



# AUTOML-AGENT: A MULTI-AGENT LLM FRAMEWORK FOR FULL-PIPELINE AUTOML

Patara Trirat<sup>1\*</sup>, Wonyong Jeong<sup>2</sup>, Sung Ju Hwang<sup>1,2</sup>

<sup>1</sup>KAIST, <sup>2</sup>DeepAuto.ai

{patara.t, sungju.hwang}@kaist.ac.kr, young@deepauto.ai

## ABSTRACT

Automated machine learning (AutoML) accelerates AI development by automating tasks in the development pipeline, such as optimal model search and hyperparameter tuning. Existing AutoML systems often require technical expertise to set up complex tools, which is in general time-consuming and requires a large amount of human effort. Therefore, recent works have started exploiting large language models (LLM) to lessen such burden and increase the usability of AutoML frameworks via a natural language interface, allowing non-expert users to build their data-driven solutions. These methods, however, are usually designed only for a particular process in the AI development pipeline and do not efficiently use the inherent capacity of the LLMs. This paper proposes *AutoML-Agent*, a novel multi-agent framework tailored for full-pipeline AutoML, i.e., from data retrieval to model deployment. *AutoML-Agent* takes user’s task descriptions, facilitates collaboration between specialized LLM agents, and delivers deployment-ready models. Unlike existing work, instead of devising a single plan, we introduce a retrieval-augmented planning strategy to enhance exploration to search for more optimal plans. We also decompose each plan into sub-tasks (e.g., data preprocessing and neural network design) each of which is solved by a specialized agent we build via prompting executing in parallel, making the search process more efficient. Moreover, we propose a multi-stage verification to verify executed results and guide the code generation LLM in implementing successful solutions. Extensive experiments on seven downstream tasks using fourteen datasets show that *AutoML-Agent* achieves a higher success rate in automating the full AutoML process, yielding systems with good performance throughout the diverse domains.

## 1 INTRODUCTION

Automated machine learning (AutoML) has significantly reduced the need for technical expertise and human labors in developing effective data-driven solutions by automating each process in the AI development pipeline (Yao et al., 2018; Ren et al., 2020; He et al., 2021), such as feature engineering, model selection, and hyperparameter optimization (HPO). However, current AutoML systems (Gijsbers et al., 2024) often necessitate programming expertise to configure complex tools and resources, potentially creating barriers for a larger pool of users with limited skills and knowledge.

To make AutoML frameworks more accessible, recent studies (Trirat et al., 2021; Viswanathan et al., 2023; Li et al., 2023; Hollmann et al., 2023b; Liu et al., 2024a; Zhang et al., 2023; Shen et al., 2023; Zhang et al., 2024a; Hong et al., 2024a; Guo et al., 2024a; Yang et al., 2024) have suggested to use natural language interfaces with large language models (LLM) for machine learning (ML) and data science (DS) tasks. Nevertheless, these previous LLM-based AutoML frameworks only considered a limited number of tasks due to their restricted designs, either only for a process in the pipeline (e.g., feature engineering (Hollmann et al., 2023b; Li et al., 2024; Malberg et al., 2024), HPO (Liu et al., 2024a; Zhang et al., 2024a), and model selection (Zhang et al., 2023; Shen et al., 2023)) or for a specific group of downstream tasks (e.g., natural language processing (Viswanathan et al., 2023) and computer vision (Yang et al., 2024)). In addition, most methods overlook the inherent capability

\*Work done during an internship at DeepAuto.ai.

of LLMs to search for promising models by performing actual training of the candidate models during the search process, making it prohibitively costly and slow.

For an AutoML framework to be truly practical, it should perform end-to-end AutoML, considering both the **data aspects** (retrieval, preprocessing, and feature engineering) and **model aspects** (selection, HPO, and deployment). This is because a process in one aspect can affect subsequent processes in the other, potentially leading to suboptimal solutions when combining results from different frameworks. Meanwhile, the AutoML framework should be computationally efficient, using strategies to minimize the computational overhead during search. However, there are two main challenges in building such a framework.

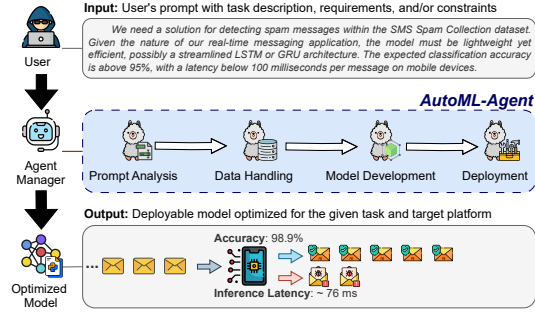


Figure 1: *AutoML-Agent* receives user’s instructions and delivers optimized deployable models.

**High Complexity of the Planning Tasks** The planning of the entire AutoML pipeline introduces additional complexities compared to task- or problem-specific planning, primarily due to the interdependencies among the steps in the pipeline. For example, types of retrieved datasets affects how to design preprocessing steps and neural networks. Then, the designed network affects which particular hyper-parameters need to be optimized depending on the given downstream task. Such inter-step dependencies result in the enlarged search space since it should consider all possible combinations of inter-related steps. Furthermore, enabling the framework to operate across various downstream tasks exacerbates these challenges, as each has task-specific requirements.

**Challenges in Accurate Implementations** To develop a modular and extendable framework that effectively handles diverse ML tasks, it is crucial to enhance the flexibility of the LLM agent in its code generation ability, such as by decoupling the template code from the code for specific datasets. However, using LLMs to autonomously generate complete ML pipelines may lead to hallucination issues, including code incompleteness, incorrect or missing dependencies, and potential undiscovered bugs (Hong et al., 2024b). Furthermore, LLMs often struggle with code generation when prompted with ambiguous task descriptions. Thus, we need accurate analysis of the requirements, and a code-generation platform that can adaptively generate code based on disambiguated requirements.

To address the above challenges, we propose a novel multi-agent framework, *AutoML-Agent*, for full-pipeline AutoML from data and model search to evaluation, with strategies to tackle the complexity of the planning problem as well as accurate implementation of code. As illustrated in Figure 1, *AutoML-Agent* accepts a user’s task description and coordinates multiple specialized agents to collaboratively identify and execute an optimal ML pipeline, ultimately delivering a deployment-ready model and its inference endpoint as the output.

Specifically, to tackle the complex planning problem, we introduce a new *retrieval-augmented planning* strategy equipped with role-specific decomposition and prompting-based execution. This strategy produces multiple plans based on retrieved knowledge for a given task description, facilitating the exploration of promising plans. Moreover, it enables LLM agents to discern global (pipeline-level) and local (process-level) relationships among steps through plan decomposition, which helps them focus on their immediate sub-tasks while aligning with the user’s goal. The retrieval-augmented component also simplifies extending LLMs to various downstream tasks using relevant knowledge. The prompting-based execution enhances search efficiency by exploiting LLMs’ in-context learning capabilities without any further training, which could introduce additional cost. To enhance the accuracy of the implementation, we adopt structure-based prompt parsing that extracts ML-relevant requirements from the user’s description and *multi-stage verification* that provides feedback between each step in the framework to ensure the quality of instructions when guiding the LLM for code generation. These modules aim to improve the correctness and clarity of the task description for code implementation. Our main **contributions** are as follows.

- We propose a novel multi-agent LLM framework for AutoML, designed to automate the entire AI development pipeline. To the best of our knowledge, this is the first attempt to employ LLMs in a task-agnostic AutoML framework that spans from data retrieval to model deployment.

