

SMARTDS-SOLVER: AGENTIC AI FOR VERTICAL DOMAIN PROBLEM SOLVING IN DATA SCIENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Automating complex, multi-step vertical domain tasks—such as Data Science (DS) workflows—presents significant challenges for large language model (LLM) agents. Existing AutoDS approaches often rely on prompt-sensitive, fragmented multi-turn interactions and costly full re-generation upon execution failure, leading to unstable workflow coherence and high token consumption. We introduce **SmartDS-Solver**, a reasoning-centric agentic system designed to enhance the stability, robustness, and cost efficiency of these workflows. Our core approach integrates rigorous workflow planning into a domain-specialized **Reasoning LLM**, which is trained using **structured methodological distillation** and a **two-stage Group Relative Policy Optimization (GRPO)** procedure. Crucially, SmartDS-Solver employs a lightweight agentic layer featuring the novel **State-Aware Refinement and Temperature Exploration (SARTE)** algorithm. SARTE dynamically adjusts the LLM’s decoding strategy based on deterministic execution feedback, enabling **minimally invasive patching** rather than costly full re-planning. We performed a comprehensive evaluation across **32 datasets** covering 11 MLE-Bench tasks, 18 AutoML-Agent benchmarks, and 3 real-world tasks, showing consistent gains while reducing inference and modification token usage. In the MLE-Bench benchmark, our 32B model attains an **81.8% win rate** over the AIDE+o1-preview baseline, and on the 18 AutoML-Agent tasks, the win rate reaches **94%**. Notably, even a **7B model** produces fully executable solutions on all evaluated tasks, demonstrating the scalability and robustness of our method. SmartDS-Solver reduces token usage by approximately **78%** on the 11 MLE-Bench tasks. The SARTE meta-control mechanism significantly boosts decoding performance—raising average accuracy by **3.9%**, lowering error rates by **12%**, and delivering an overall **75% significant improvement** on MLE-Bench tasks ($p = 0.0173$).¹

1 INTRODUCTION

Large language models (LLMs) have recently demonstrated strong potential for automating end-to-end data science workflows. Systems such as AutoML-Agent+4o (Trirat et al., 2025) and AIDE+o1-preview (Jiang et al., 2025) show that LLMs can perform dataset analysis, pipeline planning, model training, and code generation with minimal human supervision.

However, existing LLM-based AutoDS (Automated Data Science) systems still face several challenges that hinder practical deployment. AutoML-Agent adopts a multi-agent design where agents communicate to construct general-purpose AutoML pipelines. While this provides greater generality than traditional AutoML tools such as AutoGluon (Erickson et al., 2020), the architecture introduces substantial deployment overhead, and fragile inter-agent coordination often leads to pipeline instability and task-silo behavior. In contrast, reasoning-centric systems such as AIDE relies heavily on repeated frontier-model (o1-preview) calls for tree-search-based refinement. AIDE enforces a fixed search depth: after generating an initial candidate solution, subsequent steps stochastically choose between bug fixing and incremental refinement. Each step triggers two LLM calls—one for generation and one for summarization—causing computation to grow linearly with the preset depth. Although AIDE achieves strong performance, including a 16.9% bronze medal rate on the 75-task MLE-Bench benchmark (Chan et al., 2024), its heavy reliance on repeated frontier-model calls leads

¹The source code will be made publicly available upon acceptance of the paper.

054 to significant computational and monetary cost (Liu, 2024; Rodriguez, 2024), limiting its real-world
055 applicability.

056 Beyond these architectural issues, current LLM-driven data science agents still exhibit three practical
057 limitations: (i) **fragile task coherence**, where multi-step workflows easily break due to prompt
058 brittleness or locally inconsistent decisions; (ii) **high computational and monetary cost**, especially
059 for systems like AIDE+o1-preview that repeatedly invoke expensive frontier models; (iii) **task silos**
060 **and weak cross-task generalization**, where methodological knowledge is not systematically reused
061 across datasets, modalities, or objectives.

062 These limitations raise a central question: *Can we design an agentic system that, by fine-tuning*
063 *an open-source LLM, internalizes the data-science workflow into a private, domain-specialized rea-*
064 *soning model, thereby reducing dependence on frontier models and lowering operational costs?*
065 *Furthermore, during subsequent interactions with external LLMs, can the system automatically and*
066 *dynamically regulate its behavior to ensure stable agentic execution, ultimately enabling both high*
067 *performance and high efficiency?*

068 Together, these components allow SmartDS-Solver to provide strong problem-solving performance
069 while reducing inference and modification token usage by more than an order of magnitude com-
070 pared to frontier-model-based baselines.

071 **Our contributions are summarized as follows:**

- 072 • We propose **SmartDS-Solver**, a hierarchical agentic AI architecture that explicitly separates
073 domain-specific reasoning, code generation, code refinement, and strategy control, thereby main-
074 taining interaction stability and effectiveness throughout the end-to-end data science workflow.
- 075 • We develop a **domain-finetuned Data Science Reasoning LLM** that internalizes the decision-
076 making process of data science workflows through structured methodological distillation and a
077 two-stage GRPO training scheme. The resulting model can generate stable, executable, and
078 pipeline-aware code with minimal prompt engineering.
- 079 • We introduce **SARTE**, a situation-aware Meta-Learning Agent that dynamically adjusts the de-
080 coding temperature online. SARTE improves success rate and token efficiency without requiring
081 additional model training.
- 082 • SmartDS-Solver delivers 81.8% win rate over AIDE+o1-preview, 94% over AutoML-Agent, and
083 cuts inference token usage by 78% across 32 diverse tasks.

084 We next provide background and related work (§2), present the SmartDS-Solver architecture and
085 training methodology (§3), and evaluate its performance, ablations, and cost characteristics (§4).

087 2 RELATED WORK

089 2.1 TRADITIONAL AUTOML SYSTEMS

091 Automated Machine Learning (AutoML) has emerged as an important research direction, targeting
092 automation of feature engineering, algorithm selection, and hyperparameter tuning. Early systems
093 such as Auto-WEKA (Thornton et al., 2013), Auto-Sklearn (Feurer et al., 2015), and TPOT (Le
094 et al., 2020) introduced combined pipeline search and hyperparameter optimization for structured
095 data modeling using Bayesian optimization. However, these systems assume well-structured prob-
096 lems with clear objectives and static data pipelines, while real-world data science often involves
097 evolving requirements and heterogeneous data sources. To address these limitations, Auto-sklearn
098 2.0 (Feurer et al., 2022) introduced meta-feature-free meta-learning and bandit-based budget al-
099 location to achieve stable performance under time constraints, pioneering "hands-free" AutoML.
100 FLAML (Wang et al., 2021) emphasizes lightweight search strategies, particularly suited for time
101 and resource-constrained application scenarios. The field also expanded into deep learning through
102 Neural Architecture Search (NAS), where methods like NAS-RL (Zoph & Le, 2016) successfully
103 identified efficient architectures for image and language tasks. Advanced hyperparameter opti-
104 mization methods like BOHB (Falkner et al., 2018) further enhanced the scalability of architecture
105 search. Despite these advances, fundamental challenges remain. Current AutoML systems strug-
106 gle with three critical areas: interpreting ambiguous task definitions, handling unstructured data
107 effectively, and facilitating meaningful human-AI collaboration. Additionally, approaches like NAS
suffer from high computational costs and tight coupling with predefined architectural components
that hinder broader applicability. Recent efforts such as Auto-CASH (Mu et al., 2022), ML2DAC

(Treder-Tschechlov et al., 2023), and AlphaD3M (Drori et al., 2021) have attempted to enhance cross-task generalization through experience-driven optimization, achieving preliminary progress in multimodal and interpretability aspects.

2.2 THE INTERSECTION OF PROGRAM SYNTHESIS, AUTONOMOUS AGENTS, AND DATA SCIENCE AUTOML

Recent work demonstrates that large language models (LLMs) have evolved from code completion and snippet generation to end-to-end autonomous programming capabilities. AlphaCode (Li et al., 2022) achieved an average ranking in the top 54.3% of programming competitions through generating large sets of program samples, filtering them based on execution results, and clustering the remaining samples. Commercial systems like Devin demonstrate LLMs’ capability as autonomous software engineers that can iteratively edit files, execute tests, and self-repair through continuous feedback loops. However, these systems focus primarily on general programming tasks rather than comprehensive data science workflows.

To address single-agent reasoning limitations, multi-agent frameworks such as AutoML-Agent+4o (Trirat et al., 2025) and AutoGPT (Significant Gravitas) leverage goal decomposition and collaborative dialogue paradigms, combining LLMs, external tools, or human feedback to enhance planning and execution capabilities. At the tool integration level, ReAct (Yao et al., 2023) interweaves Chain-of-Thought reasoning with action sequences, while Toolformer (Schick et al., 2023) enables models to learn “when, how, and why” to call external tools through self-supervision, both demonstrating the criticality of “reasoning-action coupling” for agentic AI success. However, data science workflows require specialized tool ecosystems—such as feature stores, GPU deployment, and experiment management—that current research has not yet addressed. Meanwhile, LLM-driven data science agents are attempting to break free from traditional “static pipeline search” paradigms. AutoML-GPT and AIDE (Zhang et al., 2023; Jiang et al., 2025) can automatically generate data processing, model architecture, and hyperparameter tuning code through LLM-based optimization approaches, with AIDE employing tree search and reward mechanisms. Agent K v1.0 (Grosnit et al., 2024) incorporates long-term memory modules for cross-task knowledge transfer, achieving Grandmaster-level performance in Kaggle competitions. These works confirm that LLM agents with planning, reflection, and tool-calling capabilities can demonstrate high accuracy and rapid convergence in open-ended data science tasks.

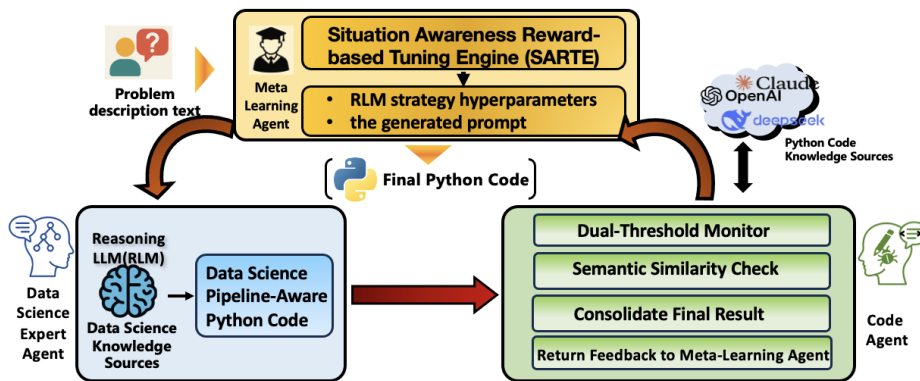


Figure 1: SmartDS-Solver: Hierarchical Multi-Agent Framework Overview

3 METHODOLOGY

3.1 OVERALL ARCHITECTURE

Figure 1 illustrates the SmartDS-Solver architecture, which is composed of three core agents. To address the above bottlenecks, SmartDS-Solver adopts a hierarchical collaborative multi-agent architecture consisting of three agents: a data science domain Expert Agent, a Meta-Learning Agent, and a Code Agent. Unlike sequential multi-agent workflows, our Meta-Learning Agent (SARTE) formulates reasoning as a strategy search problem and dynamically adjusts decoding hyperparam-

eters (e.g., temperature) using execution feedback, balancing exploration and exploitation at each step.

3.2 META-LEARNING AGENT AND SITUATION-AWARE CONTROL

3.2.1 META-LEARNING AGENT

The Meta-Learning Agent serves as the core decision module, treating agent interactions in data science tasks as a parameter optimization process. Using temperature as an example, experiments show that optimal ranges vary by task—e.g., NLP (0.8), vision (0.7), tabular (0.6), and time series (0.75)—and also shift after SFT or reinforcement learning. By dynamically adjusting temperature based on past settings and execution feedback, and providing corresponding prompts, the Meta-Learning Agent enhances the Expert Agent’s performance. Over time, it learns task-specific optimal ranges, improving both solution quality and resource efficiency through multi-turn interaction.

