond to actions. It comprises four phases: *selection*, *expansion*, *simulation/rollout*, and *backpropagation*. A common *selection* policy uses UCT (Upper Confidence Bound for Trees) (Kocsis & Szepesvari, 2006), choosing a child s of parent p such that: $s \in \arg \max_{s \in \mathbb{C}(p)} V(s) + w \sqrt{\ln N(p)/N(s)}$,

where $V(s)$ is the empirical value function, denoting the expected cumulative reward from state s , $N(p)$ is the parent’s visit count, $N(s)$ is the child’s visit count, $w > 0$ controls exploration, and $\mathbb{C}(p)$ denotes the set of children of p . Upon reaching a leaf node, it is *expanded* by selecting an action and adding the resulting next state as a child. From the newly expanded node, a *simulation* is run until the end of the episode or a fixed depth to obtain a reward r , which is then *backpropagated* along the trajectory to update values of all states along that trajectory: $V(s) \leftarrow (V(s)(N(s) - 1) + r)/N(s)$. MCTS is well-suited to large, irregular action spaces and provides a principled trade-off between exploration and exploitation. A pictorial illustration of MCTS is provided in Appendix B.

4.3 LANGUAGE AGENT TREE SEARCH (LATS)

LATS (Zhou et al., 2024a) adapts MCTS to language agents by using LLMs both to propose actions (reasoning steps or tool calls) and to evaluate node values. At each expansion, the policy LLM suggests candidates, and an evaluator LLM scores partial trajectories based on estimated progress toward the task objective. The value of a state is taken to be a weighted average of the evaluator LLM’s score and a self-consistency score (Wang et al., 2022), which upweights frequent candidates

in the expansion stage. In our tool-planning setting, we do not incorporate the self-consistency score into the value of a state. We observed that during the expansion phase, the LLM tends to propose only a small but distinct set of tool calls or reasoning steps, making the additional score unnecessary. LATS has shown improvements over purely reactive methods, such as ReAct (Yao et al., 2023b) on complex tasks. However, its value estimates can be noisy, and the effective planning depth may be limited by inconsistencies in evaluator scoring.

4.4 MCTS-SHAPED

In MCTS with shaped rewards, the agent receives intermediate credit for completing stages of the Kaggle ML challenge. The shaped-reward stages and their triggers are detailed below. Figure 2 provides an example to illustrate how rewards are provided in MCTS-Shaped.

Shaped-reward stages

1. **Train data loading:** reward when the agent successfully loads the training data.
2. **Test data loading:** reward when the agent successfully loads the test data. Note that test data does not have the target variable, that needs to be predicted.
3. **Combine train and test:** reward when the agent correctly concatenates train and test to enable consistent cleaning and feature engineering.
4. **Data cleaning:** reward when no missing values (NaNs) remain in the combined data.
5. **Feature engineering:** reward when (a) all categorical variables are properly encoded (e.g., one-hot or label encoding), and (b) the resulting feature dimensionality remains within a reasonable bound (to avoid exploding features from, e.g., high-cardinality text-like columns).
6. **Split back to train/test:** reward when the agent correctly splits the combined data back into train and test after transformations.
7. **Train features/target:** reward when the agent extracts $(X_{\text{train}}, y_{\text{train}})$ from the training dataframe using the correct target column.
8. **Test features:** reward when the agent extracts X_{test} from the test dataframe (which prior to this stage contains a dummy target), with correct arguments.
9. **Modeling:** reward when the agent successfully fits a model on the training data; the reward is proportional to cross-validation performance.
10. **Create submission:** reward when the agent generates predictions on the test data and writes a valid submission CSV to disk.

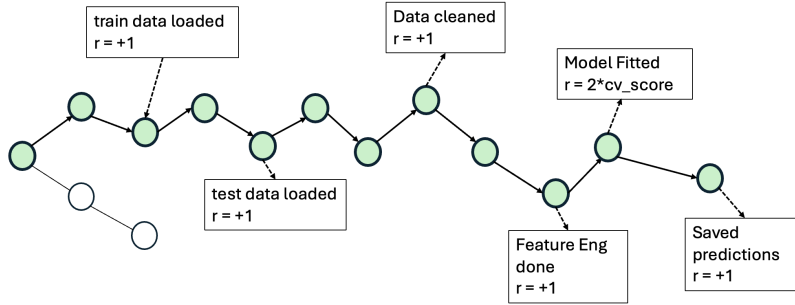


Figure 2: An example illustration of how rewards are provided in MCTS-Shaped. If a particular stage is judged to be successfully completed at a node, a reward is given, which is used to update the value of all the nodes in this trajectory. It needs to be noted that these stage-wise rewards are only provided once per trajectory and only if the earlier stages were successfully completed.

It needs to be noted that all of the stage rewards are provided to the agent only once per trajectory, and only if the earlier stages were successfully completed. The provided stage rewards are used to update the value of all the nodes in the trajectory. We verify stage completion using a reward function that inspects the node scratchpad and tool messages, confirming (i) that artifacts satisfy required properties (e.g., no NaNs for data cleaning; all columns encoded for feature engineering) and (ii) that the correct tools were invoked as evidenced by the tool logs.

4.5 HIERARCHICAL MCTS

We propose Hierarchical MCTS to improve over ReAct (Yao et al., 2023b), LATS (Zhou et al., 2024a), and classical MCTS (Kocsis & Szepesvari, 2006) in generating performant tool-use trajectories for solving Kaggle challenges within ML-Tool-Bench. Hierarchical MCTS decomposes a complex task into an ordered sequence of subtasks. We partition the available tools and assign them to relevant sub-tasks manually. For each subtask, MCTS searches its local state-action space to identify solution nodes. The solution nodes from one subtask are appended to the root of the next subtask, and the search continues. To avoid being trapped in locally optimal (but globally suboptimal) choices, we enumerate all solution nodes within each subtask up to a prescribed maximum subtask search depth. If there are no solution nodes identified after a subtask, the search terminates and we return ‘No Solution Found’. The solution node with the highest value, at the final subtask, is returned as the solution of the Hierarchical MCTS search. **Note that, when solving for each subtask in Hierarchical MCTS, we do not use any reward shaping and only check for if the subtask was solved successfully or not.** Importantly, the agent is given only the tools relevant to the current sub-task (tool masking), which reduces the branching factor and focuses the search. Figure 3 illustrates the overall procedure. Hierarchical MCTS is similar to the options framework (Sutton et al., 1999), that break down a complex problem into a hierarchy of sub-tasks, making the learning process more efficient and manageable for an agent.

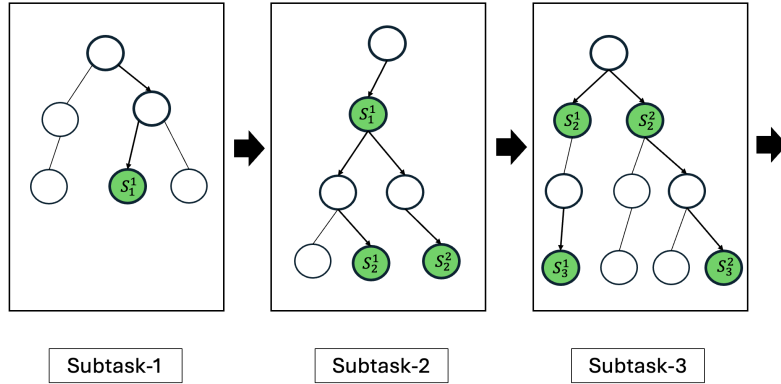


Figure 3: A schematic of Hierarchical MCTS. The task is decomposed into an ordered sequence of subtasks. For each subtask, MCTS searches for all solution nodes up to a prescribed maximum subtask depth to avoid locally optimal but globally suboptimal choices. The solution nodes from subtask t are appended to the root of subtask $t+1$, and the search resumes. In the example, the solution node from subtask 1, s_1^1 , initializes subtask 2; its solution nodes s_2^1 and s_2^2 initialize subtask 3, and so on. The highest-value solution at the final subtask is returned as the overall outcome of Hierarchical MCTS.

5 EXPERIMENTS

We evaluate the tool-planning performance of two language models—GPT-4o and GPT-4.1-mini, on ML-Tool-Bench. For each model, we compare five planning algorithms: (i) *ReAct* (Yao et al., 2023b); (ii) *LATS* (Zhou et al., 2024a); (iii) Monte Carlo Tree Search (MCTS) with outcome-based rewards, where the agent is rewarded upon successfully training a model or producing a valid submission file (denoted *MCTS-Outcome*); (iv) MCTS with shaped rewards, where the agent receives intermediate credit for completing stages of the Kaggle ML workflow (denoted *MCTS-Shaped*); and (v) *Hierarchical MCTS*: the Kaggle challenge is decomposed into subtasks. We use the reward stages defined for *MCTS-Shaped* as subtasks. A node is a solution node for a subtask, if it satisfies the reward condition for the stage corresponding to that subtask.

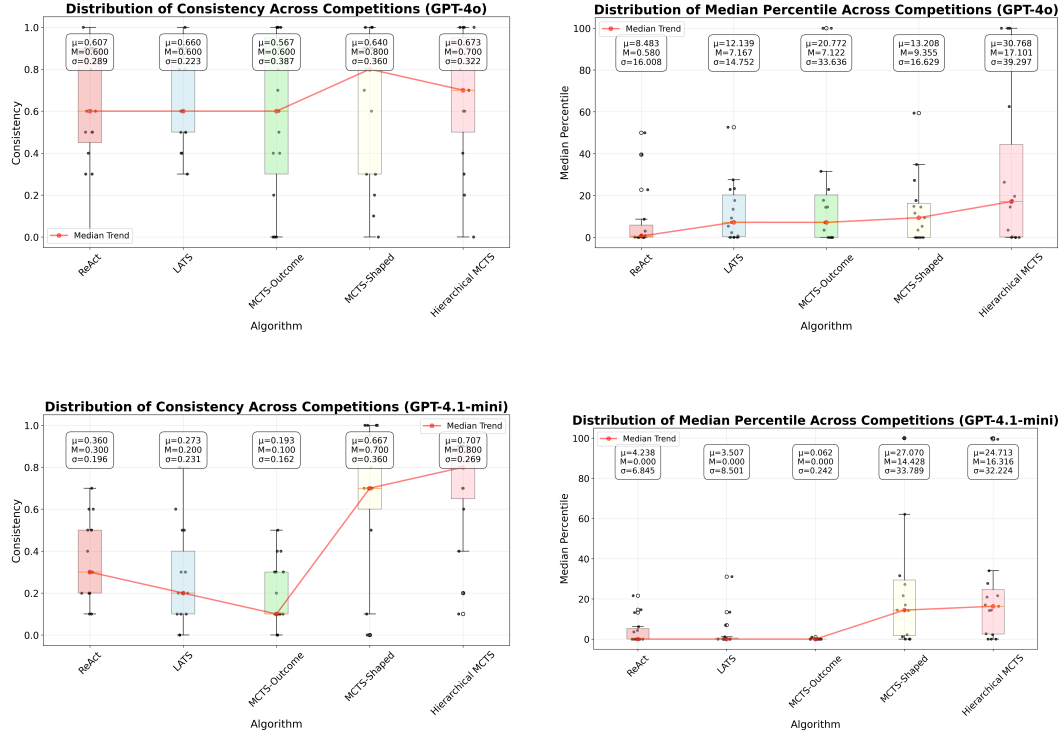


Figure 4: Plots of consistency and median leaderboard percentile across all competitions in ML-Tool-Bench, for different planning algorithms. The top row shows results for GPT-4o, with the left plot showing consistency and the right plot showing the median leaderboard percentile. The bottom row shows results for GPT-4.1-mini. Hierarchical MCTS outperforms LATS and ReAct, followed by MCTS-Shaped, in terms of leaderboard performance, for both LLMs. Also, both Hierarchical MCTS and MCTS-Shaped improve consistency over the other baselines. In the box plots, μ denotes the mean, σ denotes the standard deviation, and M denotes the Median

5.1 IMPLEMENTATION DETAILS

When using tree-search methods with our in-memory, named-object scheme, we adopt a *path-local scratchpad*, where each node v contains a scratchpad $\mathcal{S}(v)$, that stores only the objects produced by the tool call at that node. During expansion, the LLM proposes candidate actions. For a candidate that is a tool call, the accessible memory is the path union: $\mathcal{S}^*(v) = \bigcup_{u \in \text{path}(\text{root} \rightarrow v)} \mathcal{S}(u)$, and the LLM may reference any named object in $\mathcal{S}^*(v)$ as tool arguments. The tool’s outputs are written to the child’s scratchpad $\mathcal{S}(\text{child})$, preserving isolation per node while enabling reuse of intermediate artifacts along the trajectory.

LATS: To estimate the value of a state, we provide an evaluator LLM with all `AIMessage` and `ToolMessage` entries along the path from the root to the current node; it scores the trajectory by the progress made toward solving the Kaggle challenge. To propose candidate actions, we similarly pass the full trajectory history to the LLM, which returns new reasoning steps or tool calls. Unlike the original LATS formulation, we omit a self-consistency score from the value estimate, as at each expansion the agent typically proposes a small number of distinct candidates.

MCTS: We propose new candidate nodes during the expansion phase using the same approach listed in LATS. To evaluate the value of a node, we check if it produces a model or a valid submission file in the outcome rewards case. In the shaped rewards case, a node is provided a reward if it successfully completes a stage, as detailed earlier. In the case of Hierarchical MCTS, we designate a node as a solution node of the subtask, if it successfully completes the stage corresponding to that subtask. Additionally, across all MCTS variants, we apply a per-level depth penalty of 0.1 to discourage unnecessarily long trajectories that fail to make progress toward the goal.

In addition to rewards, we provide targeted textual feedback to help the agent refine its plan. When a stage fails, the agent receives an explanation of the failure. For example, in *feature engineering* we flag remaining categorical columns or an excessive increase in feature dimensionality; in *data cleaning* we report the presence of missing values. If a tool invocation fails, we return an explicit message along with the tool’s docstring to guide correct usage on the next attempt. We find that such feedback is crucial for consistently producing valid trajectories. This textual feedback is provided for all the MCTS variants (*MCTS-Outcome*, *MCTS-Shaped*, and *Hierarchical-MCTS*).