Table 1: Comparison between *AutoML-Agent* and existing LLM-based frameworks.

| Framework                               | Key Functionality |              |               |               |                      |                |
|---|-------------------|--------------|---------------|---------------|----------------------|----------------|
|   | Planning          | Verification | Full Pipeline | Task-Agnostic | Training-Free Search | With Retrieval |
| AutoML-GPT (Zhang et al., 2023)         | ×                 | ×            | ×             | ✓             | ✓                    | ×              |
| Prompt2Model (Viswanathan et al., 2023) | ×                 | ×            | ✓             | ×             | ×                    | ✓              |
| HuggingGPT (Shen et al., 2023)          | ✓                 | ×            | ×             | ✓             | ✓                    | ✓              |
| MLCopilot (Zhang et al., 2024a)         | ×                 | ×            | ×             | ✓             | ✓                    | ×              |
| AgentHPO (Liu et al., 2024a)            | ✓                 | ✓            | ×             | ✓             | ×                    | ×              |
| AutoMMLab (Yang et al., 2024)           | ×                 | ✓            | ✓             | ×             | ×                    | ×              |
| CAAFE (Hollmann et al., 2023b)          | ×                 | ✓            | ×             | ×             | ×                    | ×              |
| Data Interpreter (Hong et al., 2024a)   | ✓                 | ✓            | ×             | ✓             | ×                    | ×              |
| DS-Agent (Guo et al., 2024a)            | ✓                 | ✓            | ×             | ✓             | ×                    | ✓              |
| <i>AutoML-Agent</i> (Ours)              | ✓                 | ✓            | ✓             | ✓             | ✓                    | ✓              |

- We address the challenges due to the complexity of the planning problem in full-pipeline AutoML by introducing retrieval-augmented planning with role-specific plan decomposition and prompting-based plan execution, enhancing the flexibility and efficiency of the search process.
- To enhance the accuracy of our full-pipeline implementation, we integrate structure-based prompt parsing and multi-stage verification to ensure the quality of resulting solutions and instructions prior to actual code implementation, thereby improving overall performance.
- We demonstrate the superiority of the proposed *AutoML-Agent* framework through extensive experiments on seven downstream tasks using fourteen datasets across five application domains.

## 2 RELATED WORK

Automated machine learning (AutoML) is a transformative approach for optimizing ML workflows, enabling both practitioners and researchers to efficiently design models and preprocessing pipelines with minimal manual intervention (Ren et al., 2020; He et al., 2021; Gijssbers et al., 2024). Despite several advancements in AutoML (Jin et al., 2019; Feurer et al., 2022; Tang et al., 2024), most of them are designed only for particular elements of the ML pipeline. Only a few works (Bisong, 2019; Mukunthu et al., 2019; Microsoft, 2021) support multiple steps of the pipeline. Also, due to the traditional programming interfaces, these systems often have complex configuration procedures and steep learning curves that require substantial coding expertise and an understanding of the underlying ML concepts, limiting their accessibility to non-experts and being time-consuming even for experienced users. These limitations hinder the widespread adoption of traditional AutoML systems.

Large language models (LLM), e.g., GPT-4 (Achiam et al., 2023) and LLaMA (Touvron et al., 2023), have recently shown promise in addressing these limitations with the complex problem-solving skills across disciplines via human-friendly language interfaces, including AI problems (Xi et al., 2023). This shift towards natural language-driven interfaces democratizes access and allows users to articulate their needs in a more intuitive manner. However, existing LLM-based frameworks can only assist in a specific step of the ML pipeline, such as feature engineering (Hollmann et al., 2023b), model search (Shen et al., 2023; Hong et al., 2024a; Guo et al., 2024a), or HPO (Liu et al., 2024a; Zhang et al., 2024a). A few attempts (Viswanathan et al., 2023; Yang et al., 2024) support the entire ML production pipeline, yet only for a specific type of downstream tasks. Besides, these methods either naively use the LLMs or overlook the inherent capabilities, making their search processes time-consuming for the AutoML pipeline that requires sophisticated planning and verification.

In contrast to the existing studies, our framework aims to overcome these limitations by incorporating a new retrieval-augmented planning strategy, coupled with plan decomposition and prompting-based execution techniques, alongside structure-based prompt parsing and multi-stage verification. Through these enhancements, we can increase plan execution efficiency and support diverse ML tasks with more accurate pipeline implementation. Table 1 summarizes the key differences between *AutoML-Agent* and existing frameworks.

## 3 A MULTI-AGENT LLM FRAMEWORK FOR FULL-PIPELINE AUTOML

This section presents details of the proposed multi-agent framework, *AutoML-Agent*, including agent specifications, a prompt parsing module, a retrieval-augmented planning strategy, a prompting-based plan execution, and a multi-stage verification. As depicted in Figure 2, all agents are coordinated by an Agent Manager to complete the user’s instructions by delivering the deployment-ready model.

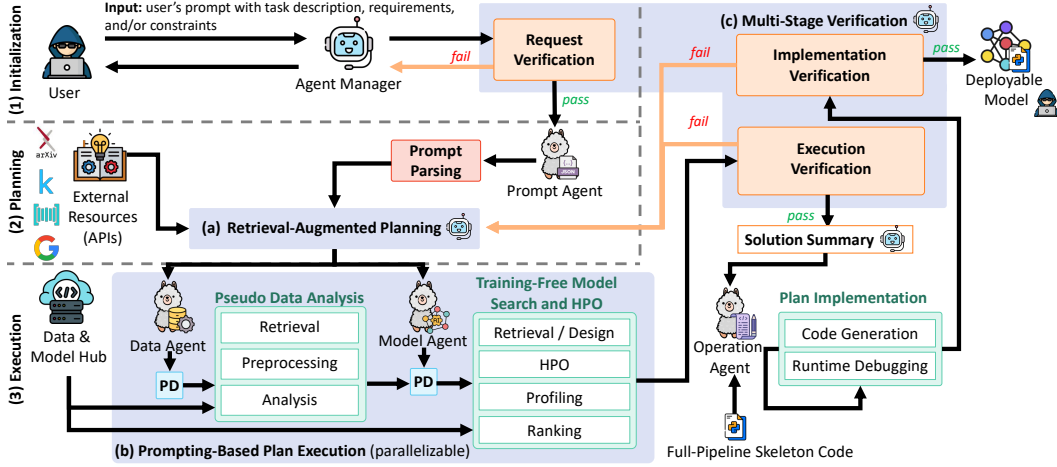


Figure 2: Overview of our *AutoML-Agent* framework. **(1) Initialization** stage aims to receive a valid user instruction using request verification. **(2) Planning** stage focuses on extracting ML related information by parsing the user instruction into a standardized form, and uses it to devise plans accordingly. **(3) Execution** stage executes each action given by the devised plans. Finally, based on the best execution results, *AutoML-Agent* outputs codes containing deployable model to the user.

### 3.1 AGENT SPECIFICATIONS

We now provide brief descriptions of the agents in our multi-agent AutoML framework.

**Agent Manager** ( $\mathcal{A}_{mgr}$ ) acts as the core interface between users and other LLM agents in the framework. It is responsible for interacting with the user, devising a set of global plans for subsequent processes with retrieved knowledge, distributing tasks to corresponding agents, verifying executed results with feedback, and tracking the system progress.

**Prompt Agent** ( $\mathcal{A}_p$ ) is an LLM specifically instruction-tuned for parsing the user’s instructions into a standardized JSON object with predefined keys. The parsed information is then shared across agents in the framework during the planning, searching, and verification phases.

**Data Agent** ( $\mathcal{A}_d$ ) is an LLM prompted for doing tasks related to data manipulation and analysis. The analysis results from the Data Agent are used to inform the Model Agent about data characteristics during the model search and HPO.

**Model Agent** ( $\mathcal{A}_m$ ) is an LLM prompted for doing tasks related to model search, HPO, model profiling, and candidate ranking. The results produced by the Model Agent are sent back to the Agent Manager for verification before proceeding to the Operation Agent.