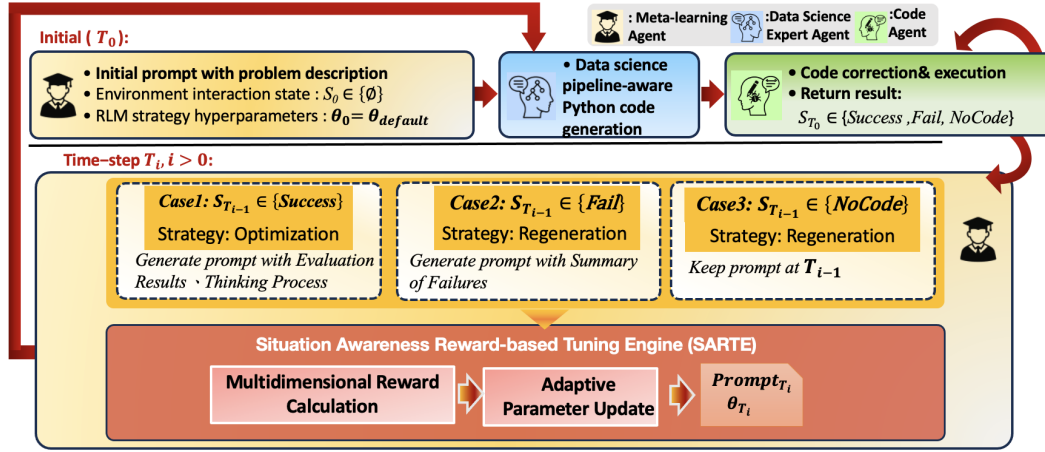


Figure 2: Multi-Turn Meta-Learning Workflow with SARTE for Adaptive Code Generation

3.2.2 SITUATION-AWARE META-CONTROL FRAMEWORK

Our Meta-Learning Agent acts as an adaptive controller for iterative problem solving in vertical domains (Fig. 2). Its core, the **Situation-Aware Reward-Tuning Engine (SARTE)**, integrates (1) execution feedback from the Code Agent—covering code executability, accuracy, and error messages, (2) hyperparameter settings, especially temperature, and (3) prompts and reasoning traces from the Expert Agent. Through SARTE, the Meta-Learning Agent serves as a global strategy regulator, analyzing the Expert Agent’s iterative reasoning and environment interaction. It then computes the suitable temperature and integrates prompt adjustments, guiding the Expert Agent’s next round of code generation. This design balances exploration and exploitation, reduces instability from full code regeneration, and improves success rates, iteration efficiency, and token usage. The SARTE algorithm proceeds as follows:

i. Multidimensional Reward Calculation We formalize the agent-environment interaction as a discrete time series $T_i, i \in \{1, 2, \dots, n\}$, where each interaction produces one of three mutually exclusive outcomes: Success, Failure, or NoCode, representing correct execution, execution failure, and absence of generated code, respectively. This tripartite classification enables precise behavioral shaping through our reward decomposition strategy. At each iteration, the outcome-specific reward is computed and aggregated to characterize dynamic behavior. Specifically, our reward function employs a time-weighted linear combination across three orthogonal dimensions:

$$R_n = \alpha(n) \cdot r_g + \beta(n) \cdot r_q + \gamma \cdot r_f \quad (1)$$

equation 1 embodies a dynamic priority shifting paradigm, where weights $\alpha(n)$ and $\beta(n)$ decay exponentially with the number of attempts n , gradually emphasizing efficiency over quality. This decay schedule prioritizes exploration early in the search and efficiency in later steps.

The weight scheduling follows the principle of diminishing marginal returns:

$$\alpha(n) = \max(0.5 - 0.1(n - 1), 0.2); \beta(n) = \max(0.4 - 0.05(n - 1), 0.2); \gamma = 0.1 \quad (2)$$

1.Efficiency Premium r_f rewards early success based on time value principles: $r_f = 10(n - i + 1)$ where n is the maximum number of attempts and i is the current attempt until success. The efficiency reward exponentially favors early success within the 10-round budget, giving the highest score to first-round solutions and the lowest to last-round ones.

2.Generation Quality Assessment r_g evaluates output format compliance: higher scores for proper structured formatting, partial scores for incomplete structure, and negative penalties when required tags are absent.

3.Execution Quality Evaluation r_q considers both execution results and generation history. Successful executions receive positive scores that decrease with correction attempts, failed executions incur severity-weighted negative penalties, while absent code generation (\emptyset) receives maximum penalty.

This reward design provides clear behavioral gradients and ensures convergence: early accurate solutions gain amplified rewards, while errors incur proportional penalties, guiding the model toward efficient and stable generation.”

ii. Adaptive Parameter Update Strategy To support dynamic adjustment of decoding hyperparameters (e.g., temperature), we adopt an adaptive update strategy that combines execution feedback from the Code Agent with boundary-aware modulation. A Control Factor determines the direction and magnitude of each update (see pseudocode in Appendix A.1).

EXECUTION STATES AND RULES.

- **Success:** increase temperature; a piecewise nonlinear amplification adjusts the step size based on distance from the target accuracy.
- **No-Code:** increase temperature; a progressive weight based on stagnation count determines the adjustment strength.
- **Failure:** decrease temperature; errors are mapped to a seven-level severity scale to determine penalty magnitude.

BOUNDARY-AWARE CONTROL. The current parameter is normalized and updated with a direction-sensitive linear decay: step size decreases by up to 70% near boundaries, full step size is preserved in central regions, and a minimum threshold (0.3) prevents stagnation. Final values are clipped to remain within valid bounds.

This constant-time update mechanism ($O(1)$ per iteration) jointly balances fast exploration, boundary safety, and stable convergence, enabling the model to adaptively manage exploration–exploitation trade-offs over multiple interactions and improving both success rate and computational efficiency.

3.3 DOMAIN-SPECIFIC REASONING MODEL AND EXPERT AGENT

3.3.1 DATA SCIENCE EXPERT AGENT

The Expert Agent acts as the domain specialist, converting ambiguous user requirements into complete data-science pipelines—covering preprocessing, feature engineering, model selection, and evaluation. Beyond outlining these workflows, it must also implement them by generating end-to-end Python code for diverse tasks. To explicitly capture such methodological capabilities, we construct domain-aware datasets and training procedures combined with GRPO-based reward design (details in later sections). This produces a domain-specialized reasoning model that serves as the Expert Agent’s “core brain,” capable of generating pipeline-aware, highly complete code with minimal post-hoc modification. Its generative behavior is further shaped by decoding parameters (e.g., temperature) provided by the Meta-Learning Agent, allowing generation tendencies to adapt to task requirements and improving overall flexibility.

3.3.2 SOURCE DATA AND DATA AUGMENTATION

We construct a multi-layered training dataset for our domain-specific reasoning model from multiple sources: (1) 3000 OpenThoughts-114 Python reasoning samples (Guha et al., 2025), (2) Kaggle

270 entries and (3) winning solutions from Hugging Face (Bigcomputer, 2024; Data-Agents, 2024), (4)
 271 cell2doc (Mondal et al., 2024) and (5) Code4ML (Drozdova et al., 2023)(see Appendix C.1 for de-
 272 tails). For sources (2)–(5), we embed structured templates of data science methodologies into our
 273 designed prompts, enabling DeepSeek R1 (Shao et al., 2024) to generate outputs with `<problem>`
 274 , `<deepseek_reasoning>`, and `<deepseek_solution>` segments. The `<deepseek_solution>` must re-
 275 tain original code structure with only minimal annotations to avoid execution drift.

276 Generated samples undergo formatting checks and are scored using Gemma3-27B based on reason-
 277 ing completeness and solution quality. We then deduplicate samples by computing semantic
 278 similarity and retaining only the highest-scoring item within each group using an 80% threshold.
 279 To ensure that the `<deepseek_solution>` systematically reflects real-world workflows and expert
 280 thinking, we explicitly require three core elements—full workflow narrative, decision logic, and an
 281 adjustment trail—and the complete prompt templates used to enforce these constraints are provided
 282 in Appendix C.2. These elements preserve alignment between reasoning and executable code and
 283 follow the agentic “planning → reflection → solution” paradigm.

284 After filtering and deduplication, the dataset is divided into three groups: general reasoning data,
 285 medal-level Kaggle solutions, and non-medal submissions, balancing high-quality exemplars with
 286 diverse real-world patterns to enhance robustness and generalization

287 3.3.3 TRAINING METHODOLOGY

288 Our training approach consists of three phases:

289
 290 **i. Capability Distillation via Supervised Finetuning (SFT)** The training data combines Group
 291 1 (general reasoning data) with a subset of Group 3 (non-medal Kaggle data), integrating broad
 292 reasoning capabilities(CoT-enhanced SFT prompt template Appendix A.2) with domain-specific
 293 knowledge to establish foundational data science problem-solving abilities.

294
 295 **ii. Reinforcement Learning with GRPO.** Building on the SFT baseline, we apply the Group
 296 Relative Policy Optimization (GRPO)(Shao et al., 2024) framework to align the model’s generation
 297 policy $\pi(y|x)$ with task-specific rewards. GRPO optimizes

$$298 \max_{\pi} \mathbb{E}_{x \sim D, y \sim \pi(\cdot|x)} [R(x, y)] \quad (3)$$

299 where x is the `<problem>` segment and y is the `<deepseek solution>` segment defined in §3.4.2,
 300 and $R(x, y)$ encourages high-quality, interpretable, and technically accurate outputs. We design a
 301 composite reward

$$302 R(x, y) = \alpha R_{\text{main}}(x, y) + \beta R_{\text{format}}(x, y) + \gamma R_{\text{reasoning}}(x, y), \quad (4)$$

303 with $\alpha = 0.5$ emphasizing technical correctness, $\beta = 0.25$ enforcing format compliance, and
 304 $\gamma = 0.25$ capturing logical structure.

305
 306 **Main task reward.** R_{main} measures semantic alignment between the generated code y and the
 307 reference y^* by extracting feature-engineering steps, algorithm specifications, and evaluation met-
 308 rics using a regex rule library, followed by verifying algorithm usage with Python AST analysis.
 309 These three dimensions produce normalized similarities ($S_{\text{feature}}, S_{\text{algorithm}}, S_{\text{metric}}$). To reflect their
 310 relative importance, we assign weights (0.4), (0.4), and (0.2) to the feature, algorithm, and metric
 311 components, respectively, and aggregate them as

$$312 R_{\text{main}} = \text{Aggregate}(S_{\text{feature}}, S_{\text{algorithm}}, S_{\text{metric}}). \quad (5)$$

313
 314 **Format reward.** R_{format} performs a binary check on the XML-like output structure using `<think`
 315 `> ...</think>` and `<answer> ...</answer>` tags.

316
 317 **Reasoning reward.** $R_{\text{reasoning}}$ scores the density of language-specific reasoning indicators and step
 318 markers, assigning higher weight to textit<think > blocks to promote clear step-by-step reasoning.
 319 Complete extraction rules and scoring formulas are listed in Appendix TablesA3 and A4.

320
 321 **iii. Two-stage GRPO for problem-solving enhancement and solution refinement.** The two
 322 GRPO phases use different training data to target distinct objectives:

- 323 • **Problem-solving enhancement (GRPO1):** trained on Group 3 non-medal Kaggle submissions
 (excluding those used for SFT).

Dataset Group	SmartDS-Solver (Ours)	AIDE	AutoML-Agent
11 MLE-Bench tasks	✓	✓	✓
18 AutoML-Agent tasks	✓	✓	✓
3 real-world tasks	✓	✓	—

Table 1: Coverage of each method across the three benchmark groups. SmartDS-Solver is evaluated on all 32 datasets, whereas AIDE and AutoML-Agent cover only subsets of the tasks.

- **Solution refinement (GRPO2):** trained on Group 2 gold, silver, and bronze Kaggle solutions.

To extend this methodology to other vertical domains, the fixed three-part reward structure—task correctness, format, and reasoning—can remain unchanged. Only the task-correctness component (R_{main}) needs to be adapted using domain-appropriate semantic matching or analysis, while the remaining reward terms and the overall training pipeline can be reused directly.

3.4 CODE AGENT: INTERACTIVE CODE REFINEMENT MANAGER

The Code Agent centers on an iterative refinement module (see Appendix Figure A1) guided by four design principles. **(1) Interactive self-repair loop:** the agent executes code, collects feedback, and augments prompts, injecting runtime signals into subsequent LLM calls. **(2) Dual-threshold monitoring with semantic early stopping:** timeouts and failure counts prevent redundant retries, while semantic similarity checks avoid ineffective correction attempts. **(3) Minimally invasive patching:** errors are fixed without altering the Expert Agent’s underlying reasoning. **(4) Hierarchical feedback and coordination:** structured execution results—no-code states, evaluations, or failure traces—are returned to the Meta-Learning Agent for the next-round strategy update. This design yields several benefits: (1) *Reliability and efficiency*—dual-threshold monitoring and error consolidation reduce deadlocks and unnecessary LLM calls; (2) *Rapid convergence*—targeted patches and focused summaries shorten iteration cycles; (3) *Enhanced observability*—structured status reporting provides clearer signals for meta-level adaptation; (4) *Scalability*—the modular components (extraction, execution, similarity analysis, consolidation) support easy extension to new domains.

4 EXPERIMENTAL RESULTS AND DISCUSSION

We evaluate SmartDS-Solver across three benchmark groups covering 32 datasets; detailed descriptions are provided in the Experimental Setup section. Our analysis focuses on: (i) overall task-solving performance, (ii) inference and code-modification cost efficiency, (iii) the effectiveness of the proposed finetuning pipeline (SFT + GRPO), and (iv) the impact of the SARTE meta-learning mechanism. Detailed task-level results are provided in the Appendix.

4.1 EXPERIMENTAL SETUP

Datasets. We evaluate on 32 datasets spanning three groups: (1) **MLE-Bench Competition** (11 Kaggle competition tasks across tabular, text, vision, and time-series); (2) **AutoML-Agent Benchmark** (18 public tasks from Kaggle Datasets, Kaggle Competitions, OpenML, Planetoid, and UCI ML); (3) **Real-world Tasks** (two ongoing Kaggle competitions and one proprietary medical dataset).