Ideally, we would run Monte Carlo rollouts to a fixed depth or until episode termination and use the return to update the value of all the nodes in the trajectory. Running to termination is impractical due to cost and compute constraints. Shallow rollouts (depth 3–5) are viable but GPT usage across many Kaggle challenges, planning algorithms, and trials, and roll outs at each state, resulted in extremely high costs and was infeasible. Learning value functions to approximate the value of states (Silver et al., 2016) is also not straightforward, on account of complex artifacts that are a part of the state. Consequently, we use the immediate reward at the current state (a depth-0 rollout), yielding a best-first search with a UCT-style exploration bonus. When budget permits, using small depth rollouts is preferred.

Hierarchical MCTS: In Hierarchical MCTS, we begin by decomposing the Kaggle challenge into a sequence of subtasks. This decomposition leverages the domain knowledge that solving a machine learning challenge typically involves data loading, data cleaning, feature engineering, modeling, evaluation, and prediction. We use a similar subtask decomposition to the reward-shaping structure used for MCTS-Shaped, described in Section 4.4.

Once we obtained textual descriptions for each subtask, a state-of-the-art coding agent was prompted with the subtask descriptions and the docstrings of the tools in our toolset, and asked to assign the relevant tools required to solve each subtask. One of the authors then manually reviewed the assignments to verify that the tool selections were sufficient and corrected minor errors made by the agent. This approach provides a general recipe for assigning tools to subtasks and can be applied in other domains, not only in the machine learning challenge-solving setting considered in this paper.

5.2 RESULTS

We evaluate GPT-4o and GPT-4.1-mini on our benchmark. For each algorithm–Kaggle challenge combination, we run 10 trials. We define *consistency* as the proportion of valid trajectories (e.g., 4 valid trajectories out of 10 trials yields a consistency of 0.4). For each trial, we evaluate predictions against the provided test labels using the competition’s official metric and compare against the leaderboard to obtain a leaderboard percentile. For each algorithm and competition, we report the median percentile across the 10 trials. Figure 4 presents boxplots for all algorithms, summarizing the distribution of leaderboard percentiles across all competitions in our benchmark. For further details on consistency and leaderboard percentiles for both models, refer to Appendix D. Additionally, for details on the prompts used, refer to Appendix H. For additional details regarding costs, tool-masking ablations, and example trajectories refer to Appendix G, F, and I respectively.

As shown in Figure 4, Hierarchical MCTS improves leaderboard performance compared to ReAct, LATS, and MCTS-Outcome, followed by MCTS-Shaped, for both GPT-4o and GPT-4.1-mini. Moreover, both Hierarchical MCTS and MCTS-Shaped achieve higher consistency than the other baselines. For GPT-4o, Hierarchical MCTS shows improvement over LATS by 9.93 percentile positions on the leaderboard and over ReAct by 16.52 percentile positions, taking the median across all competitions. For GPT-4.1-mini, Hierarchical MCTS improved over MCTS-Shaped by 1.89 percentile positions on the leaderboard, while both ReAct and LATS had a median leaderboard percentile position of 0 across all competitions. These results highlight that as toolsets become more complex and larger, it is important either to introduce hierarchy—decomposing the original task into subtasks with corresponding reward functions, or to employ shaped rewards that guide the search toward solutions. In contrast, unidirectional planning strategies like ReAct do not perform well. Similarly, tree-search methods such as LATS, that rely solely on LLM evaluation also fail, as LLMs provide inconsistent scores to nodes when trajectory lengths increase, due to the accumulation of messages and artifacts that must be considered during evaluation.

We also report the Consistency and Leaderboard percentiles (with respect to the Kaggle public leaderboard) for the five planning approaches evaluated in this paper. For this analysis, we used

the original Kaggle train and test splits rather than the smaller benchmark subsets. Since test labels were not available, we submitted our predictions to Kaggle to obtain the public leaderboard scores, which were then converted into percentile ranks. Due to cost constraints, we evaluated only a subset of six challenges from our benchmark and used GPT-4.1-mini as the underlying LLM. The results are presented in the Tables 1 and 2. We observe that these results exhibit the same trends as those seen in our benchmark evaluation: Hierarchical MCTS and MCTS-Shaped consistently outperform the other methods, while ReAct and LATS struggle, achieving a median leaderboard percentile of 0.0 across most of the six challenges evaluated.

Competition	ReAct	LATS	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS
Spaceship Titanic	0.3	0.0	<u>0.6</u>	<u>0.6</u>	0.4
BPM Prediction	0.2	0.3	0	<u>0.8</u>	0.7
Calorie Expenditure Prediction	0.4	0.3	0.2	0.3	<u>0.6</u>
california Housing Regression	0.6	0.4	0.1	0.5	<u>0.8</u>
Bank Deposit Classification	0.4	0.5	0.0	<u>1.0</u>	0.6
Bank Churn Classification	0.2	0.5	0.1	1.0	<u>0.8</u>
Overall (Median)	0.35	0.35	0.1	<u>0.7</u>	0.65

Table 1: Consistency scores for the five planning approaches across six Kaggle challenges, evaluated using the original train/test splits provided by Kaggle. [Table added to address reviewer’s questions.](#)

Competition	ReAct	LATS	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS
Spaceship Titanic	0.0	0.0	39.77	<u>41.69</u>	0.0
BPM Prediction	0	0	0	<u>5.03</u>	0.19
Calorie Expenditure Prediction	0.0	0.0	0.0	0.0	<u>16.81</u>
california Housing Regression	5.51	0.0	0.0	8.99	<u>24.78</u>
Bank Deposit Classification	0.0	13.47	0.0	<u>27.65</u>	<u>27.29</u>
Bank Churn Classification	0.0	14.91	0.0	29.87	<u>32.47</u>
Overall (Median)	0.0	0.0	0.0	18.32	<u>20.80</u>

Table 2: Median leaderboard percentiles for the five planning approaches across six Kaggle challenges. Percentiles are computed from Kaggle public leaderboard scores obtained via official submissions using the original train/test splits. [Table added to address reviewer’s questions.](#)

6 CONCLUSION

We introduced ML-Tool-Bench, a benchmark for evaluating the planning capabilities of tool-augmented LLMs on tabular Kaggle challenges. Existing tool-use benchmarks (Xu et al., 2023; Patil et al., 2025; Yao et al., 2024) primarily assess tool selection and argument grounding, rather than long-horizon planning. By contrast, many ML agents generate code directly; while flexible, this approach sacrifices modularity, reliability, and safety compared to operating within a curated toolset. Empirically, we found that ReAct and LATS struggle to consistently produce valid and performant trajectories. We proposed two improved approaches: (i) MCTS with shaped, deterministic rewards, and (ii) Hierarchical MCTS, which decomposes problems into sequenced subtasks. Across two models, Hierarchical MCTS achieved the best leaderboard performance compared to other baselines, while both Hierarchical MCTS and MCTS-Shaped improved consistency, measured as the fraction of valid trajectories. These results suggest that incorporating subtask decomposition with deterministic rewards, rather than relying on subjective LLM evaluation, yields performance gains as the set of available tools grows in size and complexity.

REFERENCES

- Altat Allah Abbassi, Leuson Da Silva, Amin Nikanjam, and Foutse Khomh. A taxonomy of inefficiencies in llm-generated python code, 2025. URL <https://arxiv.org/abs/2503.06327>.
- Rogério Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong Lu, Justin Wagle, Kazuhito Koishida, Arthur Buckner, Lawrence Jang, and Zack Hui. Windows

- agent arena: Evaluating multi-modal os agents at scale, 2024. URL <https://arxiv.org/abs/2409.08264>.
- Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Lilian Weng, and Aleksander Madry. Mle-bench: Evaluating machine learning agents on machine learning engineering, 2025. URL <https://arxiv.org/abs/2410.07095>.
- Yongchao Chen, Harsh Jhamtani, Srinagesh Sharma, Chuchu Fan, and Chi Wang. Steering large language models between code execution and textual reasoning, 2025. URL <https://arxiv.org/abs/2410.03524>.
- Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjun Zhong. Retool: Reinforcement learning for strategic tool use in llms. *ArXiv*, abs/2504.11536, 2025. URL <https://api.semanticscholar.org/CorpusID:277824366>.
- Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. Alphazero-like tree-search can guide large language model decoding and training, 2024. URL <https://arxiv.org/abs/2309.17179>.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. Tora: A tool-integrated reasoning agent for mathematical problem solving. *ArXiv*, abs/2309.17452, 2023. URL <https://api.semanticscholar.org/CorpusID:263310365>.
- Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on tabular data?, 2022. URL <https://arxiv.org/abs/2207.08815>.
- Antoine Grosnit, Alexandre Maraval, James Doran, Giuseppe Paolo, Albert Thomas, Refinath Shahul Hameed Nabeezath Beevi, Jonas Gonzalez, Khyati Khandelwal, Ignacio Iacobacci, Abdelhakim Benechehab, Hamza Cherkaoui, Youssef Attia El-Hili, Kun Shao, Jianye Hao, Jun Yao, Balazs Kegel, Haitham Bou-Ammar, and Jun Wang. Large language models orchestrating structured reasoning achieve kaggle grandmaster level, 2024. URL <https://arxiv.org/abs/2411.03562>.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model, 2023. URL <https://arxiv.org/abs/2305.14992>.
- Addison Howard, Ashley Chow, and Ryan Holbrook. Spaceship titanic. <https://kaggle.com/competitions/spaceship-titanic>, 2022. Kaggle.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Mlagentbench: Evaluating language agents on machine learning experimentation, 2024. URL <https://arxiv.org/abs/2310.03302>.
- Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and Yuxiang Wu. Aide: Ai-driven exploration in the space of code, 2025. URL <https://arxiv.org/abs/2502.13138>.
- Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming Zhang, Xinya Du, and Dong Yu. Dsbench: How far are data science agents from becoming data science experts?, 2025. URL <https://arxiv.org/abs/2409.07703>.
- Levente Kocsis and Csaba Szepesvari. Bandit based monte-carlo planning. In *European Conference on Machine Learning*, 2006. URL <https://api.semanticscholar.org/CorpusID:15184765>.
- Xuefeng Li, Haoyang Zou, and Pengfei Liu. Torl: Scaling tool-integrated rl. *ArXiv*, abs/2503.23383, 2025. URL <https://api.semanticscholar.org/CorpusID:277451754>.
- Changshu Liu, Yang Chen, and Reyhaneh Jabbarvand. Codemind: Evaluating large language models for code reasoning, 2025. URL <https://arxiv.org/abs/2402.09664>.

- Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive apis, 2023. URL <https://arxiv.org/abs/2305.15334>.
- Shishir G. Patil, Huanzhi Mao, Charlie Cheng-Jie Ji, Fanjia Yan, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.
- Martin L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Rushi Qiang, Yuchen Zhuang, Yinghao Li, Dingu Sagar V K, Rongzhi Zhang, Changhao Li, Ian Shu-Hei Wong, Sherry Yang, Percy Liang, Chao Zhang, and Bo Dai. Mle-dojo: Interactive environments for empowering llm agents in machine learning engineering, 2025. URL <https://arxiv.org/abs/2505.07782>.
- Yujia Qin, Shi Liang, Yining Ye, Kunlun Zhu, Lan Yan, Ya-Ting Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Marc H. Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+ real-world apis. *ArXiv*, abs/2307.16789, 2023. URL <https://api.semanticscholar.org/CorpusID:260334759>.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *ArXiv*, abs/2302.04761, 2023. URL <https://api.semanticscholar.org/CorpusID:256697342>.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Joykirat Singh, Raghav Magazine, Yash Pandya, and Akshay Nambi. Agentic reasoning and tool integration for llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2505.01441>.
- Richard S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3:9–44, 1988. doi: 10.1007/BF00115009.
- Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- Edan Toledo, Karen Hambardzumyan, Martin Josifoski, Rishi Hazra, Nicolas Baldwin, Alexis Audran-Reiss, Michael Kuchnik, Despoina Magka, Minqi Jiang, Alisia Maria Lupidi, Andrei Lupu, Roberta Raileanu, Kelvin Niu, Tatiana Shavrina, Jean-Christophe Gagnon-Audet, Michael Shvartsman, Shagun Sodhani, Alexander H. Miller, Abhishek Charnalia, Derek Dunfield, Carole-Jean Wu, Pontus Stenatorp, Nicola Cancedda, Jakob Nicolaus Foerster, and Yoram Bachrach. Ai research agents for machine learning: Search, exploration, and generalization in mle-bench, 2025. URL <https://arxiv.org/abs/2507.02554>.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Qinzhao Wu, Wei Liu, Jian Luan, and Bin Wang. Toolplanner: A tool augmented llm for multi granularity instructions with path planning and feedback. *ArXiv*, abs/2409.14826, 2024. URL <https://api.semanticscholar.org/CorpusID:272827086>.
- Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. On the tool manipulation capability of open-source large language models, 2023. URL <https://arxiv.org/abs/2305.16504>.

- Sherry Yang, Joy He-Yueya, and Percy Liang. Reinforcement learning for machine learning engineering agents. 2025. URL <https://api.semanticscholar.org/CorpusID:281080955>.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *ArXiv*, abs/2305.10601, 2023a. URL <https://api.semanticscholar.org/CorpusID:258762525>.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023b. URL <https://arxiv.org/abs/2210.03629>.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for tool-agent-user interaction in real-world domains, 2024. URL <https://arxiv.org/abs/2406.12045>.
- Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games with limited data. *Advances in neural information processing systems*, 34:25476–25488, 2021.
- Yuanqing Yu, Zhefan Wang, Weizhi Ma, Zhicheng Guo, Jingtao Zhan, Shuai Wang, Chuhan Wu, Zhiqiang Guo, and Min Zhang. Steptool: Enhancing multi-step tool usage in llms via step-grained reinforcement learning. 2024. URL <https://api.semanticscholar.org/CorpusID:273233670>.
- Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. Rest-mcts*: Llm self-training via process reward guided tree search, 2024. URL <https://arxiv.org/abs/2406.03816>.
- Dan Zhang, Sining Zhoubian, Min Cai, Fengzu Li, Lekang Yang, Wei Wang, Tianjiao Dong, Ziniu Hu, Jie Tang, and Yisong Yue. Datascibench: An llm agent benchmark for data science, 2025. URL <https://arxiv.org/abs/2502.13897>.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL <https://arxiv.org/abs/2306.05685>.
- Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language agent tree search unifies reasoning acting and planning in language models, 2024a. URL <https://arxiv.org/abs/2310.04406>.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents, 2024b. URL <https://arxiv.org/abs/2307.13854>.
- Yuchen Zhuang, Xiang Chen, Tong Yu, Saayan Mitra, Victor S. Bursztn, Ryan A. Rossi, Somdeb Sarkhel, and Chao Zhang. Toolchain*: Efficient action space navigation in large language models with a* search. *ArXiv*, abs/2310.13227, 2023. URL <https://api.semanticscholar.org/CorpusID:264405734>.