**Operation Agent** ( $\mathcal{A}_o$ ) is an LLM prompted for implementing the solution found by the Data and Model Agents that passes the Agent Manager’s verification. The Operation Agent is responsible for writing effective code for actual runtime execution and recording the execution results for final verification before returning the model to the user.

After we define all agents with their corresponding profiles as described above (see §B.1 for detailed prompts), the  $\mathcal{A}_{mgr}$  then assigns relevant tasks to each agent according to the user’s input. Note that we can implement  $\mathcal{A}_d$  and  $\mathcal{A}_m$  with more than one agent per task based on the degree of parallelism.

### 3.2 FRAMEWORK OVERVIEW

We present an overview of our *AutoML-Agent* framework in Figure 2 and Algorithm 1. In the **(1) initialization** stage, the Agent Manager ( $\mathcal{A}_{mgr}$ ) receives the user instruction and checks its validity through request verification (Figure 2(c) and Line 3). In the **(2) planning** stage, the Prompt Agent ( $\mathcal{A}_p$ ) parses the verified user instruction into a standardized JSON object. Then,  $\mathcal{A}_{mgr}$  generates plans to solve the given AutoML task using retrieval-augmented planning (Figure 2(a) and Line 11). In the **(3) execution** stage, the Data ( $\mathcal{A}_d$ ) and Model ( $\mathcal{A}_m$ ) Agents decompose these plans and execute them via plan decomposition (PD) and prompting-based plan execution (Figure 2(b)

**Algorithm 1** Overall Procedure of *AutoML-Agent*


---

**Initialization:** Agent Manager  $\mathcal{A}_{mgr}$ , instruction-tuned Prompt Agent  $\mathcal{A}_p$ , Data Agent  $\mathcal{A}_d$ , Model Agent  $\mathcal{A}_m$ , Operation Agent  $\mathcal{A}_o$ , deployment-ready model  $\mathcal{M}^*$ , and system state  $S$

**Input:** User instruction  $I$

```

1: while  $S \neq \text{END}$  and  $\mathcal{M}^* = \emptyset$  do
2:   if  $S = \text{INIT}$  then
3:      $F \leftarrow \mathcal{A}_{mgr}(\text{ReqVer}(I))$  ▷ run request verification (§3.6) for feedback  $F$ 
4:     if  $F = \emptyset$  then ▷ check if  $I$  is valid
5:        $R \leftarrow \mathcal{A}_p(I)$  ▷ parse user instruction  $I$  (§3.3)
6:        $S \leftarrow \text{PLAN}$ 
7:     else
8:       return  $F$  ▷ return feedback  $F$  to the user.
9:     end if
10:  else if  $S = \text{PLAN}$  then
11:     $\mathbf{P} \leftarrow \mathcal{A}_{mgr}(\text{RAP}(R))$  ▷ run retrieval-augmented planning (§3.4)
12:    for  $\mathbf{p}_i$  in  $\mathbf{P}$  do
13:       $s_i^d \leftarrow \text{PD}(R, \mathcal{A}_d, \mathbf{p}_i)$  ▷ run plan decomposition for Data Agents (§3.5)
14:       $O_i^d \leftarrow \mathcal{A}_d(s_i^d)$  ▷ run pseudo data analysis (§3.5)
15:       $s_i^m \leftarrow \text{PD}(R, \mathcal{A}_m, \mathbf{p}_i, O_i^d)$  ▷ run plan decomposition for Model Agents (§3.5)
16:       $O_i^m \leftarrow \mathcal{A}_m(s_i^m)$  ▷ run training-free model search and HPO (§3.5)
17:    end for
18:     $\mathbf{O} \leftarrow \{(O_i^d, O_i^m)\}_{i=1}^P$  ▷ aggregate execution outcomes from all agents (§3.5)
19:    if  $\mathcal{A}_{mgr}(\text{ExecVer}(\mathbf{O}))$  is pass then ▷ run execution verification (§3.6)
20:       $I^* \leftarrow \mathcal{A}_{mgr}(\mathbf{O})$  ▷ find the best plan and create corresponding instructions
21:       $\mathcal{M}^* \leftarrow \mathcal{A}_o(I^*)$  ▷ run code generation for the best plan
22:      if  $\mathcal{A}_{mgr}(\text{ImpVer}(\mathcal{M}^*))$  is pass then ▷ run implementation verification (§3.6)
23:         $S \leftarrow \text{END}$  ▷ stop the process
24:      end if
25:    end if
26:  end if
27: end while
28: return  $\mathcal{M}^*$ 

```

---

and Line 13–16), whose results are then verified against the user’s requirements via execution verification (Figure 2(c) and Line 19). Finally,  $\mathcal{A}_{mgr}$  selects the best plan and sends it to the Operation Agent ( $\mathcal{A}_o$ ) to write code (Line 21). After code generation, implementation verification (Figure 2(c) and Line 22) is conducted to ensure that the code is deployment-ready. If any of the verification steps fail, *AutoML-Agent* performs revision steps (orange lines in Figure 2) to generate new solutions. In the following subsections, we provide the descriptions of each step more in detail.

### 3.3 INSTRUCTION DATA GENERATION AND PROMPT PARSING

**Data Generation** For  $\mathcal{A}_p$  to generate accurate JSON objects, we need to instruction-tune the LLM first because it can output a valid JSON object but with incorrect keys that are irrelevant to subsequent processes. Following Xu et al. (2024), we first manually create a set of high-quality seed instructions then automatically generate a larger instruction dataset  $D = \{(I_i, R_i)\}_{i=1}^N$ , having  $N$  instruction-response pairs. Here,  $I_i$  is the  $i$ -th instruction with the corresponding response  $R_i$ . We use the JSON format substantially extended from Yang et al. (2024) for response  $R_i$  with the following top-level keys to extract the user’s requirement from various aspects of an AutoML pipeline.

- **User.** The user key represents the user intention (e.g., build, consult, or unclear) of the given instruction and their technical expertise in AI.
- **Problem.** The problem key indicates the characteristics and requirements of the given task, including area (e.g., computer vision), downstream task (e.g., image classification), application or business domain, and other constraints like expected accuracy and inference latency.
- **Dataset.** The dataset key captures the data characteristics and properties, including data modality, requested preprocessing and augmentation techniques, and potential data source.

- **Model.** The model key captures the expected model characteristics and properties, including model name (e.g., ViT), family (e.g., Transformer), and type (e.g., neural networks).
- **Knowledge.** The knowledge key extracts additional knowledge or insights helpful for solving the given problem directly provided by the user, potentially associated with the expertise level.
- **Service.** The service key is relevant to the downstream implementation and deployment. It provides information such as a target device and an inference engine.

**Prompt Parsing** Then, we can use the generated dataset  $D$  to train an LLM and use it as  $\mathcal{A}_p$ . Note that these standardized keys are important for a better control over the LLM agents’ behavior within our framework and necessary for effective communication between agents. Moreover, these keys provide contextual information for generating a high-quality AutoML pipeline from various perspectives. After the instruction tuning, we use the  $\mathcal{A}_p$  to parse the user’s instructions (or task descriptions) and return the parsed requirements  $R = \mathcal{A}_p(I)$  to  $\mathcal{A}_{mgr}$ , as shown in §C.1.

### 3.4 RETRIEVAL-AUGMENTED PLANNING

Recent studies (Guo et al., 2024b; Huang et al., 2024; Masterman et al., 2024; Zhang et al., 2024b; Hu et al., 2024) highlights that effective planning and tool utilization are essential for solving complex problems with LLMs, especially in a multi-agent framework. By bridging two techniques in a single module, we propose a retrieval-augmented planning (RAP) strategy to effectively devise a robust and up-to-date set of diverse plans for the AutoML problems.