Baselines. We compare against two representative AutoDS systems: AIDE+o1-preview (reasoning-centric) and AutoML-Agent (multi-agent AutoML).

4.2 TASK-SOLVING PERFORMANCE AND GENERALITY ACROSS 32 TASKS

Performance on 11 MLE-Bench competition tasks. To validate the generality of our approach, SmartDS-Solver is instantiated with multiple modern LLM backbones, including Qwen-32B, Qwen-7B, and Llama-8B. We compare these variants against the AutoML-Agent baseline, which is built on GPT-4o. All systems are evaluated under the unified AIDE evaluation protocol and its accompanying dataset, ensuring a controlled and fair comparison across methods.

PERFORMANCE ACROSS DIFFERENT MODEL SCALES. Table 2 reports the performance of SmartDS-Solver under different model sizes. With the full GRPO2 finetuning pipeline, SmartDS-Solver paired with Qwen-32B achieves an 81.8% win rate. In the smaller-model regime, Llama-8B

reaches a 54.5% win rate, and generates executable code for all tasks. Overall, Llama-8B performance exceeds that of the Qwen-7B variants.

COMPARISON TO AUTOML-AGENT. Using the official implementation, we evaluate AutoML-Agent on the 11 MLE-Bench tasks. It achieves only a 9% win rate and fails to generate executable code on 3 tasks.

Model	SmartDS-Solver (ours)			AutoML-Agent
	Qwen-32B	Qwen-7B	Llama-8B	GPT-4o
Win Rate vs AIDE	81.8%	54.5%	54.5%	9%
#Fail	0	1	0	3
#Rank-1	3	0	2	0

Table 2: Summary of performance across 11 MLE-Bench tasks.

Performance on 18 AutoML-Agent benchmark tasks. To assess cross-domain generality, we additionally evaluate SmartDS-Solver on 18 datasets used in the AutoML-Agent benchmark, spanning Kaggle Datasets, Kaggle Competitions, OpenML, Planetoid, and UCI ML. Table 3 summarizes the averaged improvement and win-rate results. SmartDS-Solver achieves a **94% win rate** over AutoML-Agent, with consistent advantages observed across all dataset groups. AIDE achieves only a **28% win rate** against AutoML-Agent. Among the failures, six tasks did not execute, four produced no valid outputs, and two failed to follow the required evaluation metrics.

Performance on 3 real-world tasks. Table 4 (left) reports results on two ongoing Kaggle competitions and one proprietary medical modeling task. Beyond the strong performance demonstrated earlier with Qwen-32B, we also observe that substantially smaller models—such as Qwen-7B and Llama-8B—are able to reliably produce executable code and achieve competitive results across all three real-world tasks. In contrast, AIDE fails to generate executable code for one of the tasks.

4.3 COST-EFFICIENCY EVALUATION

AIDE Configuration. We follow AIDE’s default search depth of 20 (`agent.steps = 20`). Each search step requires one generation call and one review call; thus, the code-generation model (reasoning LLM: o1-preview) and the feedback model (general LLM: GPT-4o) are each invoked 20 times. Consequently, AIDE performs exactly **40 LLM calls per task**, independent of task difficulty.

SmartDS-Solver Configuration. SmartDS-Solver adopts a different resource allocation strategy. We set the maximum number of inference rounds to 10 (`InferenceMax = 10`), meaning the reasoning LLM specialized for data science is invoked at most ten times. Within each inference round, the system allows the Code Expert Agent to apply up to five code-modification attempts (using a general LLM: DeepSeek), equipped with an early-stopping mechanism to avoid unnecessary edits. Under this configuration, SmartDS-Solver never exceeds **50 total LLM calls per task**, though the actual number varies across tasks.

Observed Resource Usage. As shown in Table 4 (right), SmartDS-Solver is substantially more efficient on the 11 MLE-Bench tasks: (1) it uses on average **26** reasoning-LLM calls vs. AIDE’s fixed **40**; and (2) overall token usage is reduced by **78%**. AutoML-Agent relies on a general-purpose LLM (e.g., GPT-4o), with LLM calls dispersed across three major stages and several substages, making fine-grained comparison with AIDE and SmartDS-Solver infeasible. In our controlled evaluation on the 11 MLE-Bench tasks, AutoML-Agent issues approximately 1k LLM calls per task, compared with 440 for AIDE and 286 for SmartDS-Solver.

4.4 ABLATION STUDY: INSTRUCT → SFT → GRPO

Figure 3(a) reports an ablation study isolating the effect of each finetuning stage while keeping SARTE and the Code Agent fixed. Starting from the domain-agnostic Instruct models, Qwen-32B already achieves a 63.6% win rate over AIDE+o1-preview with only one failure, despite the higher inference cost of o1-preview relative to 70B-class open-source models. For Qwen-32B, SFT mainly improves stability rather than win rate: it removes all failing cases and yields one best-performing

Win Rate(%)	Kaggle Comp.	Kaggle Data.	OpenML	Planetoid	UCI ML	Overall
SmartDS-Solver > AutoML-Agent	100%	86%	100%	100%	100%	94%
AIDE > AutoML-Agent	25%	28%	50%	0%	0%	28%

Table 3: Comparison of SmartDS-Solver and AIDE against AutoML-Agent on 18 AutoML-Agent tasks.

Tasks	SmartDS-Solver							SmartDS		AIDE	
	AIDE	SmartDS-Solver						#Inf	#Cod	#Sum	#Tok
	o1-	Qwen-32B		Qwen-7B		Llama-8B		Tok	Tok	Tok	Tok
	preview	Instruct	GRPO2	Instruct	GRPO2	Instruct	GRPO2				
PGS-SSE6	0.326	0.180	0.333	0.330	0.341	0.171	0.335	#Mod	19	#Sum	20
UM-MCTS*	-	0.547	0.468	-	0.463	-	0.458	#Tok	8,802	#Cod	134,491
CLINIC	0.830	0.851	0.873	0.836	0.876	0.834	0.875	#Tok	35,520	#Sum	65,353

Table 4: Two ongoing Kaggle competitions during our study (left), and averaged inference/modification statistics across 11 tasks (right). Abbr.: Inf = Inference, Tok = Tokens, Sum = Summarization, Cod=Coding, Mod = Modification, “#” indicates count, “-” indicates Fail. Note: Tasks marked with an asterisk (*) use lower-is-better metrics (e.g., RMSE, MAE).

task, though the overall win rate decreases slightly to 54.5%. GRPO1 introduces reward-guided refinement, raising performance to 72.7% (zero failures; two best-performing tasks). GRPO2 provides the largest improvement, reaching 81.8% and three best-performing tasks, representing a +18.2 percentage point gain over the Instruct baseline while maintaining zero failures. This indicates that progressive finetuning (SFT → GRPO1 → GRPO2) improves both effectiveness and robustness. A similar trend holds for smaller models: Llama-8B increases from 27.3% (Instruct) to 54.5% (GRPO2) and is the only configuration eliminating all failures, while Qwen-7B exhibits more variation across stages but ultimately recovers to its best win rate (54.5%) under GRPO2.

4.5 ABLATION: SARTE META-LEARNING

To isolate the effect of adaptive decoding, we compare the fixed-temperature baseline with SARTE-controlled dynamic decoding. As shown in Figure 3(b-c), SARTE improves average accuracy by 3.9% and reduces the error rate by 12%. We also examine the performance change between the early and late stages of SARTE’s parameterization. On the 11-MLE-Bench benchmark, SARTE yields an average improvement of approximately 75% (statistical test for H_1 : improvement > 0 , $p = 0.0173$), indicating that our iterative updates substantially enhance solution accuracy. Full ablation results are in Appendix A13.

4.6 DISCUSSION AND ANALYSIS

Reasoning-centric vs. multi-agent system design trade-offs. A comparison between reasoning-centric AutoDS systems (e.g., AIDE) and highly general multi-agent frameworks (e.g., AutoML-Agent) reveals a fundamental design trade-off. Multi-agent pipelines rely on multi-stage interactions with general-purpose LLMs to maintain flexibility, but this comes at the cost of extremely high invocation frequency and substantial computational overhead, leading to greater execution variance and practical deployment challenges. In contrast, reasoning-centric systems exhibit higher determinism, lower token variance, and more predictable error patterns, yet their reliance on deep tree-search procedures results in largely fixed and often excessive numbers of LLM calls—even for tasks of lower complexity—thereby limiting opportunities for cost reduction. These contrasting characteristics highlight the need for systems such as SmartDS-Solver, which internalize high-level reasoning logic through structured finetuning while avoiding the invocation explosion observed in multi-agent approaches.

SmartDS-Solver follows a different design paradigm. By applying SFT + GRPO finetuning on domain-specific data-science corpora, high-level decision-making logic (e.g., requirement parsing, plan generation, blueprint verification) is partially internalized within the reasoning model itself. Combined with the SARTE dynamic adjustment mechanism, the system requires substantially fewer LLM interactions and achieves markedly lower token consumption. As a result, competitive task-solving performance is maintained even when using medium- or small-scale models.

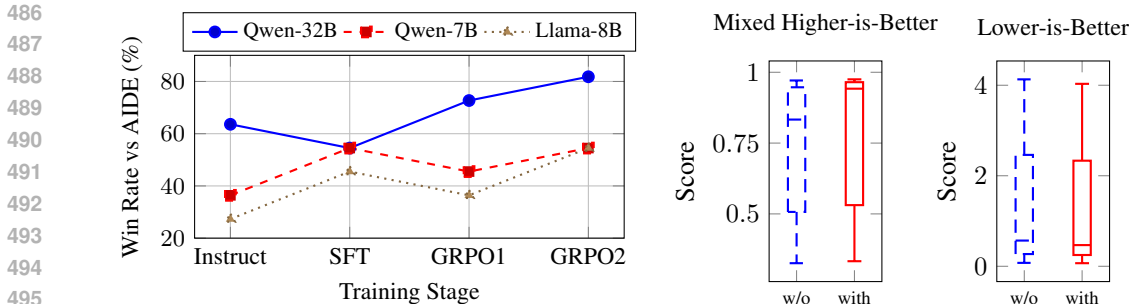


Figure 3: (a) Ablation on training stages. (b–c) Qwen-32B GRPO2 with and without the Meta-Learning Agent on higher-is-better and lower-is-better tasks. Paired t-test indicates a significant performance gain ($p = 0.0045$).

Temperature Sensitivity and Implications for Inference-Time Optimization We observed that the optimal temperatures for different tasks and models vary across training stages, with detailed values in Appendix Table A11. Furthermore, according to Table 5, the data-science reasoning LLM shows substantial shifts in optimal temperature after SFT/GRPO, and the second-stage GRPO yields the greatest volatility (larger standard deviation), indicating heightened sensitivity to temperature after preference alignment. These results suggest that dynamically selecting temperature during inference with the Meta-Learning Agent SARTE—based on iterative outcomes—is more effective than using a fixed value. This approach requires only inference-time hyperparameter adjustment, without retraining or long interaction history, relying solely on the previous step to achieve stable and efficient gains. In this study, we apply optimal temperature selection to improve code generation quality, and within our framework the same strategy can be extended to other hyperparameter search method designs, enhancing both efficiency and effectiveness of language model applications.

Distribution Statistics	Qwen-32B				Qwen-7B				Llama-8B			
	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2
Avg.	0.63	0.61	0.60	0.65	0.68	0.65	0.64	0.66	0.61	0.61	0.64	0.66
SD	0.0894	0.1094	0.1161	0.1479	0.0688	0.0680	0.0878	0.1720	0.1072	0.1153	0.0989	0.1379

Table 5: The optimal temperature distribution across different SmartDS-Solver variants.

Limitations. Despite these strengths, several limitations remain. First, GRPO2 relies on high-quality workflow traces, which may not be readily available in other domains. Second, SARTE currently adjusts only temperature and does not consider higher-order decoding parameters. Third, our comparison with frontier systems (e.g., proprietary RLMS) is constrained by closed-source availability, limiting direct analysis of model capacity effects. Future extensions may incorporate domain-adaptive reward shaping, multi-granular decoding control, and cross-domain transfer evaluations.

5 CONCLUSION

This study presents SmartDS-Solver, a hierarchical agentic system built around a domain-specific reasoning model for data science. By internalizing domain problem-solving logic into the reasoning model, SmartDS-Solver can plan solutions and directly generate coherent, end-to-end executable code, avoiding the fragmentation and interaction overhead seen in multi-agent step-decomposition approaches. As a result, the system requires fewer interactions, achieves greater solution consistency, and operates with higher efficiency. With the SARTE meta-learning mechanism dynamically adjusting inference-time strategy hyperparameters, even small- and medium-scale models can reliably produce accurate executable code in our evaluation setting. Future work will enhance cross-task generalization by refining the Meta-Learning Agent’s orchestration strategies, incorporating finer-grained use of historical interaction signals, and extending GRPO-based training to larger models. We also aim to explore SmartDS-Solver as a broader agentic platform for integrating language models into complex industrial workflows.