A LARGE LANGUAGE MODEL USAGE

Large Language Models (LLMs) were used for grammatical editing and improving writing flow. Additionally, LLMs assisted the authors in conducting literature surveys and identifying related work. LLMs were also used to aid in developing the ML-Tool-Bench toolset. LLMs were used to assist with code generation, debugging, and documentation for components of the ML-Tool-Bench toolset, based on tool descriptions provided by the authors. LLMs were also used to assign relevant tools to each subtask in the proposed Hierarchical MCTS approach. All implementations were reviewed and validated by the authors. All research methodology, experimental design, data analysis, and scientific conclusions are entirely the work of the human authors.

B APPROACHES

B.1 MONTE CARLO TREE SEARCH

Figure 5 provides a pictorial illustration of the MCTS algorithm.

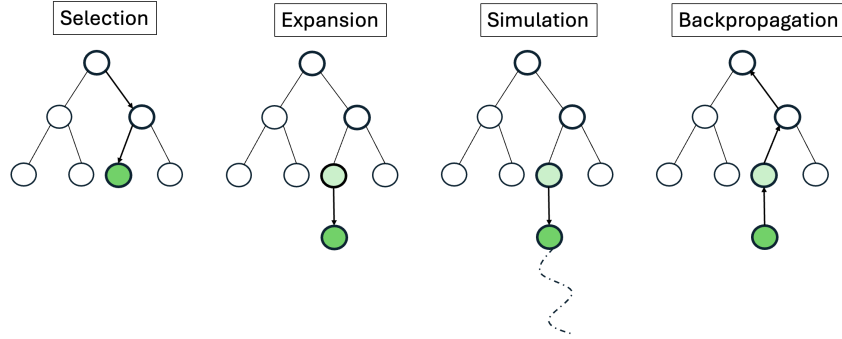


Figure 5: A pictorial illustration of Monte Carlo Tree Search

C KAGGLE CHALLENGES

The list of Kaggle challenges present in ML-Tool-Bench, and the corresponding ML problem types of each challenge are presented in Table 3

Challenge	Type
Santander Value Prediction Challenge	Regression
New York City Taxi Fare Prediction	Regression
New York City Taxi Trip Duration	Regression
Predicting the Beats-per-Minute of Songs	Regression
Predict Calorie Expenditure	Regression
Regression with a Tabular California Housing Dataset	Regression
Regression of Used Car Prices	Regression
Porto Seguro Safe Driver Prediction	Binary Classification
Costa Rican Household Poverty Prediction	Multi-Class Classification
Forest Cover Type (Kernels Only)	Multi-Class Classification
Santander Customer Transaction Prediction	Binary Classification
Binary Prediction of Poisonous Mushrooms	Binary Classification
Spaceship Titanic	Binary Classification
Binary Classification with a Bank Dataset	Binary Classification
Binary Classification with a Bank Churn Dataset	Binary Classification

Table 3: Kaggle challenges used in ML-Tool-Bench with problem type.

D RESULTS

In this section, we provide the exact consistency and performance values for each of the 15 challenges and the two models (GPT-4o and GPT-4.1-mini). Tables 4 and 5 show the consistency and leaderboard percentiles for all algorithms across all competitions in ML-Tool-Bench, for GPT-4o. Similarly, Tables 6 and 7 show the consistency and leaderboard percentiles for all algorithms across all competitions in ML-Tool-Bench, for GPT-4.1-mini.

Competition	ReAct	LATS	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS
Spaceship Titanic	0.6	0.4	0.9	0.8	0.6
Santander Value Prediction Challenge	0.9	0.6	0.5	1	1
NYC Taxi Fare Prediction	0	0.5	0.2	0.3	0.9
NYC Taxi Trip Duration	0.3	0.3	1	0.2	0.3
BPM Prediction	1	0.6	0	0.9	0.7
Calorie Expenditure Prediction	0.8	0.8	1	0.9	1
california Housing Regression	0.9	0.9	1	0.9	0.9
Used Car Prices Regression	0.9	0.4	0.4	0.9	1
Porto Seguro Safe Driver Prediction	0.3	0.5	0	0.1	0.2
Costa Rican Household Poverty Level Prediction	0.5	0.5	0.7	0.3	1
Forest Cover Type Prediction	0.6	0.9	0	0	0
Santander Customer Transaction Prediction	0.5	0.8	0.8	0.7	0.4
Poisonous Mushroom Prediction	0.9	1	1	1	0.8
Bank Deposit Classification	0.5	0.8	0.4	0.6	0.7
Bank Churn Classification	0.4	0.9	0.6	1	0.6
Overall (Median)	0.6	0.6	0.6	0.8	0.7

Table 4: Consistency across 15 competitions for five planning algorithms for GPT-4o.

Competition	ReAct	LATS	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS
Spaceship Titanic	39.54	0	22.88	59.44	62.55
Santander Value Prediction Challenge	0.09	7.17	7.12	14.87	0.27
NYC Taxi Fare Prediction	0	9.23	0	0	19.66
NYC Taxi Trip Duration	0	0	100.0	0	0
BPM Prediction	0.51	52.63	0	5.26	100
Calorie Expenditure Prediction	0.16	13.43	14.47	14.47	14.47
california Housing Regression	0.58	0.65	14.35	11.59	17.10
Used Car Prices Regression	3.0	0	0	9.35	100
Porto Seguro Safe Driver Prediction	0	5.34	0	0	0
Costa Rican Household Poverty Level Prediction	50	0	100	0	100
Forest Cover Type Prediction	0.84	23.32	0	0	0
Santander Customer Transaction Prediction	1.14	2.27	3.49	3.49	0
Poisonous Mushroom Prediction	22.75	17.61	17.69	17.62	17.62
Bank Deposit Classification	8.64	27.50	0	27.24	26.37
Bank Churn Classification	0	22.93	31.57	34.79	3.47
Overall (Median)	0.58	7.17	7.12	9.36	17.10

Table 5: Median Leaderboard percentile across 15 competitions for five planning algorithms for GPT-4o.

Competition	ReAct	LATS	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS
Spaceship Titanic	0.1	0	0.4	0.7	0.2
Santander Value Prediction Challenge	0.6	0.2	0	0.9	0.9
NYC Taxi Fare Prediction	0.2	0	0.1	0.1	0.9
NYC Taxi Trip Duration	0.1	0.1	0.3	0.5	0.9
BPM Prediction	0.4	0.1	0	1	0.8
Calorie Expenditure Prediction	0.6	0.6	0	0.9	0.9
california Housing Regression	0.7	0.3	0.3	0.8	0.8
Used Car Prices Regression	0.3	0.2	0.4	0.7	0.4
Porto Seguro Safe Driver Prediction	0.2	0.1	0.1	0.7	0.9
Costa Rican Household Poverty Level Prediction	0.2	0.2	0.1	0	0.1
Forest Cover Type Prediction	0.2	0.2	0.3	0	0.7
Santander Customer Transaction Prediction	0.3	0.5	0.1	0.7	0.7
Poisonous Mushroom Prediction	0.5	0.5	0.2	1	0.6
Bank Deposit Classification	0.5	0.3	0.1	1	0.9
Bank Churn Classification	0.5	0.8	0.5	1	0.9
Overall (Median)	0.3	0.2	0.1	0.7	0.8

Table 6: Consistency across 15 competitions for five planning algorithms for GPT-4.1-mini.

Competition	ReAct	LATS	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS
Spaceship Titanic	0	0	0	<u>62.11</u>	0
Santander Value Prediction Challenge	3.60	0	0	<u>14.24</u>	<u>14.24</u>
NYC Taxi Fare Prediction	0	0	0	0	<u>20.94</u>
NYC Taxi Trip Duration	0	0	0	1.24	<u>2.70</u>
BPM Prediction	0	0	0	<u>100</u>	<u>100</u>
Calorie Expenditure Prediction	4.31	13.41	0	<u>14.43</u>	<u>14.43</u>
california Housing Regression	<u>21.59</u>	0	0	<u>21.59</u>	<u>21.59</u>
Used Car Prices Regression	0	0	0	<u>100</u>	0
Porto Seguro Safe Driver Prediction	0	0	0	<u>17.01</u>	<u>17.01</u>
Costa Rican Household Poverty Level Prediction	0	0	0	0	0
Forest Cover Type Prediction	0	0	0	0	<u>99.44</u>
Santander Customer Transaction Prediction	0	1.17	0	<u>2.27</u>	<u>2.27</u>
Poisonous Mushroom Prediction	6.23	6.97	0	14.36	<u>16.32</u>
Bank Deposit Classification	13.19	0	0	27.24	<u>27.73</u>
Bank Churn Classification	14.66	31.09	0.94	31.57	<u>34.02</u>
Overall (Median)	0	0	0	14.43	<u>16.32</u>

Table 7: Median Leaderboard percentile across 15 competitions for five planning algorithms for GPT-4.1-mini.

Stage	Number of Tools
Data Loading	6
Data Cleaning	9
Feature Engineering	30
Modeling	10
Evaluation/Prediction	10

Table 8: Number of tools available at each stage of a Kaggle-style workflow. In total, 61 tools are provided spanning data loading, cleaning, feature engineering, and modeling. Some tools can appear in more than one stage

E TOOLS

In this section, we describe the various tools that are part of ML-Tool-Bench. Table 8 shows the number of tools in our toolset that are part of each stage in solving an ML challenge on Kaggle. Table 9 provides info about all the tools in the curated toolset provided by ML-Tool-Bench

Decorators for named references To enable tools to operate on named references rather than raw objects, we design four decorators that adapt arbitrary user-provided functions to our scratchpad interface according to their read–write behavior. We categorize tools into four types:

1. **Set tool:** saves an object to memory. Example: `read_csv` loads a dataframe and stores it under a provided name.
2. **Get tool:** reads an object from memory. Example: `get_dataframe_summary` loads a dataframe and returns a brief textual summary to guide subsequent planning.
3. **Get–Set tool:** reads an object from memory and writes a new object to memory. Example: `fit_randomforest_model` takes as input, a dataframe, and returns a fitted model.
4. **Override tool:** reads an object, returns an updated object, and overwrites the input variable binding with the returned value. Example: `cast_column` loads a dataframe and returns a modified dataframe that replaces the original.

Accordingly, we provide four decorators: `make_get_tool`, `make_set_tool`, `make_get_and_set_tool`, and `make_override_tool`, that automatically wrap user-provided tools to operate on named references and integrate with the scratchpad.

Function Signature	Description
Modeling Functions	

<code>fit_logistic_regressor(X_train, y_train, cv=5)</code>	Fit Logistic Regression model
<code>fit_linear_regressor(X_train, y_train, cv=5)</code>	Fit Linear Regression model
<code>fit_random_forest_regressor(X_train, y_train, cv=5)</code>	Fit Random Forest Regressor
<code>fit_random_forest_classifier(X_train, y_train, cv=5)</code>	Fit Random Forest Classifier
<code>fit_xgboost_regressor(X_train, y_train, cv=5)</code>	Fit XGBoost Regressor
<code>fit_xgboost_classifier(X_train, y_train, cv=5)</code>	Fit XGBoost Classifier
<code>fit_lightgbm_regressor(X_train, y_train, cv=5)</code>	Fit LightGBM Regressor
<code>fit_lightgbm_classifier(X_train, y_train, cv=5)</code>	Fit LightGBM Classifier
<code>fit_catboost_regressor(X_train, y_train, cv=5)</code>	Fit CatBoost Regressor
<code>fit_catboost_classifier(X_train, y_train, cv=5)</code>	Fit CatBoost Classifier
Data Loading Functions	
<code>read_data(filepath)</code>	Read CSV data into a pandas DataFrame
Feature Engineering and Functions to get Dataframe information	
<code>create_numeric_feature(df, name, expression)</code>	Create a numeric feature using a pandas expression
<code>create_categorical_feature(df, name, source_column, mapping)</code>	Create a categorical feature by mapping values from a source column
<code>create_conditional_feature(df, name, condition, true_value, false_value)</code>	Create a feature based on a condition
<code>extract_string_pattern(df, name, source_column, pattern, group=0)</code>	Extract pattern from string column using regex
<code>split_string_column(df, name_prefix, source_column, delimiter, max_splits=-1, indices=None)</code>	Split string column and create separate features
<code>create_group_aggregation(df, name, group_column, agg_column, agg_func)</code>	Create feature by aggregating within groups
<code>get_group_aggregation(df, group_column, agg_column, agg_func)</code>	Get aggregation result without adding it to the DataFrame
<code>create_rolling_feature(df, name, source_column, window, agg_func='mean')</code>	Create rolling window feature
<code>create_lag_feature(df, name, source_column, lag=1)</code>	Create lagged feature
<code>create_lead_feature(df, name, source_column, lead=1)</code>	Create leading feature
<code>extract_datetime_features(df, datetime_column, features=None)</code>	Extract datetime features from datetime column
<code>create_time_delta(df, name, start_column, end_column, unit='D')</code>	Create time delta feature between two date-time columns
<code>apply_custom_function(df, name, source_columns, func)</code>	Apply custom function to create feature
<code>fillna_with_value(df, columns, value)</code>	Fill missing values with a specific value
<code>fillna_with_median(df, columns=None)</code>	Fill missing values with median of the column
<code>fillna_with_mean(df, columns=None)</code>	Fill missing values with mean of the column

fillna.with.mode(df, columns=None)	Fill missing values with mode of the column
fillna.with.condition(df, target_column, condition, fill_value)	Fill missing values in a column based on a condition
fillna.with.multiple.conditions(df, target_column, conditions_and_values)	Fill missing values in a column based on multiple conditions
fillna.with.conditional.aggregation(df, target_column, condition_column, condition_values, agg_func='mean')	Fill missing values using conditional aggregation based on another column's values
fillna.with.custom.function(df, target_column, condition, custom_func)	Fill missing values using a custom function based on a condition
drop.rows.with.missing(df, columns=None, threshold=None)	Drop rows with missing values
get.missing.summary(df)	Get a summary of missing values in the DataFrame
cast.columns(df, column_type_mapping)	Cast columns to specified data types
cast.numeric.columns(df, columns=None, target_type='float')	Cast numeric columns to specified type
cast.integer.columns.to.float(df, columns=None)	Cast integer columns to float type
cast.categorical.columns(df, columns=None)	Cast categorical columns to category type
one.hot.encode(df, columns=None, drop_first=True, prefix=None)	One-hot encode categorical columns
label.encode(df, columns=None)	Label encode categorical columns
normalize.features(df, columns=None, method='standard')	Normalize numeric features
encode.all.categorical.columns(df, method='one_hot', drop_first=True)	Encode all categorical/object columns using specified method
normalize.all.numerical.columns(df, method='standard')	Normalize all numerical columns using specified method
concatenate.train.test(train_df, test_df)	Concatenate train and test data with tracking columns for proper splitting
split.combined.into.train.test(combined)	Split combined data back into train and test using tracking columns
convert.dataframe.to.features.target(df, target_column, is_train=True)	Convert DataFrame to features and target format
convert.to.dataframe(data, **kwargs)	Convert various data types to pandas DataFrame
drop.feature(df, column)	Drop feature(s) from the DataFrame
get.features(df, columns)	Extract specific features (columns) from the DataFrame
concatenate.dataframes(df1, df2, axis=0)	Concatenate two DataFrames
join.dataframes(left_df, right_df, left_on, right_on=None, how='inner', suffixes=('_x', '_y'))	Join two DataFrames using pandas merge functionality
rename.feature(df, old_name, new_name)	Rename feature(s)
get.unique.values(df, column, sort=True, include_counts=True)	Get unique values from a column as a DataFrame
get.dataframe.dtypes.summary(df)	Get comprehensive summary of the dtypes in the entire DataFrame
filter.dataframe(df, condition)	Filter DataFrame using a boolean condition
Model Utilities	
save_model(model, filepath='model.pkl')	Save the trained model to disk using pickle