Let  $\mathbf{P} = \{\mathbf{p}_1, \dots, \mathbf{p}_P\}$  be a set of plans. Based on past knowledge embedded in the LLM, knowledge retrieved via external APIs (such as arXiv papers), and  $R$ , RAP generates  $P$  multiple end-to-end plans for the entire AutoML pipeline having different scenario  $\mathbf{p}_i$ . This strategy enables *AutoML-Agent* to be aware of newer and better solutions. Specifically, we first use the parsed requirements  $R$  to acquire a summary of the relevant knowledge and insights via API calls, including web search and paper summary.  $\mathcal{A}_{mgr}$  then uses this information to devise  $P$  different plans, i.e.,  $\mathbf{P} = \mathcal{A}_{mgr}(\text{RAP}(R))$ . Note that  $\mathcal{A}_{mgr}$  devises each plan independently to make the subsequent steps parallelizable. The benefit of this strategy is that it enhances exploration for better solutions while allowing parallelization. Examples of generated plans are provided in §C.2.

### 3.5 PROMPTING-BASED PLAN EXECUTION AND IMPLEMENTATION

Given the generated  $\mathbf{P}$ , we now describe how  $\mathcal{A}_d$  and  $\mathcal{A}_m$  execute each  $\mathbf{p}_i$  using prompting techniques without actually executing the code. Examples of the execution results are in §C.4.

**Plan Decomposition** Due to the high complexity of the end-to-end plan, we first need to adaptively decompose the original plan  $\mathbf{p}_i \in \mathbf{P}$  into a smaller set of sub-tasks  $\mathbf{s}_i$  relevant to the agent’s roles and expertise to increase the effectiveness of LLMs in solving and executing the given plan (Khot et al., 2023). The plan decomposition (PD) process involves querying the agents about their understanding of the given plan specific to their roles. Formally,  $\mathbf{s}_i^d = \text{PD}(R, \mathcal{A}_d, \mathbf{p}_i)$ , where  $\mathbf{s}_i^d$  is the *decomposed* plan for Data Agent, containing sub-tasks for the given plan  $\mathbf{p}_i$ . Then, the agent executes the decomposed plan towards the user’s requirements instead of the original lengthy plan. We define the sub-tasks  $\mathbf{s}_i^m$  of  $\mathcal{A}_m$  below due to its reliance on Data Agent’s outcomes. Examples of decomposed plans are in §C.3.

**Pseudo Data Analysis** In *AutoML-Agent*,  $\mathcal{A}_d$  handles sub-tasks in  $\mathbf{s}_i^d$ , including the retrieval, pre-processing, augmentation, and analysis of the specified dataset. During the data retrieval phase, if the dataset is not directly supplied by the user, we initiate an API call to search for potential datasets in repositories, such as HuggingFace and Kaggle, using the dataset name or description. Upon locating a dataset, we augment the prompt with metadata from the dataset’s source; if no dataset is found, we rely on the inherent knowledge of the LLM. We then prompt  $\mathcal{A}_d$  to proceed by acting *as if* it actually executes  $\mathbf{s}_i^d$ , according to the dataset characteristics and user requirements  $R$ . The summarized outcomes of these sub-tasks,  $O_i^d$ , are then forwarded to the  $\mathcal{A}_m$ . Hence,  $O_i^d = \mathcal{A}_d(\mathbf{s}_i^d)$ .

**Training-Free Model Search and HPO** Like  $\mathcal{A}_d$ ,  $\mathcal{A}_m$  uses API calls to complete all sub-tasks  $\mathbf{s}_i^m$ , instead of direct code execution. However, in contrast to  $\mathcal{A}_d$ , the plan decomposition for  $\mathcal{A}_m$



incorporates outcomes from the  $\mathcal{A}_d$ , enabling it to recognize characteristics of the preprocessed dataset, i.e.,  $s_i^m = \text{PD}(R, \mathcal{A}_m, \mathbf{p}_i, O_i^d)$ . Here, the  $\mathcal{A}_m$ 's prompt is enhanced with insights gathered by  $\mathcal{A}_{mgr}$  about high-performing models and relevant hyperparameters for the specific ML task. This technique allows the  $\mathcal{A}_m$  to execute the sub-tasks in  $s_i^m$  more efficiently. Using this augmented prompt, the  $\mathcal{A}_m$  follows a similar procedure to  $\mathcal{A}_d$ , undertaking model retrieval, running HPO, and summarizing the results of these sub-tasks, which include expected numerical performance metrics such as accuracy and error, as well as model complexity factors like model size and inference time. To facilitate the subsequent verification step, we also prompt the agent to return results with the top- $k$  most promising models. Formally,  $O_i^m = \mathcal{A}_m(s_i^m)$ .

**Plan Implementation** To enhance the efficacy of  $\mathcal{A}_o$  in code generation,  $\mathcal{A}_{mgr}$  first verifies all executed results  $\mathbf{O} = \{(O_i^d, O_i^m)\}_{i=1}^P$  from  $\mathcal{A}_d$  and  $\mathcal{A}_m$ .  $\mathcal{A}_{mgr}$  then selects the best outcome  $O^* \in \mathbf{O}$  and generates the instruction  $I^*$  for  $\mathcal{A}_o$  to write the actual code accordingly. Formally,  $\mathcal{M}^* = \mathcal{A}_o(I^*)$ , where  $\mathcal{M}^*$  is the deployment-ready model.

### 3.6 MULTI-STAGE VERIFICATION

Verification, especially with refinement or feedback, is essential for maintaining the correct trajectory of LLMs (Baek et al., 2024; Madaan et al., 2023; Gou et al., 2024). Our framework incorporates three verification steps to guarantee its accuracy and effectiveness: request verification, execution verification, and implementation verification.

**Request Verification** Initially, we assess the clarity of the user's instructions to determine if they are relevant and adequate for executing ML tasks and addressing the user's objectives. If the instructions prove insufficient for progressing to the planning stage,  $\mathcal{A}_{mgr}$  will request additional information, facilitating multi-turn communication. This request verification (ReqVer in Algorithm 1 Line 3) step, however, often overlooked in existing studies, placing an undue burden on users to formulate a more detailed initial prompt—a challenging task particularly for those who are non-experts or lack experience. Prompts for ReqVer are shown in §B.4.1.

**Execution Verification** After executing the plans in §3.5,  $\mathcal{A}_{mgr}$  then verifies whether any of the pipelines produced by  $\mathcal{A}_d$  and  $\mathcal{A}_m$  (i.e.,  $\mathbf{O}$ ) satisfy the user's requirements via prompting (see §B.4.2). If the results are satisfied, the suggested solution is selected as a candidate for implementation. This execution verification (ExecVer) step effectively mitigates computational overhead in the search process by allocating resources exclusively to the most promising solution.

**Implementation Verification** This implementation verification (ImpVer) phase closely resembles the execution verification; however, it differs in that it involves validating outcomes derived from the code that has been executed and compiled by  $\mathcal{A}_o$ . We present the prompt for this verification in §B.4.3. If the outcomes meet the user's requirements,  $\mathcal{A}_{mgr}$  provides the model and deployment endpoint to the user.

Note that if any execution or implementation fails to satisfy the user requirements (i.e., does not pass the verification process), these failures are systematically documented. Subsequently, the system transitions to the plan *revision* stage. During this stage,  $\mathcal{A}_{mgr}$  formulates a revised set of plans, incorporating insights derived from the outcomes of the unsuccessful plans.

## 4 EXPERIMENTS

We validate the effectiveness of our full-pipeline AutoML framework by comparing *AutoML-Agent* with handcrafted models, state-of-the-art AutoML variants, and LLM-based frameworks across multiple downstream tasks involving different data modalities.

### 4.1 SETUP

**Downstream Tasks and Datasets** As summarized in Table 2, we select seven downstream tasks from five different data modalities, including image, text, tabular, graph, and time series. These

datasets are chosen from different sources. Also, we incorporate various evaluation metrics for these tasks, e.g., accuracy for classification and RMSLE for regression.