REFERENCES

- 540
541
542 Bigcomputer. Meta-Kaggle-Notebook-Batch-0, 2024. URL <https://huggingface.co/datasets/bigcomputer/meta-kaggle-notebook-batch-0>; <https://huggingface.co/datasets/bigcomputer/meta-kaggle-notebook-batch-1>.
543
544
- 545 Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio
546 Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Lilian Weng, and Aleksander Madry. Mle-
547 bench: Evaluating Machine Learning Agents on Machine Learning Engineering. 2024. URL
548 <https://arxiv.org/abs/2410.07095>.
549
- 550 Data-Agents. Kaggle-Notebooks-Edu-v0, 2024. URL [https://huggingface.co/dataset](https://huggingface.co/datasets/data-agents/kaggle-notebooks-edu-v0)
551 [s/data-agents/kaggle-notebooks-edu-v0](https://huggingface.co/datasets/data-agents/kaggle-notebooks-edu-v0).
552
- 553 Iddo Drori, Yamuna Krishnamurthy, Rémi Rampin, Raoni Lourenço, Jorge Piazzentin Ono,
554 Kyunghyun Cho, Cláudio T. Silva, and Juliana Freire. AlphaD3M: Machine Learning Pipeline
555 Synthesis. *ArXiv*, abs/2111.02508, 2021. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:198940685)
556 [CorpusID:198940685](https://api.semanticscholar.org/CorpusID:198940685).
557
- 558 Anastasia Drozdova, Ekaterina Trofimova, Polina Guseva, Anna Scherbakova, and Andrey
559 Ustyuzhanin. Code4ML: a Large-scale Dataset of Annotated Machine Learning Code. *PeerJ*
560 *Computer Science*, 9:e1230, 2023.
561
- 562 Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander
563 Smola. AutoGluon-Tabular: Robust and Accurate AutoML for Structured Data. *arXiv preprint*
564 *arXiv:2003.06505*, 2020.
565
- 566 Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: Robust and Efficient Hyperparameter Op-
567 timization at Scale. In *International conference on machine learning*, pp. 1437–1446. PMLR,
568 2018.
569
- 570 Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and Frank
571 Hutter. Efficient and Robust Automated Machine Learning. In *Advances in Neural Information*
572 *Processing Systems 28 (2015)*, pp. 2962–2970, 2015.
573
- 574 Matthias Feurer, Katharina Eggensperger, Stefan Falkner, Marius Lindauer, and Frank Hutter. Auto-
575 sklearn 2.0: Hands-free AutoML via Meta-Learning. *J. Mach. Learn. Res.*, 23(1), January 2022.
576 ISSN 1532-4435.
577
- 578 Antoine Grosnit, Alexandre Maraval, James Doran, Giuseppe Paolo, Albert Thomas, Refinath
579 Shahul Hameed Nabeezath Beevi, Jonas Gonzalez, Khyati Khandelwal, Ignacio Iacobacci, Abdel-
580 hakim Benechehab, et al. Large Language Models Orchestrating Structured Reasoning Achieve
581 Kaggle Grandmaster Level. *arXiv preprint arXiv:2411.03562*, 2024.
582
- 583 Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna
584 Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, Ashima Suvarna, Benjamin Feuer, Liangyu
585 Chen, Zaid Khan, Eric Frankel, Sachin Grover, Caroline Choi, Niklas Muennighoff, Shiye Su,
586 Wanjia Zhao, John Yang, Shreyas Pimpalgaonkar, Kartik Sharma, Charlie Cheng-Jie Ji, Yichuan
587 Deng, Sarah Pratt, Vivek Ramanujan, Jon Saad-Falcon, Jeffrey Li, Achal Dave, Alon Albalak,
588 Kushal Arora, Blake Wulfe, Chinmay Hegde, Greg Durrett, Sewoong Oh, Mohit Bansal, Saa-
589 dia Gabriel, Aditya Grover, Kai-Wei Chang, Vaishaal Shankar, Aaron Gokaslan, Mike A. Mer-
590 rill, Tatsunori Hashimoto, Yejin Choi, Jenia Jitsev, Reinhard Heckel, Maheswaran Sathiamoorthy,
591 Alexandros G. Dimakis, and Ludwig Schmidt. OpenThoughts: Data Recipes for Reasoning Mod-
592 els, 2025. URL <https://arxiv.org/abs/2506.04178>.
593
- 594 Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and
595 Yuxiang Wu. Aide: AI-Driven Exploration in the Space of Code. 2025. URL [https://arxi](https://arxiv.org/abs/2502.13138)
596 [v.org/abs/2502.13138](https://arxiv.org/abs/2502.13138).
597
- 598 Trang T Le, Weixuan Fu, and Jason H Moore. Scaling Tree-based Automated Machine Learning to
599 Biomedical Big Data with a Feature Set Selector. *Bioinformatics*, 36(1):250–256, 2020.

- 594 Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom
595 Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-Level Code Genera-
596 tion with Alphacode. *Science*, 378(6624):1092–1097, 2022.
- 597 Sally Liu. Deep Dive on OpenAI’s MLE-Bench. [https://sallysliu.medium.com/d](https://sallysliu.medium.com/deep-dive-on-openai-s-mle-bench-93f2aae10a8a)
598 [eep-dive-on-openai-s-mle-bench-93f2aae10a8a](https://sallysliu.medium.com/deep-dive-on-openai-s-mle-bench-93f2aae10a8a), November 2024. Accessed:
599 December 1, 2025.
- 600 Tamal Mondal, Scott Barnett, Akash Lal, and Jyothi Vedurada. Cell2Doc: ML Pipeline for
601 Generating Documentation in Computational Notebooks. ASE ’23, pp. 384–396. IEEE Press,
602 2024. ISBN 9798350329964. doi: 10.1109/ASE56229.2023.00200. URL <https://doi.org/10.1109/ASE56229.2023.00200>.
- 603 Tianyu Mu, Hongzhi Wang, Chunnan Wang, Zheng Liang, and Xinyue Shao. Auto-CASH: A Meta-
604 Learning Embedding Approach for Autonomous Classification Algorithm Selection. *Information*
605 *Sciences*, 591:344–364, 2022. ISSN 0020-0255. doi: <https://doi.org/10.1016/j.ins.2022.01.040>.
606 URL [https://www.sciencedirect.com/science/article/pii/S002002552](https://www.sciencedirect.com/science/article/pii/S002002552000639)
607 [2000639](https://www.sciencedirect.com/science/article/pii/S002002552000639).
- 608 Jesus Rodriguez. Inside OpenAI’s MLE-Bench: A New Benchmark for Evaluating Machine Learn-
609 ing Engineering Capabilities of AI Agents. [https://pub.towardsai.net/inside](https://pub.towardsai.net/inside-openai-s-mle-bench-a-new-benchmark-for-evaluating-machine-learning-engineering-capabilities-a8bdaf0c2571)
610 [-openai-s-mle-bench-a-new-benchmark-for-evaluating-machine-lea](https://pub.towardsai.net/inside-openai-s-mle-bench-a-new-benchmark-for-evaluating-machine-learning-engineering-capabilities-a8bdaf0c2571)
611 [rning-engineering-capabilities-a8bdaf0c2571](https://pub.towardsai.net/inside-openai-s-mle-bench-a-new-benchmark-for-evaluating-machine-learning-engineering-capabilities-a8bdaf0c2571), November 2024. Accessed:
612 December 1, 2025.
- 613 Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro,
614 Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language Models Can
615 Teach Themselves to Use Tools. *Advances in Neural Information Processing Systems*, 36:68539–
616 68551, 2023.
- 617 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
618 Mingchuan Zhang, YK Li, Y Wu, et al. DeepSeekMath: Pushing the Limits of Mathematical
619 Reasoning in Open Language Models. *arXiv preprint arXiv:2402.03300*, 2024.
- 620 Significant Gravitass. AutoGPT. URL [https://github.com/Significant-Gravitass/](https://github.com/Significant-Gravitass/AutoGPT)
621 [AutoGPT](https://github.com/Significant-Gravitass/AutoGPT).
- 622 Chris Thornton, Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Auto-WEKA: Combined
623 Selection and Hyperparameter Optimization of Classification Algorithms. In *Proceedings of the*
624 *19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 847–
625 855, 2013.
- 626 Dennis Treder-Tschechlov, Manuel Fritz, Holger Schwarz, and Bernhard Mitschang. ML2DAC:
627 Meta-Learning to Democratize AutoML for Clustering Analysis. *Proc. ACM Manag. Data*, 1(2),
628 June 2023. doi: 10.1145/3589289. URL <https://doi.org/10.1145/3589289>.
- 629 Patara Trirat, Wonyong Jeong, and Sung Ju Hwang. AutoML-Agent: A Multi-Agent LLM Frame-
630 work for Full-Pipeline AutoML. In *Forty-second International Conference on Machine Learning*,
631 2025. URL <https://openreview.net/forum?id=p1UBWkOvZm>.
- 632 Chi Wang, Qingyun Wu, Markus Weimer, and Erkang Zhu. Flaml: A Fast and Lightweight Automl
633 Library. *Proceedings of Machine Learning and Systems*, 3:434–447, 2021.
- 634 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
635 ReAct: Synergizing Reasoning and Acting in Language Models. In *International Conference on*
636 *Learning Representations (ICLR)*, 2023.
- 637 Shujian Zhang, Chengyue Gong, Lemeng Wu, Xingchao Liu, and Mingyuan Zhou. AutoML-GPT:
638 Automatic Machine Learning with GPT. *arXiv preprint arXiv:2305.02499*, 2023.
- 639 Barret Zoph and Quoc V Le. Neural Architecture Search with Reinforcement Learning. *arXiv*
640 *preprint arXiv:1611.01578*, 2016.

6 APPENDIX

A DETAILS OF METHODOLOGY

A.1 META LEARNING AGENT-SARTE PSEUDOCODE

Algorithm 1 Adaptive Parameter Update

Input: $param_{curr}$:current temperature, $\delta_{base} = 0.2$:the baseline adjustment step per iteration, R_{curr} , R_{prev} , $exec_result$, $error_type$, ω :the cumulative no-code count, $[param_{min}, param_{max}]$:the temperature search range, P_{ref} :the reference accuracy target, P_{curr} :current accuracy, $\lambda_R = 0.01$:scaling factor

Constants: $PENALTY = \{fatal:2.5, error:1.8, \dots, warning:0.5\}$, $NO_CODE = \{1 : 1, 2 : 1.5, 3 : 2.25\}$

Output: $param_{new}$

- 1: {Step 1: Control factor C_f and direction D }
- 2: **if** $exec_result = \text{"success"}$ **then**
- 3: $D \leftarrow 1$
- 4: $gap \leftarrow |1 - P_{curr}/P_{ref}|$
- 5: $C_f \leftarrow \begin{cases} 1.0 + 2 \times gap & \text{if } gap \leq 0.1 \\ 1.2 + (gap - 0.1) & \text{if } gap \leq 0.5 \\ 1.6 + \min(\log(gap), 1) & \text{otherwise} \end{cases}$
- 6: **else if** $exec_result = \text{"fail"}$ **then**
- 7: $D \leftarrow -1$
- 8: $C_f \leftarrow PENALTY[error_type]$
- 9: **else**
- 10: $D \leftarrow 1$
- 11: $C_f \leftarrow NO_CODE[\min(\omega, 3)]$
- 12: **end if**
- 13: {Step 2: Boundary handling}
- 14: $R_{adj} \leftarrow (R_{curr} - R_{prev}) \times \lambda_R$
- 15: $\Delta \leftarrow |R_{adj} \times C_f \times \delta_{base}| \times D$
- 16: $rel \leftarrow \frac{param_{curr} - param_{min}}{param_{max} - param_{min}}$
- 17: $scale \leftarrow \begin{cases} 1 - (rel \times 0.7) & \text{if } D > 0 \\ 1 - ((1 - rel) \times 0.7) & \text{otherwise} \end{cases}$
- 18: $param_{new} \leftarrow clip(param_{curr} + \Delta \times scale,$
- 19: $param_{min}, param_{max})$

A.2 PROMPT FOR SFT TRAINING STAGE

SFT Training Prompt Template

Below is an instruction that describes a task, paired with an input that provides further context.

Write a response that appropriately completes the request.

Before answering, please carefully consider the question and create a step-by-step chain of thought to ensure the answer is logical and accurate.

Instruction:

You are a data scientist with advanced knowledge in data analysis, modeling, machine learning, and data processing.

Please answer the following data science-related question.