<code>load_model(filepath)</code>	Load a trained model from disk using pickle
<code>save_dataframe_to_csv(df, filepath)</code>	Save a DataFrame to CSV file
Model Evaluation Functions	
<code>evaluate_regression_model(model, X_test, y_test, model_name="model", eval_data_label='test')</code>	Evaluate a trained regression model on data
<code>evaluate_classification_model(model, X_test, y_test, model_name="model", eval_data_label='test')</code>	Evaluate a trained classification model on data
<code>predict_target(model, X_data, model_name="model", return_probabilities=False)</code>	Make predictions using a trained model

Table 9: All tools in the curated toolset provided by ML-Tool-Bench.

F TOOL MASKING ABLATIONS

In this section, we perform an ablation study to investigate if tool masking contributes significantly to the performance of Hierarchical MCTS. We select a subset of five Kaggle challenges from our benchmark and evaluate all of the planning approaches, alongside a Hierarchical MCTS approach that does not use tool masking, i.e all tools are available to the agent during all the subtasks. We use GPT-4.1-mini for our experiments instead of GPT-4o for cost reasons. The results are presented in Tables 10 and 11. The results demonstrate that the performance of Hierarchical MCTS degrades substantially without tool masking. Hierarchical MCTS without tool masking, achieves a median consistency of 0.3 and a median leaderboard percentile position of 0 across the chosen subset of five challenges. In comparison, Hierarchical MCTS with tool masking achieves a median consistency and leaderboard percentile of 0.8 and 21.10 respectively. This highlights that both tool masking and subtask decomposition are critical for effectively solving long-horizon planning problems in high-dimensional action spaces using LLMs.

Competition	ReAct	LATS-Reflection	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS	Hierarchical MCTS (No Tool Masking)
Spaceship Titanic	0.1	0.1	0.1	0.7	0.5	0
Poisonous Mushroom Prediction	0.4	0.5	0.2	0.6	0.3	0.2
Bank Churn Classification	0.3	0.3	0.3	0.9	0.9	0.3
Santander Customer Transaction Prediction	0.6	0.2	0.0	0.7	0.9	0.3
NYC Taxi Fare Prediction	0.0	0.0	0.4	0.1	0.8	0.4
Overall (Median)	0.3	0.2	0.2	0.7	0.8	0.3

Table 10: Consistency across five competitions for GPT-4.1-mini with six planning algorithm variants. Maximum values per row are highlighted. The results demonstrate that the performance of Hierarchical MCTS degrades substantially without tool masking.

G COST COMPARISONS

We also provide cost comparisons for all planning algorithms using GPT-4.1-mini on the same subset of five Kaggle challenges from our benchmark that was used in the Tool Ablation Study (Appendix F). The results are reported in Table 12. LATS is the most expensive planning approach, costing 3.5× more than the more successful variants (Hierarchical-MCTS and MCTS-Shaped), while achieving only a consistency of 0.2 and a median percentile position of 0 across the five Kaggle challenges. This suggests that LATS’s search is unfocused and tends to wander due to inconsistent scoring

Competition	ReAct	LATS-Reflection	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS	Hierarchical MCTS (No Tool Masking)
Spaceship Titanic	0.0	0.0	0.0	<u>53.59</u>	28.01	0.0
Poisonous Mushroom Prediction	0.0	6.60	0.0	<u>16.99</u>	0.0	0.0
Bank Churn Classification	0.0	0.0	0.0	<u>31.57</u>	<u>31.57</u>	0.0
Santander Customer Transaction Prediction	<u>2.27</u>	0.0	0.0	<u>2.27</u>	<u>2.27</u>	0.0
NYC Taxi Fare Prediction	0.0	0.0	0.0	0.0	<u>21.10</u>	0.0
Overall (Median)	0.0	0.0	0.0	16.99	<u>21.10</u>	0.0

Table 11: Median leaderboard percentile across five competitions for GPT-4.1-mini with six planning algorithm variants. Maximum values per row are highlighted. The results demonstrate that the performance of Hierarchical MCTS degrades substantially without tool masking.

by the LLM evaluator. ReAct is the cheapest method but also performs poorly, only marginally outperforming LATS despite the latter using 10.5× more budget.

Competition	ReAct	LATS-Reflection	MCTS-Outcome	MCTS-Shaped	Hierarchical MCTS
Spaceship Titanic	1.66	<u>22.29</u>	3.22	8.44	3.93
Poisonous Mushroom Prediction	0.86	<u>13.64</u>	3.1.3	2.86	1.89
Bank Churn Classification	1.21	<u>9.98</u>	1.44	1.79	5.65
Santander Customer Transaction Prediction	1.34	<u>11.28</u>	2.36	1.43	5.38
NYC Taxi Fare Prediction	2.02	<u>17.83</u>	2.72	4.07	5.56
Overall (Sum)	7.08	<u>75.02</u>	12.86	18.59	22.42

Table 12: Total costs (\$) aggregated over 10 trajectories for each of the five competitions using GPT-4.1-mini under all planning algorithm variants examined in this study. For each row, the maximum value is highlighted

H PROMPTS

For each competition in our benchmark, we constructed a standardized instruction template that included: (i) a brief description of the Kaggle challenge, (ii) a description of the data fields, and (iii) additional requirements specified from our end regarding model training and submission. For the challenge and data field descriptions, we used only the information provided in the corresponding sections of the original Kaggle competition; no human-authored modifications or additions were introduced.

Below, we show the exact template used for the *Spaceship Titanic* challenge (Howard et al., 2022):

Welcome to the year 2912, where your data science skills are needed to solve a cosmic mystery. We’ve received a transmission from four lightyears away and things aren’t looking good.

The Spaceship Titanic was an interstellar passenger liner launched a month ago. With almost 13,000 passengers on board, the vessel set out on its maiden voyage transporting emigrants from our solar system to three newly habitable exoplanets orbiting nearby stars.

While rounding Alpha Centauri en route to its first destination—the torrid 55 Cancri E—the unwary Spaceship Titanic collided with a spacetime anomaly hidden within a dust cloud. Sadly, it met a similar fate as its namesake from 1000 years before. Though the ship stayed intact, almost half of the passengers were transported to an alternate dimension! Read the data from the train.csv file in the data folder.

To help rescue crews and retrieve the lost passengers, you are challenged to predict which passengers were transported by the anomaly using records recovered from the spaceship’s damaged computer system.

In this competition your task is to predict whether a passenger was transported to an alternate dimension during the Spaceship Titanic’s collision with the space-time anomaly. To help you make these predictions, you’re given a set of personal records recovered from the ship’s damaged computer system.

File and Data Field Descriptions

- `PassengerId` – Unique identifier in the form `gggg-pp`, where `gggg` indicates a group and `pp` a member index.
- `HomePlanet` – Planet of permanent residence.
- `CryoSleep` – Whether the passenger elected suspended animation.
- `Cabin` – Cabin number, formatted as `deck/num/side`.
- `Destination` – Planet of disembarkation.
- `Age` – Age of passenger.
- `VIP` – Whether the passenger paid for VIP services.
- `RoomService`, `FoodCourt`, `ShoppingMall`, `Spa`, `VRDeck` – Spending at onboard amenities.
- `Name` – Full passenger name.
- `Transported` – Target variable — whether the passenger was transported.

Submission File Format

- `PassengerId` – Identifier for each passenger in the test set.
- `Transported` – Prediction (True/False).

Benchmark Instructions

- The training data is located at `data/spaceship_titanic/train.csv`.
- The test data is located at `data/spaceship_titanic/test.csv`.
- Load, clean, and perform feature engineering before fitting models.
- Concatenate train and test datasets before preprocessing to ensure consistent transformations, then split back.
- Experiment with multiple models and hyperparameter tuning to find the best-performing solution.
- Report evaluation results demonstrating model fit.
- Save the best model to `model_saves/spaceship_titanic/`.
- Save predictions for the test set in CSV format to `{save_directory}/{save_file_name}.csv`.

H.1 LATS

For LATS (Zhou et al., 2024a), we used the following system prompt for the reflection step, where an LLM was asked to evaluate the quality of a given trajectory (i.e., all message contents of the trajectory were passed as input):

You are a Data Science judge, who evaluates the goodness of tool calling trajectories to solve Machine Learning tasks on Kaggle. Reflect and grade the agent’s trajectory plan for the provided challenge. The trajectories should be aimed towards solving the challenge, i.e., generating a trained model and a valid submission file. Keep your reflections concise and to the point.

For the expansion stage, we used the following system prompt. The trajectory of messages up to the current node was passed as input, and the agent was asked to propose new expansion candidates:

You are a Data Scientist tasked with solving the Kaggle competition provided below, with the tools available to you. Propose tool candidates that would help solve the problem at the current stage.

H.2 MCTS

For both MCTS-Outcome and MCTS-Shaped, we used the same system prompt in the expansion stage, as was used for LATS. The trajectory of messages up to the current node was passed as input, and the agent was asked to propose new expansion candidates (same as LATS)

H.3 HIERARCHICAL MCTS

For Hierarchical-MCTS, to propose candidates during the expansion phase, we used subtask-specific prefixes for the system prompt. Each subtask was associated with a descriptive prefix that constrained the role of the agent and defined the completion condition for that stage. The list of prefixes is shown below:

- `train_data_loading`: “You are a Data Scientist in the Data Loading stage of solving a Kaggle challenge, using only the tools available to you. This stage ends when you have loaded the train data successfully.”
- `test_data_loading`: “You are a Data Scientist in the Data Loading stage of solving a Kaggle challenge, using only the tools available to you. This stage ends when you have loaded the test data successfully.”
- `combine_train_test`: “You are a Data Scientist in the Data Loading stage of solving a Kaggle challenge, using only the tools available to you. This stage ends when you have combined the train and test data into a single dataframe successfully, to be used for downstream Data Cleaning and Feature Engineering.”
- `data_cleaning`: “You are a Data Scientist in the Data Cleaning stage of solving a Kaggle challenge, using only the tools available to you. This stage ends when there are no missing values present in the data. This also includes the column corresponding to the target variable, that may have NaNs in the test partition since the target variable is not present in the test partition. You are allowed to be innovative in filling the missing values based on your understanding of the data.”
- `feature_engineering`: “You are a Data Scientist in the Feature Engineering stage of solving a Kaggle challenge, using only the tools available to you. Create new features, or delete unimportant features or transform existing features as needed. You are not allowed to delete or modify features that indicate if the row in the data belongs to the train or test partition. You are also not allowed to augment the feature corresponding to the target variable. Use your understanding of the data to aid your decisions. This stage ends when the models feel that the features are good enough for modeling, and categorical and numerical features have been properly encoded. After the end of this stage, all the features should be (i) either int or float or (ii) int, float, category with the number of unique values in the category columns not being exorbitantly large.”
- `split_train_test`: “You are a Data Scientist in the Split Train Test stage of solving a Kaggle challenge, using only the tools available to you. Split the combined train and test data into train and test dataframes. This stage ends when the train and test dataframes are successfully split from the combined dataframe.”
- `train_data_to_features_target`: “You are a Data Scientist in the Converting the Train Data to Features and Target stage of solving a Kaggle challenge, using only the tools available to you. Convert the train data into features and target. This stage ends when the train data is successfully converted into features and target, for making downstream modeling upon.”
- `test_data_to_features`: “You are a Data Scientist in the Converting the Test Data to Features stage of solving a Kaggle challenge, using only the tools available to you. Convert the test data into features. This stage ends when the test data is successfully converted into features, for making downstream predictions upon.”
- `modeling`: “You are a Data Scientist in the Modeling stage of solving a Kaggle challenge, using only the tools available to you. Train and tune models. You might need to experiment with different model choices and properly tune your hyperparameters to get good performance. Use the provided evaluation tools to evaluate your trained models if

needed. This stage ends when the agent has successfully created a model that it considers to be the best.”

- `create_submission_dataframe`: “You are a Data Scientist in the Create Submission stage of solving a Kaggle challenge, using only the tools available to you. Make predictions on the test data, and create a submission dataframe that contains the predictions in the requested format. This stage ends when the submission dataframe in the correct format is created.”

The system prompt was then constructed as:

{`subtask_description`}. Propose tool candidates that would help solve the problem at the current stage.

The trajectory of messages up to the current node was passed as input, and the agent was asked to propose new expansion candidates (same as what was used for LATS, MCTS variants)

I VISUALIZING TRAJECTORIES

In this section, we present example trajectories for each of the planning algorithms evaluated in this paper. For every algorithm, we provide both a successful and a failed trajectory. All trajectories are collected on the Spaceship Titanic Kaggle challenge using the GPT-4.1-mini model, ensuring a consistent evaluation setting across planning methods.