For each task, we prepare *two* sets of natural language task descriptions to represent *constraint-aware* and *constraint-free* requirements (see §A) along with a full-pipeline skeleton script. As a result, we extensively evaluate **28** generated models. Note that this setting differs from previous studies (Guo et al., 2024a; Huang et al., 2023), which require dataset-specific, partially completed code preparation.

Table 2: Summary of downstream tasks and datasets.

| Data Modality                      | Downstream Task         | Dataset Name                         | Evaluation Metric |
|------------------------------------|-------------------------|--------------------------------------|-------------------|
| Image (Computer Vision)            | Image Classification    | Butterfly Image<br>Shopee-IET        | Accuracy          |
| Text (NLP)                         | Text Classification     | Ecommerce Text<br>Textual Entailment | Accuracy          |
| Tabular (Classic ML)               | Tabular Classification  | Banana Quality<br>Software Defects   | F1                |
|                                    | Tabular Regression      | Crab Age<br>Crop Price               | RMSLE             |
|                                    | Tabular Clustering      | Smoker Status<br>Student Performance | RI                |
| Time Series (Time Series Analysis) | Time-Series Forecasting | Weather<br>Electricity               | RMSLE             |
| Graph (Graph Mining)               | Node Classification     | Cora<br>Citeseer                     | Accuracy          |

**Evaluation Metrics** For a comprehensive evaluation, we measure the agent’s effectiveness in both code generation and task-specific performance aspects by using *comprehensive score* (CS) (Hong et al., 2024a) to simultaneously evaluate both the success rate (SR) of code generation and the normalized performance score (NPS) of the built pipelines. That is,  $CS = 0.5 \times SR + 0.5 \times NPS$ . Here,  $NPS = \frac{1}{1+s}$  is a transformation of loss-based performance score  $s$ , e.g., RMSLE. More detailed explanations are included in §A.4.

As described above, we evaluate all frameworks under two different settings. To measure SR of each method, we use a grading scale ranging from 0 for total failure to 1 for perfect conformity to the user’s requirements. For the *constraint-free* setting, a method can get a score of 0.5 (pass modeling) or 1.0 (pass deployment). For the *constraint-aware* setting, a method can get a score of 0.25 (pass modeling), 0.5 (pass deployment), 0.75 (partially pass the constraints), or 1.0 (pass all cases).

**Baselines** As we propose a framework for the novel task of full-pipeline AutoML with LLMs, there is no direct baseline available for comparison. We thus compare *AutoML-Agent* against the task-specific manually designed models (see §A.3): **Human Models**, the variants of state-of-the-art AutoML: **AutoGluon** (Erickson et al., 2020; Shchur et al., 2023; Tang et al., 2024), a state-of-the-art LLM for data science: **DS-Agent** (Guo et al., 2024a), and general-purpose LLMs: **GPT-3.5** (Brown et al., 2020) and **GPT-4** (Achiam et al., 2023) with zero-shot prompting.

**Implementation Details** Except for the  $\mathcal{A}_p$  that is implemented with Mixtral-8x7B (*Mixtral-8x7B-Instruct-v0.1*) (Jiang et al., 2024), we use GPT-4 (*gpt-4o-2024-05-13*) as the backbone model for all agents and LLM-based baselines to ensure an impartial performance evaluation. To instruction tune the  $\mathcal{A}_p$  (§3.3), we automatically generate about 2.3K instruction-response pairs using EvolInstruct (Xu et al., 2024). Here, we use LoRA (Hu et al., 2021) to fine-tune the model with the generated dataset. For RAP (§3.4), we set the number of plans  $P = 3$  and the number of candidate models  $k = 3$ . All experiments are conducted on an Ubuntu 22.04 LTS server equipped with eight NVIDIA A100 GPUs (CUDA 12.4) and Intel(R) Xeon(R) Platinum 8275CL CPU @ 3.00GHz. For running the generated models, we employ the same execution environment as DS-Agent (Guo et al., 2024a), with all necessary libraries included in the skeleton scripts.

## 4.2 MAIN RESULTS

We report the average scores from five independent runs for all evaluation metrics in Figure 3.

**Success Rate** Figure 3(left) and Table 5 present the results for the SR metric. For the constraint-free cases, which can be considered easier tasks, all methods have higher SR than ones in the constraint-aware setting. Notably, *AutoML-Agent* also consistently outperforms the baselines in the constraint-aware setting, achieving an average SR of 87.1%, which underscores the effectiveness of the proposed framework. We conjecture that the knowledge retrieved during the planning process helps the agents identify which areas to focus on in order to meet the given constraints. Regarding DS-Agent, although we use the provided example cases for the relevant tasks, DS-Agent appears to fail on certain tasks due to its heavy reliance on curated case banks and the inclusion of partially completed code, which is unavailable in our setting.



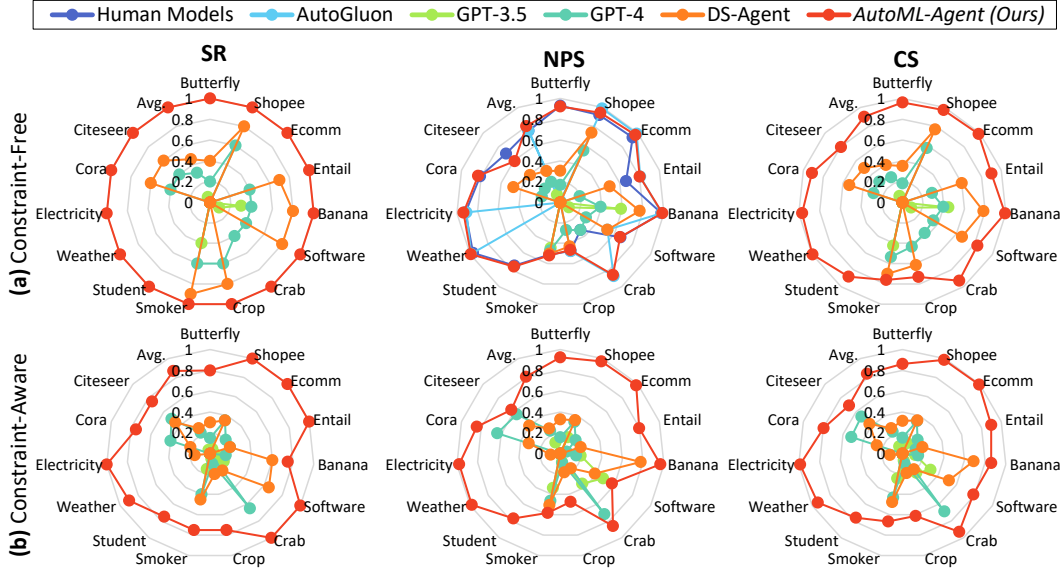


Figure 3: Performance comparison across all datasets using the SR, NPS, and CS metrics under (a) constraint-free and (b) constraint-aware settings. Higher scores indicate better results.

**Downstream Performance** We present the performance comparison between the successfully built models in Figure 3(center) and Table 6. To ensure meaningful results and to examine how the performance of LLM-generated models compares to state-of-the-art AutoML techniques and manual ML pipelines crafted by experienced experts, we select top-performing models by evaluating results reported in Papers with Code benchmarks and Kaggle notebooks for the same tasks and datasets, where applicable, as the Human Models baselines. From the results, we can observe that *AutoML-Agent* significantly outperforms other agents, including Human Models, in the NPS metric. In particular, *AutoML-Agent* achieves the best performance across all tasks under the constraint-aware setting. These findings highlight the superiority of *AutoML-Agent* in adapting to various scenarios, attributed to the retrieval-augmented planning (RAP) strategy. This approach enables agents to discover effective pipelines for given constraints. These empirical observations substantiate the efficacy of the proposed RAP, providing up-to-date solutions for various tasks.