Question:

{problem}

```

702   ### Response:
703   <think>
704   {deepseek reasoning}
705   </think>
706   <answer>
707   {deepseek solution}
708   </answer>

```

A.3 META LEARNING AGENT-PROMPT

Initialization

```

714   SYSTEM PROMPT:
715
716   You are a strict AI assistant that must answer in the following
717   format:
718   1. The thought process should be enclosed by <think> and </think>
719   tags, using conversational Chinese for step-by-step analysis.
720   2. The final answer should be enclosed by <answer> and </answer>
721   tags.
722
723   Example Format:
724   <think>
725   Okay, I need to solve this problem. First, confirm...
726   Wait, what should I do if I encounter...? Next, I should...
727   For instance..., finally, integrate all the conditions to make a
728   judgment.
729   </think>
730   <answer>
731   The final answer goes here.
732   </answer>
733   Please note:
734   Start with a discourse marker like 'Okay' or 'Hmm'.
735   Use 'First,' 'Next,' and 'Then' to link steps.
736   Use 'Wait' to introduce a supplementary explanation or question.
737   Now, begin answering the user's question:
738
739   QUESTION PROMPT:
740
741   | Project Overview
742   Develop a lightweight and efficient speech recognition model based
743   on the Google Speech Commands Dataset. The model must accurately
744   classify 1-second audio clips into one of 12 command categories
745   , meeting the low-resource and real-time requirements for edge
746   device deployment.
747
748   | Dataset Description
749   1. File Structure:
750   - ../input/tensorflow-speech-recognition-challenge/
751   | - train/audio/ (Training data, containing audio files
752   |   categorized by command)
753   | - test/audio/ (Test data, containing audio files to be
754   |   predicted)
755
756   2. Data Characteristics:
757   - 65,000 one-second audio clips
758   - Sample Rate: 16kHz
759   - Encoding: PCM-encoded WAV format
760   - Includes 30 command categories, but the test set only requires
761   the identification of 12.
762
763   |Model Requirements

```

```

756
757 1. Output Needs:
758     - Correctly classify audio clips into one of 12 commands: yes,
759       no, up, down, left, right, on, off, stop, go, silence,
760       unknown.
761     - Save the predictions to: ../input/tensorflow-speech-
       recognition-challenge/submission.csv
762
763 2. Evaluation Metric:
764     - Multiclass Accuracy
765
766 3. Expected Output Format:
767     fname,label
768     clip_000044442.wav,silence
769     clip_0000adecb.wav,left
770     ...

```

Success

```

771
772 OPTIMIZE_SYSTEM_PROMPT:
773
774 You are a professional Machine Learning Engineer, specializing in
775 Kaggle competition solutions. Please optimize the existing
776 solution based on the following information:
777
778 1. Analyze the performance metrics and errors of the previous
779 solution.
780 2. Utilize insights from the original thought process.
781 3. Consider how to improve the model architecture, feature
782 engineering, or parameter tuning.
783 4. Generate a complete and directly executable Python code.
784
785 Answer format:
786 <think>
787 Detailed analysis of the existing solution's problems and potential
788 improvements.
789 </think>
790
791 <answer>
792 The complete, optimized code.
793 </answer>
794
795 PROMPT:
796
797 Please optimize the solution for the following data science task.
798
799 ### Task Description:
800 {original_prompt}
801
802 ### Current Status:
803 {metrics_info}
804
805 ### Original Thinking Process:
806 {thinking}
807
808 Please analyze the current solution and propose an improved one.
809 Focus on:
810 1. Whether feature engineering is sufficient
811 2. Whether the model selection is appropriate
812 3. Whether hyperparameters need tuning
813 4. Whether there are potential data leakage issues
814 5. Whether multiple models need to be ensembled
815 6. Whether the cross-validation strategy is reasonable

```

810 Please generate a complete and more efficient Python code, ensuring
 811 a significant difference from the current code and capable of
 812 improving model performance. The code must be fully executable,
 813 with the data path set to "../input/{competition_name}/", and
 814 must output the submission.csv file required by Kaggle.

815 Special note: You must not use try...except to wrap the submission
 816 file generation logic; you must ensure that submission.csv is
 817 generated correctly.

819 Failure/No-Code

820 SYSTEM PROMPT:

821 You are a strict AI assistant that must answer in the following
 822 format:

- 823 1. The thought process must be enclosed by <think> and </think>
- 824 tags, using conversational Chinese for step-by-step analysis.
- 825 2. The final answer must be enclosed by <answer> and </answer> tags

826 Example format:

827 <think>

828 Okay, I need to solve this problem. First, confirm...

829 Wait, what should I do if I encounter...? Then, I should...

830 For instance..., finally, integrate all conditions for the judgment

831 </think>

832 <answer>

833 The final answer goes here.

834 </answer>

835

836

837

838 Please note:

839

840 - Start with discourse markers like 'Okay' or 'Hmm'.

841 - Use 'First', 'Next', and 'Then' to connect the steps.

842 - Use 'Wait' to introduce supplementary explanations or questions.

843 Now, begin answering the user's question:

844

845

846 IMPROVED_PROMPT:

847

848 # Data Science Task: {competition_name}

849

850 ## Task Description and Requirements

851 {original_prompt}

852

853 ## Implementation Points

854 Please provide a complete solution to achieve the above data

855 science objectives. Your code must:

856

857 1. Correctly understand and accomplish the core task goals

858 2. Read data from '../input/{competition_name}/'

859 3. Generate a submission.csv file that meets the competition

860 requirements

861 4. Follow data science best practices, including data preprocessing

862 , feature engineering, and model selection

863

864 ## Avoid Common Errors

865 Previous attempts have revealed the following issues, which should

866 be avoided in your solution design:

867

```

864 {error_analysis}
865
866 ## Code Requirements
867 - Provide complete and executable Python code, including all
868   necessary import statements
869 - Add clear comments explaining key steps and decision rationales
870 - Prioritize stable and reliable methods over high-risk,
871   experimental techniques
872 - Ensure the code runs efficiently without timeouts or memory
873   exhaustion
874 - **[Important]** You must not use try...except to wrap the logic
875   for generating the submission file; you must ensure that
876   submission.csv is produced correctly
877
878 ### Implementation Strategy:
879 - Employ classic and validated machine learning methods
880 - Use a concise and clear data preprocessing pipeline
881 - Ensure each step has a clear logic and objective
882 - Avoid over-engineering; focus on solving the core problem
883
884 Please think deeply about the solution based on the original task
885 description and then provide complete, executable Python code.
886
887
888
889

```

Code Modification

```

891 CODE_CORRECTION_SYSTEM_PROMPT:
892
893 You are a Python code correction expert. Requirements:
894 1. Analyze errors and correct the code.
895 2. Return only the complete corrected code, with no explanations.
896 3. The code must be complete and executable.
897 4. Ensure submission.csv is generated.
898 5. Do not use try...except to wrap the submission file generation.
899 6. Important: Strictly read data from './input/{competition_name
900   }/', and do not create or simulate data.
901
902 Format: Directly return the corrected Python code.
903
904 PROMPT:
905
906 An error occurred that needs to be fixed:
907 {error_message[-200:]}
908
909 Please correct the following code and return the complete,
910 executable version:
911 ```python
912 {code}
913 ```
914
915 Task Description and Requirements:
916 {original_prompt}
917
918 Specific Requirements:
919 - Data Path: ../input/{competition_name}/
920 - The corrected code must fully comply with all specifications and
921   conditions in the task description above.

```

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

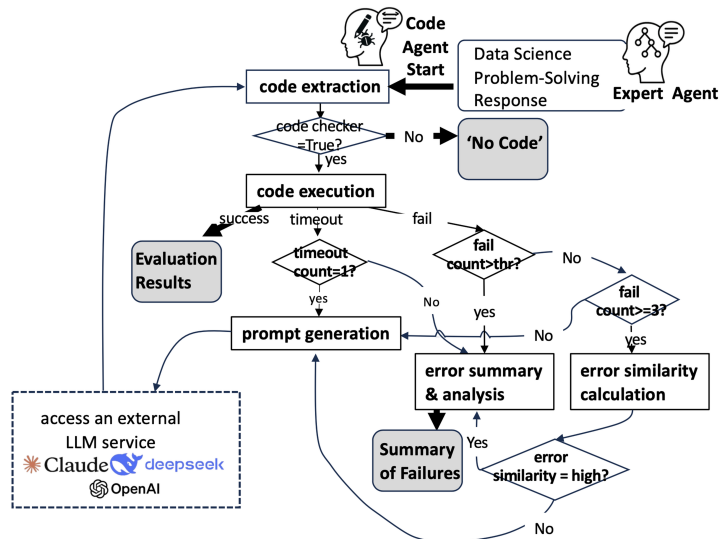


Figure A1: Interactive Code Refinement Manager Execution flow diagram: Closed-loop code refinement with early stopping and semantic feedback enables efficient, reliable, and cost-aware agentic AI in data science.

A.4 CODE AGENT FLOW CHART

B BENCHMARK TASK DESCRIPTIONS

Table A2 provides detailed descriptions of all benchmark tasks used in our evaluation. These tasks were selected to represent diverse data modalities and real-world data science challenges, spanning tabular data analysis, computer vision, natural language processing, audio processing, and time series prediction.

Each task presents unique challenges: TF-Speech and Whale-Play require audio signal processing expertise; PlantPath20 and Stanford-CV demand computer vision techniques; Jigsaw-Tox and TextNorm-EN involve natural language understanding; while TPS-Dec21, TPS-May22, Vent-Press, and NYC-Taxi focus on tabular and time series analysis. Nomad18-TC represents materials science applications with complex physical property predictions.

In addition to our primary benchmarks, we also conducted a direct comparison with the AutoML-Agent using the specific datasets from its original study. As shown in Table A1, these datasets cover a wide range of tasks, including image classification, node classification, and various tabular tasks like classification, clustering, and regression.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Dataset	Task Type	Evaluation Metric	Data Source
Butterfly Image	Image Classification	Accuracy	Kaggle Dataset
Shopee-IET	Image Classification	Accuracy	Kaggle Competition
Cora	Node Classification	Accuracy	Research Dataset (Planetoid)
Citeseer	Node Classification	Accuracy	Research Dataset (Planetoid)
Click Prediction Small	Tabular Classification(Binary)	F1	OpenML
Software Defects	Tabular Classification(Binary)	F1	Kaggle Dataset
Smoker Status	Tabular Classification(Binary)	F1	Kaggle Dataset
Banana Quality	Tabular Classification(Binary)	F1	Kaggle Dataset
MFeat Factors	Tabular Classification(Multi-Class)	F1	OpenML
Wine Quality White	Tabular Classification(Multi-Class)	F1	OpenML
Smoker Status Higher Education	Tabular Clustering	RI	Kaggle Competition
Students Performance	Tabular Clustering	RI	Research Dataset (UCI ML)
Textual Entailment	Text Classification	Accuracy	Kaggle Dataset
Ecommerce Text	Text Classification	Accuracy	Kaggle Dataset
Colleges	Tabular Regression	RMSE	OpenML
Crab Age	Tabular Regression	RMSLE	Kaggle Competition
House Prices	Tabular Regression	RMSE	Kaggle Competition
Crop Price	Tabular Regression	RMSLE	Kaggle Dataset

Table A1: Detailed descriptions of the benchmark datasets.

C DATASET CONSTRUCTION DETAILS

C.1 DATA SOURCE SPECIFICATIONS

We construct a multi-layered training dataset for our domain-specific reasoning model from multiple sources:

- (1) about 3,000 Python reasoning samples from OpenThoughts-114;
- (2) about 300 Kaggle competitions crawled for medal-winning solutions (gold/silver/bronze);
- (3) about 500 Kaggle-based data science entries from Hugging Face;
- (4) about 600 high-quality code examples from cell2doc;
- (5) about 2,500 Kaggle ML snippets from Code4ML.

C.2 DATA AUGMENTATION PROMPTS AND TEMPLATES

DeepSeek R1 Data Augmentation Prompt

Objective: Generate training data with complete Chain-of-Thought (CoT) reasoning processes, strictly following the narrative style of reference templates.