I.1 REACT

I.1.1 FAILED TRAJECTORY

The following failed ReAct trajectory contains 12 tool calls and terminates during the Feature Engineering stage, never reaching the Modeling stage. This pattern is common among failed ReAct runs: a substantial fraction are unable to complete the full end-to-end machine learning pipeline required for the Kaggle task. This snippet corresponds to the final portion of the failed ReAct trajectory.

```
{
  "step_number": 20,
  "timestamp": "2025-11-20T18:52:06.565366",
  "step_type": "tool_execution_initiation",
  "action": "tool_execution_started",
  "tools_to_execute": [
    {
      "tool_name": "apply_custom_function",
      "tool_args": {
        "bindings": { "df": "combined_data" },
        "output": "combined_data",
        "func_kwargs": {
          "name": "Cabin_deck",
          "source_columns": "Cabin",
          "func": "lambda x: x.split('/')[0] if isinstance(x, str) else
None"
        }
      }
    },
    "tool_id": "call_s5ZjVgGT9j6llr3Vg7jUQavA"
  ]
},
{
  "step_number": 21,
  "timestamp": "2025-11-20T18:52:06.599497",
  "step_type": "tool_execution_completion",
  "action": "tool_execution_completed",
  "tool_results": [
    {
```

```

1242         "message_type": "ToolMessage",
1243         "content_preview": "Error: TypeError(\"'str' object is not
1244         callable\")\n Please fix your mistakes.",
1245         "content_length": 75
1246     }
1247 ],
1248 {
1249     "step_number": 22,
1250     "timestamp": "2025-11-20T18:52:10.291750",
1251     "step_type": "tool_execution_initiation",
1252     "action": "tool_execution_started",
1253     "tools_to_execute": [
1254         {
1255             "tool_name": "apply_custom_function",
1256             "tool_args": {
1257                 "bindings": { "df": "combined_data" },
1258                 "output": "combined_data",
1259                 "func_kwargs": {
1260                     "name": "Cabin_deck",
1261                     "source_columns": "Cabin",
1262                     "func": null
1263                 }
1264             },
1265             "tool_id": "call_JiDjhYHW16xgodP4y69hOAR3"
1266         }
1267     ],
1268 },
1269 {
1270     "step_number": 23,
1271     "timestamp": "2025-11-20T18:52:10.324847",
1272     "step_type": "tool_execution_completion",
1273     "action": "tool_execution_completed",
1274     "tool_results": [
1275         {
1276             "message_type": "ToolMessage",
1277             "content_preview": "Error: TypeError(\"'NoneType' object is not
1278             callable\")\n Please fix your mistakes.",
1279             "content_length": 80
1280         }
1281     ]
1282 }

```

1279 I.1.2 SUCCESSFUL TRAJECTORY

1280 This trajectory comprises 40 tool calls and demonstrates strong performance. The feature engi-
1281 neering stage involves numerous successful and relevant tool invocations. In contrast, the modeling
1282 stage fits only a CatBoost classifier without exploring additional model families. When scored on
1283 our internal test set, this trajectory attains a percentile of 54.92 on the Kaggle leaderboard for this
1284 task.

```

1285 {
1286     "step_number": 70,
1287     "timestamp": "2025-11-20T18:53:09.461595",
1288     "step_type": "tool_execution_initiation",
1289     "action": "tool_execution_started",
1290     "tools_to_execute": [
1291         {
1292             "tool_name": "convert_dataframe_to_features_target",
1293             "tool_args": {
1294                 "bindings": {
1295                     "df": "test_df"
1296                 },
1297                 "func_kwargs": {

```



```

1296         "target_column": "Transported",
1297         "is_train": false
1298     },
1299     "output": "test_features_target"
1300 },
1301 "tool_id": "call_VTa2DIpTuPJRLydhruqaaiBK"
1302 }
1303 ],
1304 {
1305     "step_number": 71,
1306     "timestamp": "2025-11-20T18:53:09.535408",
1307     "step_type": "tool_execution_completion",
1308     "action": "tool_execution_completed",
1309     "tool_results": [
1310         {
1311             "message_type": "ToolMessage",
1312             "content_preview": "Applied convert_dataframe_to_features..",
1313             "content_length": 864
1314         }
1315     ],
1316 },
1317 {
1318     "step_number": 72,
1319     "timestamp": "2025-11-20T18:53:11.188346",
1320     "step_type": "tool_execution_initiation",
1321     "action": "tool_execution_started",
1322     "tools_to_execute": [
1323         {
1324             "tool_name": "tune_catboost_classifier",
1325             "tool_args": {
1326                 "bindings": {
1327                     "X_train": "X_train",
1328                     "y_train": "Y_train"
1329                 },
1330                 "output": "catboost_tuning_result"
1331             },
1332             "tool_id": "call_DHeRtQVO4RZb8ayZyc47Ow0t"
1333         }
1334     ],
1335 },
1336 {
1337     "step_number": 73,
1338     "timestamp": "2025-11-20T18:53:12.310197",
1339     "step_type": "tool_execution_completion",
1340     "action": "tool_execution_completed",
1341     "tool_results": [
1342         {
1343             "message_type": "ToolMessage",
1344             "content_preview": "Applied tune_catboost_classifier. The Best",
1345             "content_length": 1088
1346         }
1347     ],
1348 },
1349 {
1350     "step_number": 74,
1351     "timestamp": "2025-11-20T18:53:14.269490",
1352     "step_type": "tool_execution_initiation",
1353     "action": "tool_execution_started",
1354     "tools_to_execute": [
1355         {

```

```

1350         "tool_name": "evaluate_classification_model",
1351         "tool_args": {
1352             "bindings": {
1353                 "model": "best_estimator",
1354                 "X_test": "X_train",
1355                 "y_test": "Y_train"
1356             },
1357             "func_kwargs": {
1358                 "model_name": "CatBoost Classifier",
1359                 "eval_data_label": "train"
1360             }
1361         },
1362         "tool_id": "call_76EXOjjGyUp86wMwe3bVbES4"
1363     },
1364     {
1365         "step_number": 75,
1366         "timestamp": "2025-11-20T18:53:14.363633",
1367         "step_type": "tool_execution_completion",
1368         "action": "tool_execution_completed",
1369         "tool_results": [
1370             {
1371                 "message_type": "ToolMessage",
1372                 "content_preview": "Applied evaluate_classification_model..
1373                 [truncated]",
1374                 "content_length": 1753
1375             }
1376         ],
1377     },
1378     {
1379         "step_number": 76,
1380         "timestamp": "2025-11-20T18:53:16.071819",
1381         "step_type": "tool_execution_initiation",
1382         "action": "tool_execution_started",
1383         "tools_to_execute": [
1384             {
1385                 "tool_name": "predict_target",
1386                 "tool_args": {
1387                     "bindings": {
1388                         "model": "best_estimator",
1389                         "X_data": "X_test"
1390                     },
1391                     "func_kwargs": {
1392                         "model_name": "CatBoost Classifier"
1393                     },
1394                     "output": "test_predictions"
1395                 },
1396                 "tool_id": "call_fM4PS3vRAthZWODZw6ycGzUQ"
1397             }
1398         ],
1399     },
1400     {
1401         "step_number": 77,
1402         "timestamp": "2025-11-20T18:53:16.166858",
1403         "step_type": "tool_execution_completion",
1404         "action": "tool_execution_completed",
1405         "tool_results": [
1406             {
1407                 "message_type": "ToolMessage",
1408                 "content_preview": "Applied predict_target ..[truncated]",
1409                 "content_length": 1019
1410             }
1411         ],
1412     },

```

```

1404     {
1405         "step_number": 78,
1406         "timestamp": "2025-11-20T18:53:26.236014",
1407         "step_type": "tool_execution_initiation",
1408         "action": "tool_execution_started",
1409         "tools_to_execute": [
1410             {
1411                 "tool_name": "save_model",
1412                 "tool_args": {
1413                     "bindings": {
1414                         "model": "best_estimator"
1415                     },
1416                     "func_kwargs": {
1417                         "filepath":
1418 "model_saves/spaceship_titanic/catboost_best_model.pkl"
1419                     }
1420                 },
1421                 "tool_id": "call_PI0hV7OnLQEBSrir8PNgWSU0"
1422             }
1423         ],
1424         {
1425             "step_number": 79,
1426             "timestamp": "2025-11-20T18:53:26.322763",
1427             "step_type": "tool_execution_completion",
1428             "action": "tool_execution_completed",
1429             "tool_results": [
1430                 {
1431                     "message_type": "ToolMessage",
1432                     "content_preview": "Applied save_model ..",
1433                     "content_length": 676
1434                 }
1435             ]
1436         },
1437         {
1438             "step_number": 80,
1439             "timestamp": "2025-11-20T18:53:29.230047",
1440             "step_type": "tool_execution_initiation",
1441             "action": "tool_execution_started",
1442             "tools_to_execute": [
1443                 {
1444                     "tool_name": "save_dataframe_to_csv",
1445                     "tool_args": {
1446                         "bindings": {
1447                             "df": "test_predictions"
1448                         },
1449                         "func_kwargs": {
1450                             "filepath":
1451 "/submissions/spaceship_titanic/react/react_agent_4_47a197f8-edb6-4c43-a9ce-ae5f8ef45736.c
1452                     },
1453                     "tool_id": "call_xrfZg9Pg0FMXxgpr7a2rAGxf"
1454                 }
1455             ]
1456         },
1457         {
1458             "step_number": 81,
1459             "timestamp": "2025-11-20T18:53:29.329762",
1460             "step_type": "tool_execution_completion",
1461             "action": "tool_execution_completed",
1462             "tool_results": [
1463                 {
1464                     "message_type": "ToolMessage",
1465                     "content_preview": "Applied save_dataframe_to_csv..",
1466                     "content_length": 777
1467                 }
1468             ]
1469         }
1470     ]
1471 }

```

```

1458     }
1459   ]
1460 },
1461 {
1462   "step_number": 82,
1463   "timestamp": "2025-11-20T18:53:34.098547",
1464   "step_type": "execution_summary",
1465   "action": "agent_execution_completed",
1466   "total_execution_time": 114.67485404014587,
1467   "total_tokens": 1022157,
1468   "total_cost": 0.156462,
1469   "final_message_count": 81,
1470   "competition_name": "spaceship_titanic"
1471 }
1472 ]
1473 }

```

I.2 LATS

I.3 FAILED TRAJECTORY

This trajectory contains 17 tool calls. Although this trajectory, which denotes the agent's most promising branch, successfully fit a model, LATS ultimately prioritized other failed branches, causing it to miss this viable solution path. As a result, the final output did not include any valid trajectory that solved the challenge. This failure stems from inconsistencies in the reflection-based rating process, which incorrectly ranked failed nodes above successful progress.

```

1482 {
1483   "step_number": 37,
1484   "message_type": "AIMessage",
1485   "timestamp": "2025-11-21T00:42:26.668798",
1486   "step_type": "tool_selection",
1487   "action": "selected_tools_for_execution",
1488   "tools_selected": [
1489     "tune_lightgbm_classifier"
1490   ],
1491   "tool_calls_detail": [
1492     {
1493       "name": "tune_lightgbm_classifier",
1494       "args": {
1495         "bindings": {
1496           "X_train": "X_train",
1497           "y_train": "Y_train"
1498         },
1499         "output": "best_lgbm_model"
1500       },
1501       "id": "call_iff8Jho9Tc91alv1bEPfxIsZ",
1502       "type": "tool_call"
1503     }
1504   ],
1505   "content": "Features and target variables are now separated. I will proceed with hyperparameter tuning for a classification model to find the best model for this task. I will start with a LightGBM Classifier tuned using cross-validation."
1506 },
1507 {
1508   "step_number": 38,
1509   "message_type": "ToolMessage",
1510   "timestamp": "2025-11-21T00:42:26.668802",
1511   "step_type": "tool_result",
1512   "action": "received_tool_output",
1513   "content_preview": "Error: ValueError('Input y contains NaN.')\nPlease fix your mistakes.",
1514   "content_length": 69

```

```

1512     },
1513     {
1514         "step_number": 39,
1515         "message_type": "HumanMessage",
1516         "timestamp": "2025-11-21T00:42:26.668807",
1517         "step_type": "reflection",
1518         "action": "llm_reflection",
1519         "content_preview": "Reasoning: The agent failed to check for and
1520 handle missing target values before model tuning, resulting in a
1521 ValueError due to NaNs in the target variable. This indicates a gap
1522 in data cleaning and validation steps before model training. Proper
1523 imputation or removal of missing target rows is necess...",
1524         "content_length": 704,
1525         "extracted_score": 4.0,
1526         "full_reflection_content": "Reasoning: The agent failed to check
1527 for and handle missing target values before model tuning, resulting
1528 in a ValueError due to NaNs in the target variable. This indicates a
1529 gap in data cleaning and validation steps before model training.
1530 Proper imputation or removal of missing target rows is necessary.
1531 The agent should have performed or verified missing value handling
1532 comprehensively before tuning models. This oversight blocks progress
1533 in model fitting and prediction generation, which are critical to
1534 solving the task. The trajectory is incomplete and insufficient as
1535 no valid model training was achieved due to missing data issues.
1536 Addressing this gap is essential for a correct solution.\nScore: 4"
1537     },
1538     {
1539         "step_number": 40,
1540         "message_type": "AIMessage",
1541         "timestamp": "2025-11-21T00:42:26.668821",
1542         "step_type": "tool_selection",
1543         "action": "selected_tools_for_execution",
1544         "tools_selected": [
1545             "filter_dataframe"
1546         ],
1547         "tool_calls_detail": [
1548             {
1549                 "name": "filter_dataframe",
1550                 "args": {
1551                     "bindings": {
1552                         "df": "combined_df"
1553                     },
1554                     "output": "combined_train_df",
1555                     "func_kwargs": {
1556                         "condition": "Transported.notna()"
1557                     }
1558                 },
1559                 "id": "call_seR8LUmqkK389bcayn0QbxWI",
1560                 "type": "tool_call"
1561             }
1562         ],
1563         "content": "The model tuning failed because there are missing
1564 values in the target column \"Transported\". I need to clean the
1565 data by removing rows with missing target values before proceeding
1566 with model training. I will filter out such rows from the combined
1567 dataset and then separate features and target again."
1568     },
1569     {
1570         "step_number": 41,
1571         "message_type": "ToolMessage",
1572         "timestamp": "2025-11-21T00:42:26.668824",
1573         "step_type": "tool_result",
1574         "action": "received_tool_output",
1575         "content_preview": "Applied filter_dataframe with docstring: \n
1576 Filter DataFrame using a boolean condition.\n      \n      Parameters:\n