**Comprehensive Score** Figure 3(right) and Table 7 present the weighted quality of each agent based on the CS metric. Overall, *AutoML-Agent* outperforms all other baselines, especially in more complicated tasks. Interestingly, it is evident that general-purpose LLMs still work relatively well on classical tasks like tabular classification and regression, while more sophisticated methods, such as DS-Agent and our *AutoML-Agent* work significantly better in complex tasks. This finding aligns with previous research (Guo et al., 2024a), which suggests that tabular tasks typically involve straightforward function calls from the sklearn library (Pedregosa et al., 2011), and therefore do not demand advanced reasoning or coding abilities from LLM agents, unlike more complex tasks.

#### 4.3 RESOURCE COST

As we primarily use closed-source LLMs in this paper, we analyze the resource costs in terms of time and money. Figure 4 presents the average time and monetary costs across different tasks and datasets for a single run, under the constraint-free (upper) and constraint-aware (lower) settings. On average, it takes around 525 seconds and costs 0.30 USD (using GPT-4o) to search for a single model that will be deployable after training. The significant amount of time spent in the planning stage also suggests the difficulty in devising plans for full-pipeline AutoML.

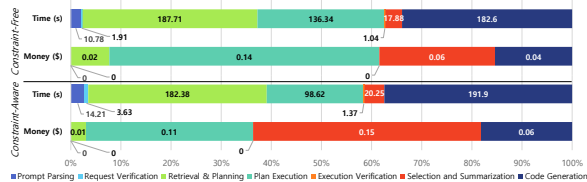


Figure 4: Average time and monetary cost breakdown.

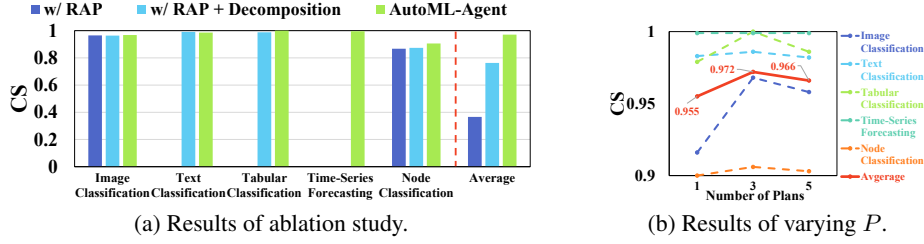


Figure 5: Results of (a) ablation and (b) hyperparameter studies in the CS metric.

#### 4.4 ABLATION AND HYPERPARAMETER STUDIES

To validate the effectiveness of each component in *AutoML-Agent*, we conduct the following ablation studies. The results are presented in Figure 5a and Table 8. First, we investigate *retrieval-augmented planning (RAP)* alone, where retrieved knowledge from external APIs is directly used without plan decomposition and multi-stage verification. As expected, this ablation leads to a decline in performance, and in some cases, even fails to generate a runnable model. This outcome highlights the importance of the decomposition and verification modules. Second, we evaluate *RAP with plan decomposition*, where the plan is decomposed for each specific agent. While this variant demonstrates better downstream performance, it still fails to produce runnable models in certain downstream tasks due to the lack of code verification. Finally, we assess the *full framework with multi-stage verification*, which provides feedback to agents, thereby enhancing both their planning and coding capabilities. Integrating all components significantly empowers LLM agents to effectively incorporate external knowledge from various sources to build a full-pipeline AutoML system.

To further verify the effectiveness of devising multiple plans in our retrieval-augmented planning strategy (§3.4), we conduct a hyperparameter study by varying the number of plans  $P$  in the constraint-free setting. As shown in Figure 5b and Table 9, the number of plans does not significantly affect the success rate, likely due to GPT-4’s robust planning capabilities. However, based on the NPS and CS metrics, we observe that the number of plans has a notable impact on downstream task performance. Also, these results also suggest that adding more plans does not necessarily lead to better results, as the model may generate multiple similar plans, resulting in similar outcomes. Consequently, we select 3 as the default number of plans.

## 5 CONCLUSION

This paper presents *AutoML-Agent*, a novel LLM-based multi-agent framework designed for AutoML, covering the entire pipeline from data retrieval to model deployment. *AutoML-Agent* tackles the full-pipeline planning complexity and implementation accuracy challenges in the LLMs for task-agnostic AutoML by leveraging the newly proposed retrieval-augmented planning strategy and multi-stage verification. In addition, we enhance the plan execution efficiency by integrating role-specific decomposition and prompting-based execution techniques into the framework. Our experiments on seven ML tasks demonstrate that *AutoML-Agent* outperforms existing methods in terms of success rate and downstream task performance.

**Limitations and Future Work** Even though we offer a flexible module to accommodate various ML tasks and data modalities, the absence of a skeleton code for completely new tasks could increase the risk of code hallucination problems. Additionally, in the current version, there is still a gap in code generation quality when using different backbones, e.g., GPT-4 vs. GPT-3.5, which is not unique to our approach but a common challenge faced by existing LLM-based frameworks. Developing a more robust framework that can effectively provide reasonable solutions with less reliance on the LLM backbone is very promising future work.

In addition, our work still faces code generation failures when applied to machine learning tasks that require significantly different development pipelines from those tested in our experiments, which focused on general supervised and unsupervised settings. Tasks such as reinforcement learning and recommendation systems pose particular challenges. Consequently, extending *AutoML-Agent* to these tasks will require the development of additional agents to handle specific steps in the target pipeline, such as actor-environment interaction and reward modeling in reinforcement learning.

## REPRODUCIBILITY STATEMENT

We present the complete prompts and showcase results in §B and §D to facilitate reproduction. More experimental and implementation details are provided in §A, along with detailed results in §E.

## ETHICS STATEMENT

We expect *AutoML-Agent* to offer significant advantages by promoting AI-driven innovation and enabling individuals with limited AI expertise to effectively utilize AI capabilities. However, we acknowledge the potential misuse of *AutoML-Agent* by malicious users, such as generating offensive content, malicious software, or invasive surveillance tools when exposed to harmful inputs. This vulnerability is not unique to our approach but represents a common challenge faced by existing LLMs with substantial creative and reasoning capabilities, which can occasionally produce undesirable outputs.

Although we strictly instruct the LLM to focus on generating positive results for machine learning tasks, there is a possibility of unforeseen glitches that could introduce security issues within the system. Therefore, we recommend running *AutoML-Agent* within a Docker container to ensure isolation from the host’s file system. Additionally, due to its integration with external services for retrieval-augmented generation and API-based LLMs like GPT-4, privacy concerns may arise. Users should carefully review any data included in API prompts to prevent unintended data disclosures.

## REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. Researchagent: Iterative research idea generation over scientific literature with large language models. *arXiv preprint arXiv:2404.07738*, 2024.
- Ekaba Bisong. *Building machine learning and deep learning models on Google cloud platform*. Springer, 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Advances in neural information processing systems*, pp. 1877–1901, 2020.
- Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander Smola. Autogluon-tabular: Robust and accurate automl for structured data. *arXiv preprint arXiv:2003.06505*, 2020.
- Matthias Feurer, Katharina Eggensperger, Stefan Falkner, Marius Lindauer, and Frank Hutter. Auto-sklearn 2.0: Hands-free automl via meta-learning. *Journal of Machine Learning Research*, 23(261):1–61, 2022.
- Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- Pieter Gijsbers, Marcos LP Bueno, Stefan Coors, Erin LeDell, Sébastien Poirier, Janek Thomas, Bernd Bischl, and Joaquin Vanschoren. Amlb: an automl benchmark. *Journal of Machine Learning Research*, 25(101):1–65, 2024.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujiu Yang, Nan Duan, and Weizhu Chen. CRITIC: Large language models can self-correct with tool-interactive critiquing. In *The Twelfth International Conference on Learning Representations*, 2024.
- Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. Ds-agent: Automated data science by empowering large language models with case-based reasoning. In *International Conference on Machine Learning*, 2024a.

- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*, 2024b.
- Xin He, Kaiyong Zhao, and Xiaowen Chu. Auttml: A survey of the state-of-the-art. *Knowledge-Based Systems*, 212:106622, 2021.
- Noah Hollmann, Samuel Müller, Katharina Eggersperger, and Frank Hutter. TabPFN: A transformer that solves small tabular classification problems in a second. In *ICLR*, 2023a.
- Noah Hollmann, Samuel Müller, and Frank Hutter. Large language models for automated data science: Introducing caafe for context-aware automated feature engineering. In *Advances in Neural Information Processing Systems*, volume 36, pp. 44753–44775, 2023b.
- Sirui Hong, Yizhang Lin, Bangbang Liu, Binhao Wu, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Lingyao Zhang, Mingchen Zhuge, et al. Data interpreter: An llm agent for data science. *arXiv preprint arXiv:2402.18679*, 2024a.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024b.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Mengkang Hu, Yao Mu, Xinmiao Chelsey Yu, Mingyu Ding, Shiguang Wu, Wenqi Shao, Qiguang Chen, Bin Wang, Yu Qiao, and Ping Luo. Tree-Planner: Efficient close-loop task planning with large language models. In *ICLR*, 2024.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Benchmarking large language models as ai research agents. *arXiv preprint arXiv:2310.03302*, 2023.
- Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*, 2024.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Haifeng Jin, Qingquan Song, and Xia Hu. Auto-keras: An efficient neural architecture search system. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1946–1956, 2019.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*, 2023.
- Dawei Li, Zhen Tan, and Huan Liu. Exploring large language models for feature selection: A data-centric perspective. *arXiv preprint arXiv:2408.12025*, 2024.
- Haoyuan Li, Hao Jiang, Tianke Zhang, Zhelun Yu, Aoxiong Yin, Hao Cheng, Siming Fu, Yuhao Zhang, and Wanggui He. Traineragent: Customizable and efficient model training through llm-powered multi-agent system. *arXiv preprint arXiv:2311.06622*, 2023.
- Siyi Liu, Chen Gao, and Yong Li. Large language model agent for hyper-parameter optimization. *arXiv preprint arXiv:2402.01881*, 2024a.
- Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. In *ICLR*, 2024b.

- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Simon Malberg, Edoardo Mosca, and Georg Groh. FELIX: Automatic and interpretable feature engineering using llms. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 230–246. Springer, 2024.
- Tula Masterman, Sandi Besen, Mason Sawtell, and Alex Chao. The landscape of emerging ai agent architectures for reasoning, planning, and tool calling: A survey. *arXiv preprint arXiv:2404.11584*, 2024.
- Microsoft. Neural Network Intelligence, 1 2021. URL <https://github.com/microsoft/nni>.
- D. Mukunthu, P. Shah, and W.H. Tok. *Practical Automated Machine Learning on Azure: Using Azure Machine Learning to Quickly Build AI Solutions*. O’Reilly Media, 2019. ISBN 9781492055549.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12(85): 2825–2830, 2011.
- Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Xiaojiang Chen, and Xin Wang. A comprehensive survey of neural architecture search: Challenges and solutions. *arXiv preprint arXiv:2006.02903*, 2020.
- Oleksandr Shchur, Ali Caner Turkmen, Nick Erickson, Huibin Shen, Alexander Shirkov, Tony Hu, and Bernie Wang. Autogluon-timeseries: Automl for probabilistic time series forecasting. In *International Conference on Automated Machine Learning*. PMLR, 2023.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. In *Advances in Neural Information Processing Systems*, volume 36, pp. 38154–38180, 2023.
- Zhiqiang Tang, Haoyang Fang, Su Zhou, Taojiannan Yang, Zihan Zhong, Tony Hu, Katrin Kirchhoff, and George Karypis. Autogluon-multimodal (automm): Supercharging multimodal automl with foundation models. *arXiv preprint arXiv:2404.16233*, 2024.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Patara Trirat, Yooju Shin, Sejin Kim, and Minseok Kim. Generating a machine learning model with a few sentences. *Korea Software Congress*, pp. 688–690, 2021.
- Vijay Viswanathan, Chenyang Zhao, Amanda Bertsch, Tongshuang Wu, and Graham Neubig. Prompt2model: Generating deployable models from natural language instructions. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 413–421, 2023.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. WizardLM: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*, 2024.

- Chenxiao Yang, Qitian Wu, Jiahua Wang, and Junchi Yan. Graph neural networks are inherently good generalizers: Insights by bridging GNNs and MLPs. In *ICLR*, 2023.
- Zekang Yang, Wang Zeng, Sheng Jin, Chen Qian, Ping Luo, and Wentao Liu. Autommlab: Automatically generating deployable models from language instructions for computer vision tasks. *arXiv preprint arXiv:2402.15351*, 2024.
- Quanming Yao, Mengshuo Wang, Yuqiang Chen, Wenyuan Dai, Yu-Feng Li, Wei-Wei Tu, Qiang Yang, and Yang Yu. Taking human out of learning applications: A survey on automated machine learning. *arXiv preprint arXiv:1810.13306*, 2018.
- Lei Zhang, Yuge Zhang, Kan Ren, Dongsheng Li, and Yuqing Yang. Mlcpilot: Unleashing the power of large language models in solving machine learning tasks. In *The 18th Conference of the European Chapter of the Association for Computational Linguistics*, 2024a.
- Shujian Zhang, Chengyue Gong, Lemeng Wu, Xingchao Liu, and Mingyuan Zhou. Automl-gpt: Automatic machine learning with gpt. *arXiv preprint arXiv:2305.02499*, 2023.
- Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weiming Lu. Self-Contrast: Better reflection through inconsistent solving perspectives. In *ACL*, pp. 3602–3622, 2024b.

## A DETAILS OF EXPERIMENTAL SETUP

This section outlines the detailed experimental setup used in this paper, including the complete instruction prompts for both constraint-free (Table 3) and constraint-aware (Table 4) settings, a full-pipeline skeleton script (§A.1), dataset and baseline descriptions, as well as evaluation metrics.

Table 3: User instruction (i.e., task description) for experiments under the *constraint-free* setting.