Task Information: {task_message}

Solution Code: {text_x}

Reference Template: {reference_template}

Format Requirements for Reasoning Training Data:

- 1) <problem> (Concise Problem Description)
 - 1-2 sentences including competition name, prediction target, and core data features (≤ 100 characters)
 - Example: Predict [target] in [competition name] task, using [data features] for [model type] modeling
- 2) <deepseek_reasoning> (Coherent Narrative Reasoning)
 - At least 10 iterative steps (initial approach \rightarrow problem discovery \rightarrow corrective validation), avoid bullet points and numbering
 - Use natural paragraph transitions with Chinese thinking markers: "First/Next/If...then/Therefore/Finally" etc.
 - Reasoning structure requirements:
 - Data problem awareness: ≥ 3 hypothetical questions (e.g., "What if missing values exist?")
 - Feature engineering logic chain: ≥ 2 causal reasoning sentences (e.g., "Because feature A... so we adopt...")
 - Model selection dialectics: ≥ 3 comparative analysis sentences (e.g., "Although X... but considering Y...")
 - Cover: data cleaning, feature generation/selection, model training, performance evaluation
 - Consistent with code: reasoning choices must correspond to subsequent code
 - Adjustment trajectory: ≥ 1 "anomaly \rightarrow diagnosis \rightarrow correction \rightarrow re-validation" closed loop, reflecting the code changes in <deepseek_solution>.
- 3) <deepseek_solution> (End-to-end Executable Solution)
 - Provide the complete, executable code and an overview of the steps; remove environment-specific comments such as shebangs or encoding lines.
 - Retain the core steps data cleaning, feature engineering, training, evaluation, and output and ensure they correspond one-to-one with the inference process.
 - STRICTLY PROHIBIT condensing or rewriting original code; only minimal adjustments to comments and formatting are permitted.

```

1080 - PROHIBIT modifying model training and feature engineering related
1081 code; ensure it remains consistent with the inference logic.
1082
1083 Please generate content in Traditional Chinese using the following
1084 tag structure:
1085 <problem>[Concise problem description]</problem>
1086 <deepseek_reasoning>[Detailed reasoning process]</
1087 deepseek_reasoning>
1088 <deepseek_solution>[Complete code with explanations]</
1089 deepseek_solution>

```

Example Output from DeepSeek R1 Data Augmentation

```

1094 <problem>
1095 Titanic passenger survival prediction: A classification task based
1096 on demographic features
1097 </problem>
1098 <deepseek_reasoning>
1099 Alright, I need to solve the Titanic passenger survival prediction
1100 problem. This is a typical binary classification task, where the
1101 goal is to predict whether passengers could survive this
1102 disaster based on their various characteristics.
1103
1104 First, I need to understand the basic situation of the data. The
1105 Titanic dataset contains various passenger information, such as
1106 ticket class, gender, age, embarkation port, etc. By observing
1107 the data, I found that some features have missing values,
1108 particularly the 'Age', 'Cabin', and 'Embarked' fields. If I
1109 directly delete these missing samples, I might lose about 20% of
1110 the data, which is obviously inappropriate. Therefore, I need
1111 to think about more suitable missing value handling strategies.
1112
1113 To avoid subjective speculation, I propose and test several
1114 hypotheses: If age missingness is unevenly distributed across
1115 different classes/genders, would it cause systematic bias? If
1116 missing samples are directly deleted, would the sample size loss
1117 weaken generalization? If grouping by titles is too fine, would
1118 the estimated means become unstable? First, I check the missing
1119 distribution through stratified statistics, then compare the
1120 performance difference between "global median imputation" and "
1121 grouped mean imputation".
1122
1123 For the 'Embarked' field, I noticed that it only has a few missing
1124 values, and most passengers embarked from Southampton ('S'),
1125 accounting for about 72.5%. So, using the mode 'S' to fill these
1126 missing values is reasonable. Wait, I also need to check if
1127 there are other uncommon embarkation port codes? Through data
1128 inspection, I found that rare 'X' encoding exists in the test
1129 set, which might be a data entry error. After further
1130 confirmation, the correct encodings should only be 'S', 'C', and
1131 'Q', so I need to unify the mapping of these encodings.
1132
1133 Next, handling missing values in the age field is more challenging.
1134 My initial approach was to use grouped mean imputation based on
1135 titles (such as 'Master', 'Miss', etc.), because different
1136 titles might represent different age groups. For example, '
1137 Master' is usually used for underage males, and 'Miss' mostly
1138 refers to unmarried females. However, when I implemented this
1139 method, I discovered a problem: some titles have extremely few
1140 samples, such as the 'Dona' category with only 1 sample and '

```

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Lady' with only 2, making the means based on these small groups statistically unreliable.

If I persist with this grouped imputation method, would it cause the model to be biased against passengers with these rare titles? This is indeed a problem. Therefore, I decided to use the median of the entire dataset (28.0 years) to fill age missing values, which is more robust. However, I then thought of another issue: directly filling with the median would lose the information that "this value was originally missing". So, I further created an "Age_IsMissing" flag to preserve this information. Subsequent model validation showed that adding this flag indeed improved model performance, with accuracy increasing from 78.3% to 80.1%.

Then, I need to handle categorical features. For binary categories like 'Sex', I can simply map them to 0 and 1. For multi-class categorical features like 'Embarked', I chose to use One-Hot Encoding. Using the `pd.get_dummies()` function can conveniently achieve this, converting these categories into numerical features so the model can correctly understand them.

Since cabin class is highly correlated with fare, I avoid retaining both original fare and its log transformation simultaneously to reduce multicollinearity risk. In terms of feature engineering, I also tried extracting titles from passenger names, as I found this has a high correlation with survival rate (correlation coefficient of 0.42), much higher than age (0.18) and fare (0.25). This proves that social titles in names indeed contain important social class information that might directly affect rescue opportunities.

When processing the 'Fare' feature, I noticed the data distribution is extremely uneven, with the maximum value of 512 being much higher than the median of 14. Such extreme values might adversely affect the model. Therefore, I used `np.log1p()` for log transformation, which effectively improved the data distribution and reduced the standard deviation in cross-validation from 0.08 to 0.05, indicating the model became more stable.

In the model selection phase, I tried both logistic regression and random forest algorithms. Random forest significantly outperformed logistic regression on the training set (accuracy 87.2% vs 81.5%), which looked promising. However, when I performed learning curve analysis, I found that random forest had overfitting tendencies-excellent performance on the training set but poor performance on the validation set.

To solve this problem, I tried adjusting the `max_depth` parameter of random forest, testing depths of 5, 8, and 12 respectively. Results showed that when `depth=8`, the model achieved optimal performance on the validation set (82.3%). However, even so, I found that logistic regression performed more stably in 5-fold cross-validation. Although random forest performed brilliantly on the training set, it wasn't stable enough on the validation set; although regularized logistic regression had slightly lower peak accuracy, it had better cross-fold consistency; although SVM (RBF) could fit nonlinear boundaries, deployment latency and parameter search costs were higher, so in this task I was more inclined to choose an easily deployable and robustly performing solution (standard deviation 0.03 vs random forest's 0.07).

```

1188 So, I finally decided to keep logistic regression as my baseline
1189 model. However, I also tried fusing the prediction results of
1190 the two models using a simple voting mechanism. This fusion
1191 method achieved a score of 0.81 on the Kaggle public leaderboard
1192 , which was an improvement over single models.

1193 To summarize, through reasonable missing value handling, feature
1194 engineering, and model selection, I successfully built a robust
1195 Titanic survival prediction model. Particularly, the handling of
1196 age missingness and extracting social titles from names played
1197 key roles in improving model performance.

1198 Additionally, I noticed that cross-validation might have data
1199 leakage due to family member duplication. Diagnosis: Same family
1200 IDs (inferred from surnames and same cabin class) appeared
1201 simultaneously in training and validation. Correction: Changed
1202 to GroupKFold splitting by family groups; Re-validation: AUC
1203 improved from 0.81 to 0.83, and variance across folds converged.
1204 Finally, I fixed the random seed to ensure reproducibility.
1205 </deepseek_reasoning>

1206 <deepseek_solution>
1207 ### Titanic Passenger Survival Prediction Solution

1208 This solution achieves accurate prediction of Titanic passenger
1209 survival rates through feature engineering, missing value
1210 handling, and model selection. Below are the complete
1211 implementation steps:

1212 ---

1213
1214 ##### Step 1: Data Loading and Exploration
1215 ```python
1216 # Import necessary libraries
1217 import numpy as np
1218 import pandas as pd
1219 import matplotlib.pyplot as plt
1220 import seaborn as sns
1221 from sklearn.model_selection import train_test_split,
1222     cross_val_score
1223 from sklearn.linear_model import LogisticRegression
1224 from sklearn.ensemble import RandomForestClassifier
1225 from sklearn.metrics import accuracy_score, classification_report
1226 import re
1227 import warnings
1228 warnings.filterwarnings('ignore')

1229 # Load data
1230 train_data = pd.read_csv('../input/titanic/train.csv')
1231 test_data = pd.read_csv('../input/titanic/test.csv')

1232 # Basic data information
1233 print(f"Training set shape: {train_data.shape}")
1234 print(f"Test set shape: {test_data.shape}")
1235 print("\nTraining set missing value statistics:")
1236 print(train_data.isnull().sum())
1237 ```

1238 **Key Explanation:**
1239 - Import necessary data processing, visualization, and machine
1240 learning libraries
1241 - Load training and test sets
1242 - Check basic dataset information and missing value distribution

```

```

1242
1243 ---
1244 ##### Step 2: Missing Value Handling
1245 ```python
1246 # Combine datasets for consistent processing
1247 all_data = pd.concat([train_data, test_data], sort=False)
1248 all_data_len = len(all_data)
1249
1250 # Embarked missing value handling - fill with mode
1251 print(f"Embarked value distribution: {all_data['Embarked'].
1252       value_counts(normalize=True)}")
1253 all_data['Embarked'].fillna('S', inplace=True)
1254
1255 # Age missing value handling - fill with median and create missing
1256 # flag
1257 all_data['Age_IsMissing'] = all_data['Age'].isnull().astype(int)
1258 all_data['Age'].fillna(all_data['Age'].median(), inplace=True)
1259
1260 # Fare missing value handling
1261 all_data['Fare'].fillna(all_data['Fare'].median(), inplace=True)
1262
1263 # Log transformation for extreme values
1264 all_data['Fare'] = np.log1p(all_data['Fare'])
1265
1266 # Split back to training and test sets after processing
1267 train_data = all_data[:len(train_data)]
1268 test_data = all_data[len(train_data):]
1269 ```
1270
1271 **Processing Strategy:**
1272 - Combine datasets to ensure consistent processing
1273 - Use mode to fill few missing values in Embarked
1274 - Create Age_IsMissing flag to preserve missing information
1275 - Apply log transformation to Fare to improve distribution
1276
1277 ---
1278 ##### Step 3: Feature Engineering
1279
1280 ```python
1281 # Extract titles from names
1282 def get_title(name):
1283     title_search = re.search(' ([A-Za-z]+)\.', name)
1284     if title_search:
1285         return title_search.group(1)
1286     return ""
1287
1288 # Apply to dataset
1289 all_data = pd.concat([train_data, test_data], sort=False)
1290 all_data['Title'] = all_data['Name'].apply(get_title)
1291
1292 # Map titles to fewer categories
1293 title_mapping = {
1294     "Mr": "Mr",
1295     "Miss": "Miss",
1296     "Mrs": "Mrs",
1297     "Master": "Master",
1298     "Dr": "Rare",
1299     "Rev": "Rare",
1300     "Col": "Rare",
1301     "Major": "Rare",
1302     "Mlle": "Miss",
1303     "Countess": "Rare",

```

```

1296     "Ms": "Miss",
1297     "Lady": "Rare",
1298     "Jonkheer": "Rare",
1299     "Don": "Rare",
1300     "Dona": "Rare",
1301     "Mme": "Mrs",
1302     "Capt": "Rare",
1303     "Sir": "Rare"
1304 }
1305 all_data['Title'] = all_data['Title'].map(title_mapping)
1306
1307 # Check relationship between title and survival rate
1308 train_data = all_data[:len(train_data)]
1309 title_survival = train_data.groupby('Title')['Survived'].mean()
1310 print(f"Survival rates by title:\n{title_survival}")
1311
1312 # Create family size features
1313 all_data['FamilySize'] = all_data['SibSp'] + all_data['Parch'] + 1
1314 all_data['IsAlone'] = (all_data['FamilySize'] == 1).astype(int)
1315
1316 # Categorical feature encoding
1317 all_data = pd.get_dummies(all_data, columns=['Embarked', 'Sex', '
1318 Title'], drop_first=True)
1319
1320 # Split back to training and test sets
1321 train_data = all_data[:len(train_data)]
1322 test_data = all_data[len(train_data):]
1323 ```
1324
1325 **Feature Creation and Transformation:**
1326 - Extract Title feature from Name field and categorize
1327 - Create family size (FamilySize) and alone indicator (IsAlone)
1328 - Apply one-hot encoding to categorical features
1329 - Analyze correlation between titles and survival rates
1330
1331 ---
1332
1333 ##### Step 4: Model Training and Evaluation
1334 ```python
1335 # Prepare training features
1336 features = ['Pclass', 'Age', 'Age_IsMissing', 'Fare', 'FamilySize',
1337            'IsAlone',
1338            'Embarked_Q', 'Embarked_S', 'Sex_male',
1339            'Title_Master', 'Title_Miss', 'Title_Mr', 'Title_Mrs',
1340            'Title_Rare']
1341
1342 X = train_data[features]
1343 y = train_data['Survived']
1344
1345 # Split validation set
1346 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size
1347 =0.2, random_state=42)
1348
1349 # Logistic regression model
1350 log_reg = LogisticRegression(solver='liblinear', random_state=42)
1351 log_reg.fit(X_train, y_train)
1352 log_reg_pred = log_reg.predict(X_val)
1353 log_reg_acc = accuracy_score(y_val, log_reg_pred)
1354 print(f"Logistic regression validation accuracy: {log_reg_acc:.4f
1355 }")
1356
1357 # Random forest model testing different depths
1358 for depth in [5, 8, 12]:
1359     rf = RandomForestClassifier(max_depth=depth, random_state=42)