```

```

1566         -----\n      df : pd.DataFrame\n                Input DataFrame\n1567 condition : str\n                Boolean condition expression (e.g., 'coll >
1568 0', 'coll == \"value\"', 'coll > 0 and coll < 100')\n      \n...",
1569 "content_length": 1141
1570 },
1571 {
1572     "step_number": 42,
1573     "message_type": "HumanMessage",
1574     "timestamp": "2025-11-21T00:42:26.668830",
1575     "step_type": "reflection",
1576     "action": "llm_reflection",
1577     "content_preview": "Reasoning: The agent initially struggled with
1578 the correct function call syntax for loading and manipulating data,
1579 causing delays in early steps. It successfully loaded and
1580 concatenated the train and test datasets, summarized missing values,
1581 cast categorical columns, and separated features and target...",
1582     "content_length": 850,
1583     "extracted_score": 5.0,
1584     "full_reflection_content": "Reasoning: The agent initially
1585 struggled with the correct function call syntax for loading and
1586 manipulating data, causing delays in early steps. It successfully
1587 loaded and concatenated the train and test datasets, summarized
1588 missing values, cast categorical columns, and separated features and
1589 target variables. However, the agent overlooked checking and
1590 handling missing target values before model training, leading to a
1591 model fitting failure. Once identified, it correctly filtered out
1592 rows with missing targets to fix this issue. The trajectory shows
1593 gradual improvement and understanding but is currently incomplete as
1594 it has not proceeded to feature engineering, model training,
1595 evaluation, or prediction saving yet. The agent needs to continue
1596 systematically through these steps to fully solve the
1597 task.\n\nScore: 5\nFound_solution: False\nScore: 5"
1598 },
1599 {
1600     "step_number": 43,
1601     "message_type": "AIMessage",
1602     "timestamp": "2025-11-21T00:42:26.668845",
1603     "step_type": "tool_selection",
1604     "action": "selected_tools_for_execution",
1605     "tools_selected": [
1606         "convert_dataframe_to_features_target"
1607     ],
1608     "tool_calls_detail": [
1609         {
1610             "name": "convert_dataframe_to_features_target",
1611             "args": {
1612                 "bindings": {
1613                     "df": "combined_train_df"
1614                 },
1615                 "output": "features_target_tuple_cleaned",
1616                 "func_kwargs": {
1617                     "target_column": "Transported",
1618                     "is_train": true
1619                 }
1620             },
1621             "id": "call_rm2jT91kE90EoOciGAwUXnco",
1622             "type": "tool_call"
1623         }
1624     ],
1625     "content": ""
1626 },
1627 {
1628     "step_number": 44,
1629     "message_type": "ToolMessage",
1630     "timestamp": "2025-11-21T00:42:26.668849",

```

```

1620         "step_type": "tool_result",
1621         "action": "received_tool_output",
1622         "content_preview": "Applied convert_dataframe_to_features_target
1623 with docstring: \n      Convert DataFrame to features and target
1624 format.\n      \n      Parameters:\n      -----\n      df :
1625 pd.DataFrame\n      Input DataFrame\n      target_column : str\n
1626 Name of the target column\n      is_train : bool\n      Whether
1627 the Data...",
1628         "content_length": 888
1629     },
1630     {
1631         "step_number": 45,
1632         "message_type": "HumanMessage",
1633         "timestamp": "2025-11-21T00:42:26.668854",
1634         "step_type": "reflection",
1635         "action": "llm_reflection",
1636         "content_preview": "Reasoning: The trajectory initially suffered
1637 from repeated syntax errors in function calls, blocking early
1638 progress on data loading. After correcting these, the agent
1639 effectively loaded train and test data, summarized missing values,
1640 concatenated datasets, and cast categorical columns properly. A
1641 ke...",
1642         "content_length": 949,
1643         "extracted_score": 6.0,
1644         "full_reflection_content": "Reasoning: The trajectory initially
1645 suffered from repeated syntax errors in function calls, blocking
1646 early progress on data loading. After correcting these, the agent
1647 effectively loaded train and test data, summarized missing values,
1648 concatenated datasets, and cast categorical columns properly. A key
1649 oversight was the presence of missing target values, which caused
1650 model tuning to fail. The agent correctly identified this and
1651 filtered out rows with missing targets before proceeding. The agent
1652 has now prepared clean training data and is ready for model tuning.
1653 However, feature engineering and model evaluation steps are still
1654 pending. The trajectory shows a solid recovery from early mistakes
1655 but remains incomplete in solving the full task of generating a
1656 trained model and valid predictions. Continued systematic progress
1657 through cleaning, feature engineering, modeling, and prediction
1658 saving is needed to fully solve the challenge.\nScore: 6"
1659     },
1660     {
1661         "step_number": 46,
1662         "message_type": "AIMessage",
1663         "timestamp": "2025-11-21T00:42:26.668870",
1664         "step_type": "tool_selection",
1665         "action": "selected_tools_for_execution",
1666         "tools_selected": [
1667             "tune_lightgbm_classifier"
1668         ],
1669         "tool_calls_detail": [
1670             {
1671                 "name": "tune_lightgbm_classifier",
1672                 "args": {
1673                     "bindings": {
1674                         "X_train": "X_train",
1675                         "y_train": "Y_train"
1676                     },
1677                     "output": "best_lgbm_model"
1678                 },
1679                 "id": "call_SraYlT5IClOasjkTW9oju8xt",
1680                 "type": "tool_call"
1681             }
1682         ],
1683         "content": ""
1684     },

```

```

1674 {
1675     "step_number": 47,
1676     "message_type": "ToolMessage",
1677     "timestamp": "2025-11-21T00:42:26.668873",
1678     "step_type": "tool_result",
1679     "action": "received_tool_output",
1680     "content_preview": "Applied tune_lightgbm_classifier with
1681 docstring: \n    Perform hyperparameter tuning for LightGBM
1682 Classifier using GridSearchCV.\n\n    Parameters\n    -----\n
1683 X_train : pd.DataFrame\n    Training features\n    y_train :
1684 pd.Series\n    Training target\n    cv : int, default=5\n
1685 Number o...",
1686     "content_length": 1092
1687 },
1688 {
1689     "step_number": 48,
1690     "message_type": "HumanMessage",
1691     "timestamp": "2025-11-21T00:42:26.668878",
1692     "step_type": "reflection",
1693     "action": "llm_reflection",
1694     "content_preview": "Reasoning: The agent initially struggled with
1695 the correct syntax for function calls, which delayed progress in
1696 loading and concatenating data. However, the agent quickly recovered
1697 by successfully loading both train and test datasets, summarizing
1698 missing values, concatenating datasets, and correctly ...",
1699     "content_length": 1163,
1700     "extracted_score": 6.0,
1701     "full_reflection_content": "Reasoning: The agent initially
1702 struggled with the correct syntax for function calls, which delayed
1703 progress in loading and concatenating data. However, the agent
1704 quickly recovered by successfully loading both train and test
1705 datasets, summarizing missing values, concatenating datasets, and
1706 correctly casting categorical columns. A notable oversight was the
1707 failure to check for missing target values before model tuning,
1708 which resulted in an error. The agent fixed this by filtering out
1709 rows with missing targets. Subsequently, the agent performed
1710 hyperparameter tuning for a LightGBM classifier and obtained a best
1711 CV score of about 0.795, which is a good baseline for the
1712 classification task. The agent has established a good foundation but
1713 has not clearly shown subsequent steps such as evaluation on a
1714 validation set, feature engineering enhancements, saving the best
1715 model, or generating and saving predictions for the test set
1716 submission file. To fully solve the challenge, these remaining steps
1717 are needed. Overall, the trajectory shows reasonable understanding
1718 and recovery from initial mistakes but remains incomplete in
1719 finalizing the solution.\nScore: 6"
1720 },
1721 },

```

1715 I.3.1 SUCCESSFUL TRAJECTORY

1717 This successful trajectory comprises 27 tool calls and attains a leaderboard percentile of 58.21 on
1718 our test split. It is worth noting that LATS exclusively employed the LightGBM classifier and did
1719 not consider other modeling alternatives.

```

1720 {
1721     "step_number": 61,
1722     "message_type": "AIMessage",
1723     "timestamp": "2025-11-20T23:52:58.042763",
1724     "step_type": "tool_selection",
1725     "action": "selected_tools_for_execution",
1726     "tools_selected": [
1727         "convert_dataframe_to_features_target"
1728     ],
1729     "tool_calls_detail": [

```



```

1728     {
1729         "name": "convert_dataframe_to_features_target",
1730         "args": {
1731             "bindings": {
1732                 "df": "test_df"
1733             },
1734             "func_kwargs": {
1735                 "target_column": "Transported",
1736                 "is_train": false
1737             },
1738             "output": "X_test"
1739         },
1740         "id": "call_LNz1Q2BWh1hYxqkPIbk29ojD",
1741         "type": "tool_call"
1742     }
1743 ],
1744 "content": ""
1745 },
1746 {
1747     "step_number": 62,
1748     "message_type": "ToolMessage",
1749     "timestamp": "2025-11-20T23:52:58.042766",
1750     "step_type": "tool_result",
1751     "action": "received_tool_output",
1752     "content_preview": "Applied convert_dataframe_to_features_target
1753 with docstring: \n    Convert DataFrame to features and target
1754 format.\n    \n    Parameters:\n    -----\n    df :
1755 pd.DataFrame\n    Input DataFrame\n    target_column : str\n
1756 Name of the target column\n    is_train : bool\n    Whether
1757 the Data...",
1758     "content_length": 863
1759 },
1760 {
1761     "step_number": 63,
1762     "message_type": "HumanMessage",
1763     "timestamp": "2025-11-20T23:52:58.042770",
1764     "step_type": "reflection",
1765     "action": "llm_reflection",
1766     "content_preview": "Reasoning: The agent followed a generally
1767 appropriate and methodical trajectory, starting with loading and
1768 concatenating train and test datasets to allow consistent
1769 preprocessing. They identified and handled missing values using
1770 sensible imputation methods for categorical and numerical features.
1771 The...",
1772     "content_length": 1404,
1773     "extracted_score": 6.0,
1774     "full_reflection_content": "Reasoning: The agent followed a
1775 generally appropriate and methodical trajectory, starting with
1776 loading and concatenating train and test datasets to allow
1777 consistent preprocessing. They identified and handled missing values
1778 using sensible imputation methods for categorical and numerical
1779 features. The agent cast categorical columns correctly and performed
1780 meaningful feature engineering, including extracting cabin features,
1781 aggregating spending features, and computing group features from
1782 PassengerId, demonstrating good domain understanding.\n\nHowever,
1783 the agent made some critical mistakes in the execution, such as not
1784 splitting the combined dataset back into train and test before model
1785 training, causing NaNs in the target labels and failing model tuning
1786 initially. They corrected this splitting issue later but also had
1787 minor syntax errors in tool usage. The agent has yet to demonstrate
1788 the crucial steps of fitting models, hyperparameter tuning,
1789 evaluation, saving the best model, and generating and saving
1790 predictions in the required submission format and
1791 location.\n\nOverall, the trajectory is solid for data loading,
1792 cleaning, and feature engineering but incomplete for the full

```

```

1782 challenge solution. The agent needs to proceed with modeling and
1783 output saving to fully solve the task. The response quality is good
1784 for preliminary steps but insufficient for a complete
1785 solution.\n\nScore: 6\nScore: 6"
1786 },
1787 {
1788     "step_number": 64,
1789     "message_type": "AIMessage",
1790     "timestamp": "2025-11-20T23:52:58.042793",
1791     "step_type": "tool_selection",
1792     "action": "selected_tools_for_execution",
1793     "tools_selected": [
1794         "tune_lightgbm_classifier"
1795     ],
1796     "tool_calls_detail": [
1797         {
1798             "name": "tune_lightgbm_classifier",
1799             "args": {
1800                 "bindings": {
1801                     "X_train": "X_train",
1802                     "y_train": "Y_train"
1803                 },
1804                 "output": "best_lgbm_model"
1805             },
1806             "id": "call_WZihuqBOZNl9QcPJNoHlPhVo",
1807             "type": "tool_call"
1808         }
1809     ],
1810     "content": ""
1811 },
1812 {
1813     "step_number": 65,
1814     "message_type": "ToolMessage",
1815     "timestamp": "2025-11-20T23:52:58.042796",
1816     "step_type": "tool_result",
1817     "action": "received_tool_output",
1818     "content_preview": "Applied tune_lightgbm_classifier with
1819 docstring: \n    Perform hyperparameter tuning for LightGBM
1820 Classifier using GridSearchCV.\n\n    Parameters\n    -----\n
1821 X_train : pd.DataFrame\n    Training features\n    y_train :
1822 pd.Series\n    Training target\n    cv : int, default=5\n
1823 Number o...",
1824     "content_length": 1092
1825 },
1826 {
1827     "step_number": 66,
1828     "message_type": "HumanMessage",
1829     "timestamp": "2025-11-20T23:52:58.042802",
1830     "step_type": "reflection",
1831     "action": "llm_reflection",
1832     "content_preview": "Reasoning: The agent has followed a generally
1833 sound and methodical approach suitable for the Spaceship Titanic
1834 classification task. They correctly loaded and concatenated the
1835 datasets, handled missing values with sensible imputation, cast
1836 categorical variables, and performed meaningful feature engin...",
1837     "content_length": 1162,
1838     "extracted_score": 7.0,
1839     "full_reflection_content": "Reasoning: The agent has followed a
1840 generally sound and methodical approach suitable for the Spaceship
1841 Titanic classification task. They correctly loaded and concatenated
1842 the datasets, handled missing values with sensible imputation, cast
1843 categorical variables, and performed meaningful feature engineering
1844 such as extracting cabin features and group-related variables. The
1845 agent made some mistakes in tool usage and workflow order, notably
1846 not splitting combined data before modeling, causing initial tuning

```

errors due to NaNs in the target. These were corrected later. The agent successfully performed hyperparameter tuning on a LightGBM classifier and obtained a reasonable CV score (~0.79), indicating a decent model fit. However, the agent has not yet shown the steps to save the best model or generate and save predictions on the test set as required. The trajectory is strong in data preparation and modeling but incomplete on final output generation and saving, which is essential for solving the task fully. Overall, the response is effective but partial, scoring high on foundational steps but lacking completion of the full solution pipeline.\nScore: 7"

```

},
{
  "step_number": 67,
  "message_type": "AIMessage",
  "timestamp": "2025-11-20T23:52:58.042822",
  "step_type": "tool_selection",
  "action": "selected_tools_for_execution",
  "tools_selected": [
    "evaluate_classification_model"
  ],
  "tool_calls_detail": [
    {
      "name": "evaluate_classification_model",
      "args": {
        "bindings": {
          "model": "best_estimator",
          "X_test": "X_train",
          "y_test": "Y_train"
        },
        "func_kwargs": {
          "model_name": "LightGBM",
          "eval_data_label": "train"
        }
      },
      "id": "call_18GzPmWkNs6zNLuY9s0774Cm",
      "type": "tool_call"
    }
  ],
  "content": ""
},
{
  "step_number": 68,
  "message_type": "ToolMessage",
  "timestamp": "2025-11-20T23:52:58.042825",
  "step_type": "tool_result",
  "action": "received_tool_output",
  "content_preview": "Applied evaluate_classification_model with
docstring: \n    Evaluate a trained classification model on
data.\n\n    This function is used to evaluate the performance of a
trained classification model on a given dataset.\n\n    The dataset
can be the anything, such as the test set or the train set.\n\n
U...",
  "content_length": 1820
},
{
  "step_number": 69,
  "message_type": "HumanMessage",
  "timestamp": "2025-11-20T23:52:58.042829",
  "step_type": "reflection",
  "action": "llm_reflection",
  "content_preview": "Reasoning: The agent followed a standard and
appropriate pipeline for the Spaceship Titanic classification task
by loading and concatenating train and test datasets, imputing
missing values, casting categorical variables, and executing
meaningful feature engineering (extracting deck and side from Ca..."