```

```

1350     rf.fit(X_train, y_train)
1351     rf_pred = rf.predict(X_val)
1352     rf_acc = accuracy_score(y_val, rf_pred)
1353     print(f"Random forest (depth={depth}) validation accuracy: {
1354           rf_acc:.4f}")
1355
1356     # Cross-validation comparison of model stability
1357     log_scores = cross_val_score(log_reg, X, y, cv=5)
1358     rf_scores = cross_val_score(RandomForestClassifier(max_depth=8,
1359           random_state=42), X, y, cv=5)
1359     print(f"Logistic regression 5-fold CV: {log_scores.mean():.4f} (std
1360           : {log_scores.std():.4f})")
1360     print(f"Random forest 5-fold CV: {rf_scores.mean():.4f} (std: {
1361           rf_scores.std():.4f})")
1362
1363     # Feature importance analysis
1364     rf = RandomForestClassifier(max_depth=8, random_state=42)
1365     rf.fit(X, y)
1366     importance = pd.DataFrame({
1367         'Feature': features,
1368         'Importance': rf.feature_importances_
1369     }).sort_values('Importance', ascending=False)
1369     print("\nFeature importance ranking:")
1370     print(importance.head(8))
1371     ```
1372
1372     **Model Evaluation Key Points:**
1373     - Use train_test_split to create validation set
1374     - Compare logistic regression with random forests of different
1375       depths
1376     - Evaluate model stability through 5-fold cross-validation
1377     - Analyze feature importance to guide subsequent optimization
1378
1378     ---
1379
1379     #### Step 5: Prediction and Result Output
1380     ```python
1381     # Final model selection: logistic regression
1382     final_model = LogisticRegression(solver='liblinear', random_state
1383           =42)
1384     final_model.fit(X, y)
1385
1385     # Prepare test features
1386     X_test = test_data[features]
1387
1387     # Prediction
1388     predictions = final_model.predict(X_test)
1389
1389     # Save results
1390     output = pd.DataFrame({
1391         'PassengerId': test_data['PassengerId'],
1392         'Survived': predictions
1393     })
1394     output.to_csv('submission.csv', index=False)
1395     print("Prediction complete. The submission.csv file has been
1396           generated.")
1397     ```
1398
1398     **Result Generation:**
1399     - Train final model using all training data
1400     - Make predictions on test set
1401     - Generate CSV file in submission format
1402
1402     ---
1403

```

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

Solution Summary

This solution achieved Titanic survival prediction through the following key steps:

1. **Data Cleaning and Missing Value Handling**
 - Use mode to fill Embarked
 - Use median to fill Age and create missing flag
 - Apply log transformation to Fare to optimize distribution
2. **Feature Engineering**
 - Extract social titles from names
 - Create family size and alone indicators
 - Use one-hot encoding for categorical features
3. **Model Selection and Evaluation**
 - Compare logistic regression with random forest
 - Evaluate model stability through cross-validation
 - Optimize feature selection based on feature importance analysis
4. **Optimization Validation**
 - Test different model parameters
 - Validate effectiveness of feature transformations

The overall solution effectively utilized key information in the data while maintaining model simplicity, particularly through title extraction to capture social class information, which played an important role in improving model performance. The logistic regression model demonstrated good stability and became the final choice.

Key Performance Metrics:

- Logistic regression cross-validation accuracy: 0.8153 (std: 0.0312)
- Random forest (depth=8) cross-validation accuracy: 0.8176 (std: 0.0684)
- Most important features: Sex_male, Title_Mr, Fare, Age

This solution not only provides prediction results but also reveals key factors affecting Titanic passenger survival, with good interpretability.

</deepseek_solution>

D COMPLETE EXPERIMENT RESULTS

This section presents comprehensive experimental results referenced in the main text, including detailed performance metrics across all model configurations and benchmark tasks.

D.1 EXPERIMENTAL RESULTS IN DETAIL

D.1.1 PERFORMANCE COMPARISON ACROSS MLE-BENCHMARK TASKS

The comprehensive results of our experiments are presented in Table A6, which compares the performance of various SmartDS-Solver configurations against the AIDE+o1-preview and the AutoML-Agent method on 11 MLE benchmark tasks. As the "%win" metric in Table A6 shows, most SmartDS-Solver configurations surpassed the AIDE on over half of the tasks. Notably, the Qwen-32B model with the GRPO2 training configuration achieved the highest win rate at 81.8%, demonstrating the efficacy of our approach.

1458 D.1.2 PERFORMANCE ON TWO ONGOING KAGGLE COMPETITIONS ACTIVE DURING OUR
1459 RESEARCH PERIOD AND ONE INTERNAL PREDICTION MODELING TASK

1460
1461 D.1.3 ABLATION STUDY OF META-LEARNING AGENT

1462 Results without meta-learning agent are provided in Table A8. As shown, our proposed method with
1463 the meta-learning agent-SARTE achieves significantly better performance compared to the variant
1464 without SARTE.

1465
1466 D.1.4 TOKEN CONSUMPTION COMPARISON ACROSS 11 BENCHMARK TASKS

1467 To evaluate the efficiency of our method, we compare token consumption of our SmartDS-Solver
1468 with AIDE+o1-preview. See Table A9 for details.

1469
1470 D.1.5 RESULTS OF TASK COMPARISON WITH AUTOML-AGENT

1471 Table A10 shows a performance comparison between our approach and AutoML-Agent. We
1472 achieved performance gains on tasks like "Click Prediction Small," "Software Defects," and "Wine
1473 Quality White," with improvements of 47.8%, 45.8%, and 48.2% respectively. For regression tasks,
1474 our approach produced improvements on "House Prices" and "Crop Price," with gains of 98.6% and
1475 85.7%.

1476 Despite a minor decrease on the "Ecommerce Text" dataset, our overall results validate the effec-
1477 tiveness of the SmartDS-Solver's design and its potential to significantly enhance automated data
1478 science workflows.

1479
1480 D.1.6 OPTIMAL TEMPERATURE VALUES ACROSS TASKS AND MODEL STAGES

1481 We observed the optimal temperatures identified across different tasks and models, with detailed
1482 values provided in Table A11.

1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Abbreviation	Competition Name	Task Description	Task Type
Historical Competitions			
Nomad18-TC	nomad2018-predict-transparent-conductors	Predict formation energy and band gap from material composition and structure to select films with both conductivity and transparency	Regression
Stanford-CV	stanford-covid-vaccine	Predict chemical reactivity and stability of each base from RNA sequences	Multi-output Regression
Vent-Press	ventilator-pressure-prediction	Real-time prediction of patient airway pressure from ventilator settings	Time Series Regression
TF-Speech	tensorflow-speech-recognition-challenge	Recognize 1-second audio clips into 30 different commands	Audio Classification
TextNorm-EN	text-normalization-challenge-english-language	Convert text with dates, currency symbols, etc. into spoken form	Sequence-to-Sequence
Jigsaw-Tox	jigsaw-toxic-comment-classification-challenge	Identify whether online comments contain six types of toxic speech	Multi-label Classification
PlantPath20	plant-pathology-2020-fgvc7	Identify healthy apple leaves or specific diseases from leaf photographs	Image Classification
Whale-Play	whale-categorization-playground	Identify same individual/species of whales from tail fin photographs	Image Retrieval / Metric Learning
NYC-Taxi	new-york-city-taxi-fare-prediction	Predict taxi fare based on origin-destination GPS coordinates and time	Regression
TPS-May22	tabular-playground-series-may-2022	Classify three categories using tabular data features	Multi-class Classification
TPS-Dec21	tabular-playground-series-dec-2021	Predict binary outcomes using tabular data features	Binary Classification
Ongoing Competitions (July 2025)			
PGS-S5E6	playground-series-s5e6	Predict optimal fertilizer types for crops based on soil and weather conditions. Alternative variations: Predict suitable fertilizers from crop/soil/weather data	Multi-class classification
UM-MCTS	um-game-playing-strength-of-mcts-variants	Predict the relative performance between two MCTS algorithm variants based on game features.	Regression
Private Tasks			
CLINIC	clinic-data	Predict cancer recurrence after treatment using tabular clinical data	Binary Classification

Table A2: Detailed descriptions of benchmark tasks used in evaluation

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583

Rule Type	Scoring Logic
Similarity	Feature (exact/stage/technique = 0.5/0.25/0.25), Algorithm (exact/paradigm/type/subtype = 0.5/0.25/0.15/0.1), Metric (exact/purpose/name = 0.5/0.25/0.25).
Format	<code><think>...</think><answer>...</answer></code> binary compliance check.
Executability	Presence of <code>python</code> block, >10 lines, containing <code>def/return</code> yields a small bonus (final score capped at 1.0).
Reasoning	Density of language-specific reasoning markers & step indicators (e.g., connectors, “if...then”, causal terms, enumerations, bullet points). <code><think></code> block weighted $\times 2$. Formula: $\min(1.0, 0.15 \cdot (\text{global} + 2 \cdot \text{think}))$, rounded to 2 decimals.

Table A3: Composite Reward Function – Additional Matching and Scoring Rules

1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612

Sub-reward	Extraction Rule Category	Representative Keywords Extracted
Main Reward: Feature Engineering	Cleaning / Generation / Selection & Preprocessing / Time-related / Text / Image	fillna, dropna, astype, query; groupby, rolling, shift; StandardScaler, PCA, SelectKBest; year, month, hour; TfidfVectorizer, word2vec; cv2, augmentation, resize
Main Reward: Algorithm Extraction	Supervised / Unsupervised / Time-series / Reinforcement Learning / Neural Networks / AutoML	LinearRegression, RandomForest, XGB*; KMeans, DBSCAN, PCA, TSNE, UMAP; ARIMA, Prophet; DQN, PPO; Conv2D, LSTM, Transformer; GridSearchCV, Optuna
Main Reward: Metric Extraction	Classification / Regression / Clustering / Cross-validation / Explainability / Domain-specific	accuracy, f1, roc_auc; mse, mae, r2; silhouette, Calinski–Harabasz, Davies–Bouldin; KFold, StratifiedKFold; feature_importances_, SHAP, LIME; BLEU, ROUGE, SSIM, PSNR, FID, WRMSSE

Table A4: Composite Reward Function – Extraction Rules and Representative Keywords

1613
1614
1615
1616
1617
1618
1619

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

		SmartDS-Solver																								
		Qwen-32B				Qwen-7B				Qwen-72B				Llama-8B				Llama-70B								
AIDE		Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	
Tasks	o1-preview																									
TF-Speech	0.312	0.642	0.697	0.574	0.687	0.611	0.316	0.341	0.422	0.710	0.627	0.306	0.341	0.643	0.574											
Nomad18-TC*	0.0715	0.075	0.075	0.071	0.068	0.0721	0.070	0.075	0.0719	0.070	0.081	0.071	0.075	0.0719	0.071											
PlantPath20	0.941	0.792	0.884	0.956	0.942	0.897	0.846	0.928	0.916	0.940	0.854	0.928	0.533	0.895	0.897											
Jigsaw-Tox	0.976	0.9819	0.974	0.970	0.975	0.973	0.972	0.9823	0.980	0.973	0.970	0.974	0.508	0.975	0.973											
TPS-Dec21	0.9493	0.950	0.9541	0.951	0.951	0.952	0.951	0.9543	0.9486	0.950	0.950	0.951	0.928	0.916	0.951											
TPS-May22	0.894	0.932	0.918	0.935	0.959	0.934	0.903	0.968	0.960	0.932	0.888	0.942	0.908	0.924	0.916											
Vent-Press*	2.823	2.147	1.045	1.191	0.635	3.825	4.062	0.770	1.294	0.807	3.694	6.340	2.217	0.681	1.853											
NYC-Taxi*	4.868	4.107	4.475	3.465	4.031	4.268	3.661	5.783	4.263	3.588	3.539	4.337	3.599	3.391	4.399											
Whale-Play	0.3279	0.124	0.404	0.404	0.375	0.3272	0.325	0.097	0.388	0.327	0.325	0.325	0.197	0.338	0.3277											
Stanford-CV*	0.480	-	0.640	0.663	0.430	-	-	-	-	0.473	-	-	-	0.556	0.438											
TextNorm-EN	0.978	0.980	0.930	0.938	0.971	0.933	0.980	0.930	0.936	0.981	0.919	0.980	0.925	0.982	0.932											
%win	63.6%	54.5%	54.5%	72.7%	81.8%	36.4%	54.5%	45.5%	54.5%	72.7%	27.3%	45.5%	36.4%	54.5%	63.6%											
#Fail	1	0	1	0	0	1	1	1	1	0	1	1	1	0	0											
#Rank-1		1	1	2	3	3	3	3	3	1	1	1	1	1	2											

Table A5: Performance comparison across MLE-benchmark tasks. Tasks use domain-appropriate metrics (higher-is-better: accuracy, AUC; lower-is-better: RMSE, MAE, loss). We report three summary statistics: %win (proportion of tasks outperforming the AIDE+o1-preview), Fail (number of tasks producing invalid outputs), and Rank-1 (number of tasks achieving the best performance). Boxed values indicate results better than the AIDE baseline; underlined values denote the best configuration for each task.

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

		SmartDS-Solver												AutoML-Agent							
AIDE		Qwen-32B				Qwen-7B				Qwen-72B				Llama-8B			Llama-70B			GPT-4o	
Tasks	o1-preview	Instruct	SFT	GRPOI	GRPO2	Instruct	SFT	GRPOI	GRPO2	SFT	Instruct	SFT	GRPOI	GRPO2	SFT	Instruct	SFT	GRPOI	GRPO2	SFT	
TF-Speech	0.312	0.642	0.697	0.574	0.687	0.611	0.316	0.341	0.422	0.710	0.627	0.306	0.341	0.643	0.574	0.627	0.306	0.341	0.643	0.574	0.146
Nomad18-TC*	0.0715	0.075	0.075	0.071	0.068	0.0721	0.070	0.075	0.0719	0.070	0.081	0.071	0.075	0.0719	0.071	0.081	0.071	0.075	0.0719	0.071	0.077
PlantPath20	0.941	0.792	0.884	0.956	0.942	0.897	0.846	0.928	0.916	0.940	0.854	0.928	0.533	0.895	0.897	0.854	0.928	0.533	0.895	0.897	0.519
Jigsaw-Tox	0.976	0.9819	0.974	0.970	0.975	0.973	0.972	0.9823	0.980	0.973	0.970	0.974	0.508	0.975	0.973	0.970	0.974	0.508	0.975	0.973	0.981
TPS-Dec21	0.9493	0.950	0.9541	0.951	0.951	0.952	0.951	0.9543	0.9486	0.950	0.950	0.951	0.928	0.916	0.951	0.950	0.951	0.928	0.916	0.951	-
TPS-May22	0.894	0.932	0.918	0.935	0.959	0.934	0.903	0.968	0.960	0.932	0.888	0.942	0.908	0.924	0.916	0.888	0.942	0.908	0.924	0.916	0.892
Vent-Press*	2.823	2.147	1.045	1.191	0.635	3.825	4.062	0.770	1.294	0.807	3.694	6.340	2.217	0.681	1.853	3.694	6.340	2.217	0.681	1.853	10.822
NYC-Taxi*	4.868	4.107	4.475	3.465	4.031	4.268	3.661	5.783	4.263	3.588	3.539	4.337	3.599	3.391	4.399	3.539	4.337	3.599	3.391	4.399	-
Whale-Play	0.3279	0.124	0.404	0.404	0.375	0.3272	0.325	0.097	0.388	0.327	0.325	0.325	0.197	0.338	0.3277	0.325	0.325	0.197	0.338	0.3277	0.001
Stanford-CV*	0.480	-	0.640	0.663	0.430	-	-	-	-	0.473	-	-	-	0.556	0.438	-	-	-	0.556	0.438	0.578
TextNorm-EN	0.978	0.980	0.930	0.938	0.971	0.933	0.980	0.930	0.936	0.981	0.919	0.980	0.925	0.982	0.932	0.919	0.980	0.925	0.982	0.932	-
%win	63.6%	63.6%	54.5%	72.7%	81.8%	36.4%	54.5%	45.5%	54.5%	72.7%	27.3%	45.5%	36.4%	54.5%	63.6%	27.3%	45.5%	36.4%	54.5%	63.6%	9%
#Fail	1	1	0	0	0	1	1	1	1	0	1	1	1	0	0	1	1	1	0	0	3
#Rank-1			1	2	3	1	1	3	1	1	1	1	1	2	2	1	1	1	2	2	

Table A6: Performance comparison across MLE-benchmark tasks. Tasks use domain-appropriate metrics (higher-is-better: accuracy, AUC; lower-is-better: RMSE, MAE, loss). We report three summary statistics: %win, #Fail, and #Rank-1. Boxed values indicate results better than the AIDE baseline; underlined values denote the best configuration for each task.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737

Tasks	SmartDS-Solver											
	AIDE											
	o1-preview	Qwen-32B			Qwen-7B			Qwen-72B	llama_8B			llama-70B
	Instruct	SFT	GRPO2	Instruct	SFT	GRPO2	SFT	Instruct	SFT	GRPO2	SFT	
PGS-SSE6	0.326	0.180	0.330	0.333	0.330	0.325	0.341	0.330	0.171	0.334	0.335	0.335
UM-MCTS*	-	0.547	0.569	0.468	-	-	0.463	0.456	-	0.540	0.458	0.544
CLINIC	0.830	0.851	0.856	0.873	0.836	0.831	0.876	0.875	0.834	0.853	0.875	0.854

Table A7: Performance on two ongoing Kaggle competitions active during our research period and one internal prediction modeling task using proprietary dataset. Notes: "-" indicates Fail.

1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759

SmartDS-Solver w/o Meta Learning Agent	
Dataset	Performance
TF-Speech	0.640
Plant Path20	0.830
Jigsaw-Tox	0.971
TPS-Dec21	0.951
TPS-May22	0.944
Whale-Play	0.3747
TextNorm-EN	0.923
PGS-SSE6	0.326
CLINIC	0.833
Nomad18-TC*	0.075
Vent-Press*	0.792
NYC-Taxi*	4.130
Stanford-CV*	0.568
UM-MCTS*	0.490

Table A8: Performance without meta-learning agent across all 14 benchmark tasks.

1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779

Tasks	AIDE+o1-preview				SmartDS-Solver			
	Coding		Summarization		Inference		Modification	
	Cnt	Tokens	Cnt	Tokens	Cnt	Tokens	Cnt	Tokens
Jigsaw-Tox	20	100579	20	36364	9	12354	23	43950
NYC-Taxi	20	145340	20	74390	4	3843	9	17212
Nomad18-TC	20	144415	20	74782	8	10469	22	45257
PlantPath20	20	144258	20	69041	9	11144	23	60352
Stanford-CV	20	191440	20	106559	9	12195	31	60871
TPS-Dec21	20	79883	20	36564	7	7047	16	21415
TPS-May22	20	108133	20	44556	4	3207	8	11050
TF-Speech	20	194569	20	113589	7	8101	17	30597
TextNorm-EN	20	131068	20	59906	8	12353	19	32662
Vent-Press	20	102519	20	47333	6	7170	16	25822
Whale-Play	20	137201	20	55803	7	8940	24	41537

Table A9: Token consumption comparison across 11 benchmark tasks.

1780
1781

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

Dataset	AutoML-Agent	SmartDS-Solver (qwen-32b-GRPO2)	AIDE (o1-preview)	AutoGluon	Human Models	SELA (MCTS)
Butterfly Image	0.926	0.937(+1.2%)	0.8566(-7.6%)	0.014	0.931	-
Shopee-IET	0.972	0.976(+0.4%)	0.977(+0.5%)	0.988	0.921	-
Cora	0.843	0.919(+9.0%)	-	-	0.811	-
Citeseer	0.632	0.777(+22.9%)	-	-	0.702	-
Click Prediction Small	0.352	0.520(+47.8%)	0.140(-60.1%)	-	-	0.238
Software Defects	0.573	0.836(+45.8%)	0.5531(-3.5%)	0.524	0.669	-
Smoker Status	0.762	0.763(+0.2%)	0.742(-2.6%)	-	-	0.785
Banana Quality	0.967	0.980(+1.3%)	0.9744(+0.8%)	0.980	0.976	-
MFeat Factors	0.940	0.977(+3.9%)	0.968(+3.0%)	-	-	0.957
Wine Quality White	0.652	0.967(+48.2%)	0.654 (+0.3%)	-	-	0.650
Smoker Status Higher Education	0.582	0.590(+1.3%)	-	-	0.513	-
Students Performance	0.769	0.887(+15.4%)	-	-	0.750	-
Textual Entailment	0.796	0.868(+9.0%)	-	0.807	0.664	-
Ecommerce Text	0.982	0.980(-0.2%)	0.9764(-0.6%)	0.987	0.935	-
Colleges	0.139	0.123(+11.5%)	-	-	-	0.142
Crab Age	0.161	0.130(+19.7%)	0.1619(-0.3%)	0.143	2.049	-
House Prices	10.111	0.137(+98.6%)	23908.9956(-236365.2%)	-	-	10.111
Crop Price	1.114	0.159(+85.7%)	1.0223(8.2%)	1.088	1.101	-

Table A10: Performance Comparison on data science tasks. The table shows the performance of our proposed SmartDS-Solver (using Qwen-32B-GRPO2) compared to other data science agents. The values in parentheses represent the percentage improvement of our method over AutoML-Agent, where positive values indicate improvement and negative values indicate degradation. The results of other methods (AutoGluon, Human Models, SELA) are cited from paper of AutoML-Agent.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

Tasks	SmartDS-Solver																					
	Qwen-32B				Qwen-7B				Qwen-72B				Llama-8B				Llama-70B					
	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	Instruct	SFT	GRPO1	GRPO2	SFT	
TF-Speech	0.58	0.66	0.52	0.59	0.60	0.60	0.60	0.61	0.55	0.46	0.48	0.60	0.55	0.44								
Nomad18-TC*	0.60	0.54	0.53	0.74	0.77	0.57	0.79	0.79	0.76	0.60	0.60	0.58	0.85	0.55								
PlantPath20	0.65	0.62	0.68	0.78	0.67	0.62	0.89	0.89	0.92	0.86	0.72	0.70	0.54	0.81								
Jigsaw-Tox	0.57	0.73	0.77	0.78	0.67	0.67	0.62	0.62	0.60	0.51	0.60	0.60	0.62	0.56								
TPS-Dec21	0.70	0.55	0.52	0.48	0.60	0.66	0.52	0.66	0.60	0.60	0.50	0.56	0.76	0.60								
TPS-May22	0.60	0.60	0.60	0.49	0.66	0.60	0.67	0.60	0.60	0.60	0.58	0.66	0.57	0.60								
Vent-Press*	0.84	0.48	0.50	0.62	0.69	0.60	0.66	0.60	0.80	0.64	0.58	0.73	0.66	0.74								
NYC-Taxi*	0.62	0.60	0.52	0.55	0.60	0.60	0.63	0.53	0.60	0.60	0.55	0.52	0.60	0.60								
Whale-Play	0.51	0.56	0.65	0.80	0.70	0.68	0.70	0.38	0.63	0.58	0.62	0.63	0.59	0.70								
Stanford-CV*	-	0.50	0.48	0.47	-	-	-	-	0.37	-	-	-	0.60	0.63								
TexNorm-EN	0.62	0.86	0.82	0.88	0.80	0.80	0.80	0.96	0.47	0.69	0.88	0.86	0.97	0.67								

Table A11: Optimal Temperature for different tasks of each model

D.1.7 STATISTICAL SIGNIFICANCE ANALYSIS

To further validate our findings beyond win rate comparisons, we conducted t-tests to assess statistical significance across our best-performing configurations. Table A12 presents the results comparing our top SmartDS-Solver configurations against the AIDE across all 14 benchmark tasks.

The analysis reveals statistically significant improvements ($p \leq 0.05$) for most configurations. Larger models with SFT-only training (Qwen-72B and Llama-70B) demonstrate particularly strong significance ($p \leq 0.04$), while models with complete GRPO2 training also show significant improvements. These results provide strong statistical evidence supporting the effectiveness of our approach.

Model Config.	t-value	p-value	Significance
Qwen-32B GRPO2	1.842	0.044	Significant*
Qwen-7B GRPO2	1.353	0.0995	Marginal†
Llama-8B GRPO2	1.777	0.049	Significant*
Qwen-72B SFT	1.902	0.0398	Significant*
Llama-70B SFT	1.978	0.0348	Significant*

Table A12: Statistical significance analysis compared to AIDE using paired t-tests across 14 tasks. * $p < 0.05$, † $p < 0.10$.

Task	Improvement Ratio
TF-Speech	3.69
Nomad18-TC*	0.13
PlantPath20	0.53
Jigsaw-Tox	0.01
TPS-Dec21	0.09
TPS-May22	0.03
Vent-Pres*	0.83
NYC-Taxi*	0.00
Whale-Play	2.66
Stanford-CV*	0.04
TextNorm-EN	0.40

Table A13: Task-wise improvement ratios comparing late-stage to early-stage performance of Qwen-32B across the 11 MLE-Bench tasks.

The improvement ratio is defined as $(\text{best} - \text{first success}) / \text{first success}$, where "first success" represents the score at the earliest successful iteration and "best" represents the highest score achieved across all iterations. A value of 0 indicates no gain from the meta mechanism. A one-sided test for $H_1: \text{improvement} > 0$ yields $p = 0.0173$, demonstrating significantly improves model accuracy.

E LLM ASSISTANCE IN MANUSCRIPT PREPARATION

In the preparation of this manuscript, we utilized Large Language Models (LLMs) as an assistive tool to enhance the quality and clarity of our writing. Our use of LLMs was confined to the following aspects:

- **Language and Prose Refinement:** We employed an LLM for proofreading, refining sentence structures, and polishing the prose. This was done to improve the overall readability and ensure our research contributions were communicated clearly and effectively.
- **Pseudo-code Generation:** To articulate our algorithms, we first described the logic and steps in natural language. An LLM was then prompted to translate these descriptions into formatted pseudo-code. All generated pseudo-code was meticulously reviewed and manually revised by the authors to ensure its correctness and complete fidelity to our proposed methods.
- **Clarity Verification:** To confirm that our explanations were unambiguous, we provided key paragraphs to an LLM and requested it to paraphrase or summarize the content. This

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

process served as a check to validate whether our intended meaning was conveyed without ambiguity, akin to having a digital reader provide feedback.

We wish to emphasize that all core ideas, experimental designs, data analysis, and conclusions presented in this paper are the original work of the authors. The LLM served strictly as a writing and formatting aid and was not involved in any part of the conceptual or analytical research process.