```

```

1890     "content_length": 1303,
1891     "extracted_score": 7.0,
1892     "full_reflection_content": "Reasoning: The agent followed a
1893 standard and appropriate pipeline for the Spaceship Titanic
1894 classification task by loading and concatenating train and test
1895 datasets, imputing missing values, casting categorical variables,
1896 and executing meaningful feature engineering (extracting deck and
1897 side from Cabin, aggregating spending, and creating group features).
1898 They correctly split combined data back to train and test before
1899 modeling, which fixed an earlier critical workflow mistake. The
1900 agent performed hyperparameter tuning on a LightGBM classifier and
1901 achieved a reasonable CV score (~0.79) and training accuracy (~81%),
1902 indicating a decent model fit.\n\nHowever, despite the solid data
1903 preparation and modeling, the agent did not proceed to generate
1904 predictions on the test set or save the best model and predictions
1905 to the required directories, which is mandatory to fully solve the
1906 task. There were also minor syntax errors and API usage issues that
1907 caused some delays. Overall, the trajectory is methodical and
1908 technically sound in data handling, feature engineering, and
1909 modeling, but incomplete as it lacks the final step of output
1910 generation. The agent should now predict on the test set, save the
1911 submission file as specified, and also save the best model to
1912 conclude the task.\n\nScore: 7\nScore: 7"
1913 },
1914 {
1915     "step_number": 70,
1916     "message_type": "AIMessage",
1917     "timestamp": "2025-11-20T23:52:58.042851",
1918     "step_type": "tool_selection",
1919     "action": "selected_tools_for_execution",
1920     "tools_selected": [
1921         "predict_target"
1922     ],
1923     "tool_calls_detail": [
1924         {
1925             "name": "predict_target",
1926             "args": {
1927                 "bindings": {
1928                     "model": "best_estimator",
1929                     "X_data": "X_test"
1930                 },
1931                 "func_kwargs": {
1932                     "model_name": "LightGBM",
1933                     "return_probabilities": false
1934                 },
1935                 "output": "test_predictions"
1936             },
1937             "id": "call_9j7gLB7G1cH6WGnnRpygvWVb",
1938             "type": "tool_call"
1939         }
1940     ],
1941     "content": ""
1942 },
1943 {
1944     "step_number": 71,
1945     "message_type": "ToolMessage",
1946     "timestamp": "2025-11-20T23:52:58.042854",
1947     "step_type": "tool_result",
1948     "action": "received_tool_output",
1949     "content_preview": "Applied predict_target with docstring: \n
1950 Make predictions using a trained model.\n\n\nParameters\n
1951 -----\n\n    model : Any\n\n        Trained model (regression or
1952 classification)\n\n    X_data : pd.DataFrame\n\n        Features to make
1953 predictions on\n\n    model_name : str, default=\"model\"\n\n
1954 ...",

```

```

1944     "content_length": 1018
1945 },
1946 {
1947     "step_number": 72,
1948     "message_type": "HumanMessage",
1949     "timestamp": "2025-11-20T23:52:58.042858",
1950     "step_type": "reflection",
1951     "action": "llm_reflection",
1952     "content_preview": "Reasoning: The agent followed a comprehensive
1953 and well-structured approach to the Spaceship Titanic classification
1954 task. They started by loading and concatenating the train and test
1955 datasets to ensure consistent preprocessing. They successfully
1956 identified missing data and applied appropriate imputat...",
1957     "content_length": 1588,
1958     "extracted_score": 7.0,
1959     "full_reflection_content": "Reasoning: The agent followed a
1960 comprehensive and well-structured approach to the Spaceship Titanic
1961 classification task. They started by loading and concatenating the
1962 train and test datasets to ensure consistent preprocessing. They
1963 successfully identified missing data and applied appropriate
1964 imputation strategies for categorical and numerical features.
1965 Categorical columns were correctly cast, and meaningful feature
1966 engineering was performed, including extracting cabin deck and side,
1967 aggregating spending features, and creating group features from
1968 PassengerId, demonstrating good domain understanding. The agent then
1969 corrected a critical oversight by properly splitting the combined
1970 dataset back into train and test before modeling, which allowed
1971 successful hyperparameter tuning of a LightGBM model, yielding a
1972 reasonable cross-validation score (~0.79) and training accuracy
1973 (~81%). Evaluation results were clearly reported. Finally, the agent
1974 generated predictions on the test data. However, the agent did not
1975 proceed to save the best model or the predictions to the specified
1976 submission folder, which is essential to fully solve the task as per
1977 instructions. There were also some minor API usage errors and
1978 workflow missteps early on, but these were quickly corrected.
1979 Overall, the agent demonstrated good understanding and execution of
1980 data preparation, feature engineering, and modeling. The main
1981 shortcoming is not completing the final output saving step, leaving
1982 the task only partially solved. The trajectory is methodical, mostly
1983 sufficient, but incomplete. Score: 7\nScore: 7"
1984 },
1985 {
1986     "step_number": 73,
1987     "message_type": "AIMessage",
1988     "timestamp": "2025-11-20T23:52:58.042883",
1989     "step_type": "tool_selection",
1990     "action": "selected_tools_for_execution",
1991     "tools_selected": [
1992         "save_dataframe_to_csv"
1993     ],
1994     "tool_calls_detail": [
1995         {
1996             "name": "save_dataframe_to_csv",
1997             "args": {
1998                 "df": "test_predictions",
1999                 "filepath":
2000 "submissions/spaceship-titanic/lats-reflection/lats-reflection.csv"
2001             },
2002             "id": "call_BTgHyO2TJNUOWcyx4YTiNwIr",
2003             "type": "tool_call"
2004         }
2005     ],
2006     "content": ""
2007 },

```

I.4 MCTS-OUTCOME

I.4.1 FAILED TRAJECTORY

The most promising trajectory identified consisted of 18 tool calls but repeatedly encountered failures when attempting to load the data. Consequently, the entire search tree remained unproductive and failed to reach the modeling phase. This suggests that the system lacked sufficient guidance to advance beyond the initial stages of the pipeline.

```
{
  "step_number": 22,
  "message_type": "AIMessage",
  "timestamp": "2025-11-20T16:12:41.763597",
  "step_type": "tool_selection",
  "action": "selected_tools_for_execution",
  "tools_selected": [
    "read_data"
  ],
  "tool_calls_detail": [
    {
      "name": "read_data",
      "args": {
        "output": "train_df"
      },
      "id": "call_3did4EVkLGDyw9QRFoiAnU1A",
      "type": "tool_call"
    }
  ],
  "content": ""
},
{
  "step_number": 23,
  "message_type": "ToolMessage",
  "timestamp": "2025-11-20T16:12:41.763599",
  "step_type": "tool_result",
  "action": "received_tool_output",
  "content_preview": "Error: TypeError(\"read_data() missing 1
required positional argument: 'filepath'\")\n Please fix your
mistakes.",
  "content_length": 109
},
```

I.4.2 SUCCESSFUL TRAJECTORY

This trajectory contains 32 tool calls and fits only a CatBoost model. Notably, the submission file produced by this run exhibited a shape mismatch with the hidden test set, resulting in a percentile score of 0. The mismatch arose from errors in the feature engineering stage, which filtered out a portion of the test data.

```
{
  "step_number": 63,
  "message_type": "ToolMessage",
  "timestamp": "2025-11-20T16:13:57.197079",
  "step_type": "tool_result",
  "action": "received_tool_output",
  "content_preview": "Applied label_encode with docstring: \n
Label encode categorical columns.\n      \n      Parameters:\n
-----\n      df : pd.DataFrame\n      Input DataFrame\n
columns : str or List[str], optional\n      Column(s) to encode.\n
If None, encodes all object/category columns.\n      \n      Returns:\n
----...",
  "content_length": 701
},
{
```

```

2052     "step_number": 64,
2053     "message_type": "HumanMessage",
2054     "timestamp": "2025-11-20T16:13:57.197084",
2055     "step_type": "reward_feedback",
2056     "action": "generated_reward_feedback",
2057     "content_preview": "Modeling is still in progress",
2058     "content_length": 29
2059 },
2060 {
2061     "step_number": 65,
2062     "message_type": "AIMessage",
2063     "timestamp": "2025-11-20T16:13:57.197089",
2064     "step_type": "tool_selection",
2065     "action": "selected_tools_for_execution",
2066     "tools_selected": [
2067         "tune_catboost_classifier"
2068     ],
2069     "tool_calls_detail": [
2070         {
2071             "name": "tune_catboost_classifier",
2072             "args": {
2073                 "bindings": {
2074                     "X_train": "X_train",
2075                     "y_train": "Y_train"
2076                 },
2077                 "output": "best_catboost_model"
2078             },
2079             "id": "call_WbNtNSt3M3UXFDk5JBj49Nkm",
2080             "type": "tool_call"
2081         }
2082     ],
2083     "content": ""
2084 },
2085 {
2086     "step_number": 66,
2087     "message_type": "ToolMessage",
2088     "timestamp": "2025-11-20T16:13:57.197092",
2089     "step_type": "tool_result",
2090     "action": "received_tool_output",
2091     "content_preview": "Applied tune_catboost_classifier with
2092 docstring: \n    Perform hyperparameter tuning for CatBoost
2093 Classifier using GridSearchCV.\n\n    Parameters\n    -----\n
2094 X_train : pd.DataFrame\n    Training features\n    y_train :
2095 pd.Series\n    Training target\n    cv : int, default=5\n
2096 Number o...",
2097     "content_length": 1087
2098 },
2099 {
2100     "step_number": 67,
2101     "message_type": "HumanMessage",
2102     "timestamp": "2025-11-20T16:13:57.197097",
2103     "step_type": "reward_feedback",
2104     "action": "generated_reward_feedback",
2105     "content_preview": "Human Feedback: Verified that the modeling was
successful",
2106     "content_length": 57
2107 },
2108 {
2109     "step_number": 68,
2110     "message_type": "AIMessage",
2111     "timestamp": "2025-11-20T16:13:57.197102",
2112     "step_type": "tool_selection",
2113     "action": "selected_tools_for_execution",
2114     "tools_selected": [
2115         "evaluate_classification_model"

```

```

2106 ],
2107 "tool_calls_detail": [
2108   {
2109     "name": "evaluate_classification_model",
2110     "args": {
2111       "bindings": {
2112         "model": "best_estimator",
2113         "X_test": "X_train",
2114         "y_test": "Y_train"
2115       },
2116       "func_kwargs": {
2117         "model_name": "CatBoostClassifier",
2118         "eval_data_label": "train"
2119       }
2120     },
2121     "id": "call_hKlPSrF4ZrPiof5KJTUGpsOb",
2122     "type": "tool_call"
2123   }
2124 ],
2125 "content": ""
2126 },
2127 {
2128   "step_number": 69,
2129   "message_type": "ToolMessage",
2130   "timestamp": "2025-11-20T16:13:57.197105",
2131   "step_type": "tool_result",
2132   "action": "received_tool_output",
2133   "content_preview": "Applied evaluate_classification_model with
2134 docstring: \n    Evaluate a trained classification model on
2135 data.\n\n    This function is used to evaluate the performance of a
2136 trained classification model on a given dataset.\n\n    The dataset
2137 can be the anything, such as the test set or the train set.\n\n
2138 U...",
2139   "content_length": 1830
2140 },
2141 {
2142   "step_number": 70,
2143   "message_type": "HumanMessage",
2144   "timestamp": "2025-11-20T16:13:57.197112",
2145   "step_type": "reward_feedback",
2146   "action": "generated_reward_feedback",
2147   "content_preview": "Human Feedback: submission DataFrame was NOT
2148 created successfully. Please check the signature of the wrapped
2149 function if any, and call it with the correct arguments. The tool
2150 signature is This tool reads arguments from the scratchpad using
2151 'bindings', passes them to the internal function.\n\n    A...",
2152   "content_length": 1979
2153 },
2154 {

```

2149 I.5 MCTS-SHAPED

2151 I.5.1 FAILED TRAJECTORY

2152 The most promising trajectory consisted of 19 steps and terminated at the stage of splitting the
 2153 combined data into train and test. This indicates that the search never reached the modeling phase.

```

2154 {
2155   "step_number": 49,
2156   "message_type": "AIMessage",
2157   "timestamp": "2025-11-20T17:41:18.892989",
2158   "step_type": "tool_selection",
2159   "action": "selected_tools_for_execution",
2160   "tools_selected": [

```



```

2160         "drop_feature"
2161     ],
2162     "tool_calls_detail": [
2163         {
2164             "name": "drop_feature",
2165             "args": {
2166                 "bindings": {
2167                     "df": "combined_df"
2168                 },
2169                 "output": "combined_df",
2170                 "func_kwargs": {
2171                     "column": [
2172                         "PassengerId",
2173                         "Name",
2174                         "Cabin"
2175                     ]
2176                 }
2177             },
2178             "id": "call_W8BQW8V16Jch32pPwgSkJXIN",
2179             "type": "tool_call"
2180         }
2181     ],
2182     "content": ""
2183 },
2184 {
2185     "step_number": 50,
2186     "message_type": "ToolMessage",
2187     "timestamp": "2025-11-20T17:41:18.892992",
2188     "step_type": "tool_result",
2189     "action": "received_tool_output",
2190     "content_preview": "Applied drop_feature with docstring: \n
2191 Drop feature(s) from the DataFrame.\n      \n      Parameters:\n
2192 -----\n      df : pd.DataFrame\n      Input DataFrame\n
2193 column : str or List[str]\n      Name(s) of the column(s) to
2194 drop\n      \n      Returns:\n      -----\n      pd.DataFrame\n
2195 DataFrame wi...",
2196     "content_length": 662
2197 },
2198 {
2199     "step_number": 51,
2200     "message_type": "HumanMessage",
2201     "timestamp": "2025-11-20T17:41:18.892997",
2202     "step_type": "reward_feedback",
2203     "action": "generated_reward_feedback",
2204     "content_preview": "Categorical columns found: [{'column':
2205 'HomePlanet', 'unique_count': 3, 'dtype': 'category'}, {'column':
2206 'CryoSleep', 'unique_count': 3, 'dtype': 'category'}, {'column':
2207 'Destination', 'unique_count': 4, 'dtype': 'category'}]. Columns
2208 with dtypes that are not numeric/categorical found: []. Please
2209 co...",
2210     "content_length": 525
2211 },
2212 {
2213     "step_number": 52,
2214     "message_type": "AIMessage",
2215     "timestamp": "2025-11-20T17:41:18.893002",
2216     "step_type": "tool_selection",
2217     "action": "selected_tools_for_execution",
2218     "tools_selected": [
2219         "label_encode"
2220     ],
2221     "tool_calls_detail": [
2222         {
2223             "name": "label_encode",
2224             "args": {

```

```

2214         "bindings": {
2215             "df": "combined_df"
2216         }
2217     },
2218     "id": "call_LZuEYLrfsTWPCfWsh5Ungld1",
2219     "type": "tool_call"
2220 }
2221 ],
2222 "content": ""
2223 },
2224 {
2225     "step_number": 53,
2226     "message_type": "ToolMessage",
2227     "timestamp": "2025-11-20T17:41:18.893005",
2228     "step_type": "tool_result",
2229     "action": "received_tool_output",
2230     "content_preview": "Applied label_encode with docstring: \n
Label encode categorical columns.\n      \n      Parameters:\n
-----\n      df : pd.DataFrame\n      Input DataFrame\n
columns : str or List[str], optional\n      Column(s) to encode.\n
If None, encodes all object/category columns.\n      \n      Returns:\n
---...",
2231     "content_length": 709
2232 },
2233 {
2234     "step_number": 54,
2235     "message_type": "HumanMessage",
2236     "timestamp": "2025-11-20T17:41:18.893010",
2237     "step_type": "reward_feedback",
2238     "action": "generated_reward_feedback",
2239     "content_preview": "Human Feedback: Verified that the feature
engineering is complete and all columns are numeric (int/float)",
2240     "content_length": 105
2241 },
2242 {
2243     "step_number": 55,
2244     "message_type": "AIMessage",
2245     "timestamp": "2025-11-20T17:41:18.893015",
2246     "step_type": "tool_selection",
2247     "action": "selected_tools_for_execution",
2248     "tools_selected": [
2249         "split_combined_into_train_test"
2250     ],
2251     "tool_calls_detail": [
2252         {
2253             "name": "split_combined_into_train_test",
2254             "args": {
2255                 "bindings": {
2256                     "combined": "combined_df"
2257                 },
2258                 "output": "splitted_dfs"
2259             },
2260             "id": "call_LsWLDkHhkrEQ8Ja0rrBMQOXh",
2261             "type": "tool_call"
2262         }
2263     ],
2264     "content": ""
2265 },
2266 {
2267     "step_number": 56,
2268     "message_type": "ToolMessage",
2269     "timestamp": "2025-11-20T17:41:18.893018",
2270     "step_type": "tool_result",
2271     "action": "received_tool_output",

```

```

2268         "content_preview": "Applied split_combined_into_train_test with
2269         docstring: \n      Split combined data back into train and test using
2270         tracking columns.\n      \n      Parameters:\n      ----- \n
2271         combined : pd.DataFrame\n      Combined DataFrame with tracking
2272         columns\n      \n      Returns:\n      ----- \n      tuple\n
2273         (train_d...",
2274         "content_length": 607
2275     },
2276     {
2277         "step_number": 57,
2278         "message_type": "HumanMessage",
2279         "timestamp": "2025-11-20T17:41:18.893022",
2280         "step_type": "reward_feedback",
2281         "action": "generated_reward_feedback",
2282         "content_preview": "Human Feedback: Verified that the train and
2283         test data were split successfully",
2284         "content_length": 77
2285     }
2286 ],

```

I.5.2 SUCCESSFUL TRAJECTORY

The trajectory spans 23 steps and fits exclusively a LightGBM classifier in the modeling stage. When evaluated on our test split, it attains a leaderboard percentile of 33.11.

```

2290     {
2291         "step_number": 61,
2292         "message_type": "AIMessage",
2293         "timestamp": "2025-11-20T17:45:34.047940",
2294         "step_type": "tool_selection",
2295         "action": "selected_tools_for_execution",
2296         "tools_selected": [
2297             "tune_lightgbm_classifier"
2298         ],
2299         "tool_calls_detail": [
2300             {
2301                 "name": "tune_lightgbm_classifier",
2302                 "args": {
2303                     "bindings": {
2304                         "X_train": "X_train",
2305                         "y_train": "Y_train"
2306                     },
2307                     "output": "tuned_model_info"
2308                 },
2309                 "id": "call_LfolkEfImAXxyb5Ha3nNW4cp",
2310                 "type": "tool_call"
2311             }
2312         ],
2313         "content": ""
2314     },
2315     {
2316         "step_number": 62,
2317         "message_type": "ToolMessage",
2318         "timestamp": "2025-11-20T17:45:34.047944",
2319         "step_type": "tool_result",
2320         "action": "received_tool_output",
2321         "content_preview": "Applied tune_lightgbm_classifier with
2322         docstring: \n      Perform hyperparameter tuning for LightGBM
2323         Classifier using GridSearchCV.\n      \n      Parameters\n      ----- \n
2324         X_train : pd.DataFrame\n      Training features\n      y_train :
2325         pd.Series\n      Training target\n      cv : int, default=5\n
2326         Number o...",
2327         "content_length": 1092
2328     },

```

```

2322     {
2323         "step_number": 63,
2324         "message_type": "HumanMessage",
2325         "timestamp": "2025-11-20T17:45:34.047949",
2326         "step_type": "reward_feedback",
2327         "action": "generated_reward_feedback",
2328         "content_preview": "Human Feedback: Verified that the modeling was
2329 successful",
2329         "content_length": 57
2330     },
2331     {
2332         "step_number": 64,
2333         "message_type": "AIMessage",
2334         "timestamp": "2025-11-20T17:45:34.047955",
2335         "step_type": "tool_selection",
2336         "action": "selected_tools_for_execution",
2337         "tools_selected": [
2338             "evaluate_classification_model"
2339         ],
2340         "tool_calls_detail": [
2341             {
2342                 "name": "evaluate_classification_model",
2343                 "args": {
2344                     "bindings": {
2345                         "model": "best_estimator",
2346                         "X_test": "X_train",
2347                         "y_test": "Y_train"
2348                     },
2349                     "func_kwargs": {
2350                         "model_name": "LightGBM",
2351                         "eval_data_label": "train"
2352                     }
2353                 },
2354                 "id": "call_iLtqzfkOBwtWmROl9WDHKO75",
2355                 "type": "tool_call"
2356             }
2357         ],
2358         "content": ""
2359     },
2360     {
2361         "step_number": 65,
2362         "message_type": "ToolMessage",
2363         "timestamp": "2025-11-20T17:45:34.047959",
2364         "step_type": "tool_result",
2365         "action": "received_tool_output",
2366         "content_preview": "Applied evaluate_classification_model with
2367 docstring: \n    Evaluate a trained classification model on
2368 data.\n\n    This function is used to evaluate the performance of a
2369 trained classification model on a given dataset.\n\n    The dataset
2370 can be the anything, such as the test set or the train set.\n\n
2371 U...",
2372         "content_length": 1820
2373     },
2374     {
2375         "step_number": 66,
2376         "message_type": "HumanMessage",
2377         "timestamp": "2025-11-20T17:45:34.047965",
2378         "step_type": "reward_feedback",
2379         "action": "generated_reward_feedback",
2380         "content_preview": "Human Feedback: submission DataFrame was NOT
2381 created successfully. Please check the signature of the wrapped
2382 function if any, and call it with the correct arguments. The tool
2383 signature is This tool reads arguments from the scratchpad using
2384 `bindings`, passes them to the internal function.\n\n    A...",
2385         "content_length": 1979
2386     }

```

```

2376     },
2377     {
2378         "step_number": 67,
2379         "message_type": "AIMessage",
2380         "timestamp": "2025-11-20T17:45:34.047972",
2381         "step_type": "tool_selection",
2382         "action": "selected_tools_for_execution",
2383         "tools_selected": [
2384             "predict_target"
2385         ],
2386         "tool_calls_detail": [
2387             {
2388                 "name": "predict_target",
2389                 "args": {
2390                     "bindings": {
2391                         "model": "best_estimator",
2392                         "X_data": "X_test"
2393                     },
2394                     "output": "test_predictions",
2395                     "func_kwargs": {
2396                         "model_name": "LightGBM"
2397                     }
2398                 },
2399                 "id": "call_lZQQw4OoT0LANXYCUvjXh1Bp",
2400                 "type": "tool_call"
2401             }
2402         ],
2403         "content": ""
2404     },
2405     {
2406         "step_number": 68,
2407         "message_type": "ToolMessage",
2408         "timestamp": "2025-11-20T17:45:34.047975",
2409         "step_type": "tool_result",
2410         "action": "received_tool_output",
2411         "content_preview": "Applied predict_target with docstring: \n
2412         Make predictions using a trained model.\n      \n      Parameters\n
2413         -----\n      model : Any\n      Trained model (regression or
2414         classification)\n      X_data : pd.DataFrame\n      Features to make
2415         predictions on\n      model_name : str, default=\"model\"\n
2416         ...",
2417         "content_length": 1018
2418     },
2419     {
2420         "step_number": 69,
2421         "message_type": "HumanMessage",
2422         "timestamp": "2025-11-20T17:45:34.047981",
2423         "step_type": "reward_feedback",
2424         "action": "generated_reward_feedback",
2425         "content_preview": "Human Feedback: Verified that the submission
2426         DataFrame was created successfully with 1 columns (boolean) and no
2427         missing values.",
2428         "content_length": 127
2429     }
2430 ],

```

I.6 HIERARCHICAL MCTS

I.6.1 FAILED TRAJECTORY

This trajectory consists of 55 tool calls and fails at the train data to features and target stage, prior to modeling. Nothing of note in this trajectory, and hence is not reported.

I.6.2 SUCCESSFUL TRAJECTORY

The trajectory spans 80 steps and evaluates multiple modeling approaches, attaining a leaderboard percentile of 62.65. Hierarchical MCTS demonstrates the ability to execute substantially longer trajectories, owing to its use of subtask decomposition and task masking.

```
{
  "step_type": "tool_selection",
  "action": "selected_tools_for_execution",
  "tools_selected": [
    "tune_lightgbm_classifier"
  ],
  "tool_calls_detail": [
    {
      "name": "tune_lightgbm_classifier",
      "args": {
        "bindings": {
          "X_train": "X_train",
          "y_train": "Y_train"
        },
        "output": "tuned_lgbm_classifier"
      },
      "id": "call_fM1ezojD8q9aj9f052GtppXm",
      "type": "tool_call"
    }
  ],
  "content": ""
},
{
  "step_type": "tool_result",
  "action": "received_tool_output",
  "content_preview": "Applied tune_lightgbm_classifier with
docstring: \n    Perform hyperparameter tuning for LightGBM
Classifier using GridSearchCV.\n\n    Parameters\n    -----\n    X_train : pd.DataFrame\n    Training features\n    y_train : pd.Series\n    Training target\n    cv : int, default=5\n    Number o...",
  "content_length": 1092
},
{
  "step_type": "reward_feedback",
  "action": "generated_reward_feedback",
  "content_preview": "Human Feedback: Verified that the modeling
was successful",
  "content_length": 57
},
{
  "step_type": "tool_selection",
  "action": "selected_tools_for_execution",
  "tools_selected": [
    "tune_random_forest_classifier"
  ],
  "tool_calls_detail": [
    {
      "name": "tune_random_forest_classifier",
      "args": {
        "bindings": {
          "X_train": "X_train",
          "y_train": "Y_train"
        },
        "func_kwargs": {
          "cv": 5
        },
        "output": "rf_tuned_model_output"
      }
    }
  ],
}
```

```

    "id": "call_30sXBYeYfbCnehW6uvlfA99X",
    "type": "tool_call"
  }
],
  "content": "I cannot evaluate the model on the test data
because the true target labels are not available for the test set.
Instead, I will generate predictions for the test data using the
best Random Forest classifier model and save the predictions to the
required submission file path."
},
{
  "step_type": "tool_result",
  "action": "received_tool_output",
  "content_preview": "Applied tune_random_forest_classifier with
docstring: \n    Perform hyperparameter tuning for Random Forest
Classifier using GridSearchCV.\n\n    Parameters\n    -----\n    X_train : pd.DataFrame\n    Training features\n    y_train :
pd.Series\n    Training target\n    cv : int, default=5\n
...",
  "content_length": 1134
},
{
  "step_type": "reward_feedback",
  "action": "generated_reward_feedback",
  "content_preview": "Human Feedback: Verified that the modeling
was successful",
  "content_length": 57
}

```