| Task                    | Dataset                               | Instruction Prompt   |
|-------------------------|---------------------------------------|--|
| Image Classification    | Butterfly Image                       | I need a very accurate model to classify images in the Butterfly Image Classification dataset into their respective categories. The dataset has been uploaded with its label information in the labels.csv file.   |
|                         | Shopee-IET                            | Please provide a classification model that categorizes images into one of four clothing categories. The image path, along with its label information, can be found in the files train_labels.csv and test_labels.csv.  |
| Text Classification     | Ecommerce Text                        | We need a state-of-the-art model for text classification based on the Ecommerce Text dataset. The model should be capable of accurately classifying text into four categories: Electronics, Household, Books, and Clothing & Accessories. We have uploaded the entire dataset without splitting it here.   |
|                         | Textual Entailment                    | You are solving this machine learning tasks of classification: The dataset presented here (the Textual Entailment) comprises a series of labeled text pairs. Given two texts (text1 and text2), your task is to predict the relationship of the text pair of neutral (0), contradiction (1) or entailment (2). The evaluation metric is accuracy. Build a language model to get a good performance.  |
| Tabular Classification  | Banana Quality                        | Build a model to classify banana quality as Good or Bad based on their numerical information about bananas of different quality (size, weight, sweetness, softness, harvest time, ripeness, and acidity). We have uploaded the entire dataset for you here in the banana_quality.csv file.   |
|                         | Software Defects                      | You are solving this data science tasks of binary classification: The dataset presented here (the Software Defects Dataset) comprises a lot of numerical features. Please split the dataset into three parts of train, valid and test. Your task is to predict the defects item, which is a binary label with 0 and 1. The evaluation metric is the F1 score. Please train a binary classification model to get a good performance on this task.   |
| Tabular Regression      | Crab Age                              | You are solving this data science tasks of regression: The dataset presented here (the Crab Age Dataset) comprises a lot of both categorical and numerical features. Please split the dataset into three parts of train, valid and test. Your task is to predict the age item. The evaluation metric is the RMSLE (root mean squared log error). Now train a regression model to get a good performance on this task.  |
|                         | Crop Price                            | I need a regression model to predict crop prices based on features like soil composition, environmental factors, historical yield data, and crop management practices from the dataset I uploaded here.  |
| Tabular Clustering      | Smoker Status                         | You are solving this data science tasks of unsupervised clustering: The dataset presented here (the Smoker Status Dataset) comprises a lot of numerical features. Please use the features in the test.csv file. Your task is to create the clustered items, which is a binary label with 0 and 1 (two clusters). The evaluation metric is the Rand index or Rand score, can be tested against 'smoking' labels. Now train an unsupervised clustering model to get a good performance on this task.   |
|                         | Higher Education Students Performance | I want an unsupervised clustering model to group student performances into eight groups. The dataset named 'Higher Education Students Performance Evaluation' (id=856) can be downloaded via ucimlrepo library. The clustering quality can be check against target variable OUTPUT Grade.  |
| Time-Series Forecasting | Weather                               | I want you to create a model for node classification on the Cora dataset to predict the category of each paper. You need to directly find the Cora dataset from a relevant library.  |
|                         | Electricity                           | I want you to create a model for node classification on the Citeseer dataset to predict the category of each paper. You need to directly find the Citeseer dataset from a relevant library.  |
| Node Classification     | Cora                                  | Build a model to perform time-series forecasting using the Weather dataset uploaded here, evaluating its accuracy with the RMSLE metric. Note that the input is a sequence of past observations with fixed size (INPUT_SEQ_LEN=96, INPUT_DIM=21). The model should predict the next future sequence with a fixed size (PRED_SEQ_LEN=96, PRED_DIM=21).  |
|                         | Citeseer                              | You are solving this machine learning tasks of time series forecasting: The dataset presented here (the Electricity dataset) comprises real-world time series data. Please split the dataset into three parts of train, valid and test. The input is a sequence of past observation with fixed size (INPUT_SEQ_LEN=96, INPUT_DIM=321). Your task is to predict the next future sequence with fixed size (PRED_SEQ_LEN=96, PRED_DIM=321). The evaluation metric is root mean squared log error (RMSLE). Now train a time series forecasting model to get a good performance on the given fixed sequences. |



Table 4: User instruction (i.e., task description) for experiments under the *constraint-aware* setting. **Bold** texts indicate constraints used for evaluation.

| Task                    | Dataset                               | Instruction Prompt  |
|-------------------------|---------------------------------------|---|
| Image Classification    | Butterfly Image                       | I need a highly accurate machine learning model developed to classify images within the Butterfly Image Classification dataset into their correct species categories. The dataset has been uploaded with its label information in the labels.csv file. Please use a convolutional neural network (CNN) architecture for this task, leveraging transfer learning from a <b>pre-trained ResNet-50 model</b> to improve accuracy. Optimize the model using cross-validation on the training split to fine-tune hyperparameters, and aim for an <b>accuracy of at least 0.95 on the test split</b> . Provide the final trained model, a detailed report of the training process, hyperparameter settings, accuracy metrics, and a confusion matrix to evaluate performance across different categories. |
|                         | Shopee-IET                            | Please provide a classification model that categorizes images into one of four clothing categories. The image path, along with its label information, can be found in the files train.labels.csv and test.labels.csv. The model should <b>achieve at least 85% accuracy on the test set</b> and be implemented using PyTorch. Additionally, please include data augmentation techniques and a confusion matrix in the evaluation.   |
| Text Classification     | Ecommerce Text                        | We require the development of an advanced <b>neural network model</b> for text classification tailored to the Ecommerce Text dataset, with the objective of <b>achieving at least 0.95 classification accuracy</b> . The model should be specifically trained to distinguish text into four defined categories: Electronics, Household, Books, and Clothing & Accessories. To facilitate this, we have uploaded the complete dataset in its entirety, without any prior division into training, validation, or test sets.   |
|                         | Textual Entailment                    | You are solving this machine learning task of classification: The dataset presented here (the Textual Entailment) comprises a series of labeled text pairs. Given two texts, your task is to predict the relationship of the text pair as neutral (0), contradiction (1), or entailment (2). The evaluation metric is accuracy. Build a language model to get good performance, ensuring the <b>model size does not exceed 200 million parameters</b> and the <b>inference time is less than 200 milliseconds per prediction</b> .  |
| Tabular Classification  | Banana Quality                        | Build a machine learning model, potentially <b>XGBoost or LightGBM</b> , to classify banana quality as Good or Bad based on their numerical information about bananas of different quality (size, weight, sweetness, softness, harvest time, ripeness, and acidity). We have uploaded the entire dataset for you here in the banana_quality.csv file. The model must <b>achieve at least 0.98 accuracy</b> .  |
|                         | Software Defects                      | You are solving this data science task of binary classification: The dataset presented here (the Software Defects Dataset) comprises a lot of numerical features. Please split the dataset into three parts of train, valid, and test. Your task is to predict the defects item, which is a binary label with 0 and 1. The evaluation metric is the F1 score. Please train a binary classification model to get a good performance on this task, ensuring that the model <b>training time does not exceed 30 minutes</b> and the <b>prediction time for each instance is under 5 milliseconds</b> .   |
| Tabular Regression      | Crab Age                              | You are solving this data science task of regression: The dataset presented here (the Crab Age Dataset) comprises a lot of both categorical and numerical features. Please split the dataset into three parts of train, valid, and test. Your task is to predict the age item. The evaluation metric is the RMSLE (root mean squared log error). Now train a regression model to get a good performance on this task, ensuring that the model's <b>training time does not exceed 30 minutes</b> and that it can make <b>predictions on the test set within 5 seconds</b> .  |
|                         | Crop Price                            | I need an accurate regression model to predict crop prices based on features like soil composition, environmental factors, historical yield data, and crop management practices from the dataset I uploaded here. You should optimize the model to achieve <b>RMSLE less than 1.0</b>   |
| Tabular Clustering      | Smoker Status                         | You are solving this data science task of unsupervised clustering: The dataset presented here (the Smoker Status Dataset) comprises a lot of numerical features. Please use the features in test.csv. Your task is to create the clustered items, which is a binary label with 0 and 1 (two clusters). The evaluation metric is the Rand index or Rand score, which can be tested against 'smoking' labels. Now train an unsupervised clustering model to get a good performance on this task, ensuring that the <b>Rand index is at least 0.75</b> and the model <b>training time does not exceed 10 minutes</b> .   |
|                         | Higher Education Students Performance | I want an unsupervised clustering model to group student performances into eight groups. The dataset named 'Higher Education Students Performance Evaluation' (id=856) can be downloaded via uci/mlrepo library. The clustering quality can be checked against the target variable OUTPUT Grade. The model should achieve a <b>Rand Score of at least 0.8</b> and <b>complete clustering within 10 minutes</b> .  |
| Time-Series Forecasting | Weather                               | Build a state-of-the-art time-series forecasting model for the Weather dataset uploaded here, evaluating its accuracy with the RMSLE metric. Note that the input is a sequence of past observations with fixed size (INPUT_SEQ_LEN=96, INPUT_DIM=21). The model should predict the next future sequence with a fixed size (PRED_SEQ_LEN=96, PRED_DIM=21). We target <b>RMSLE lower than 0.05</b> .  |
|                         | Electricity                           | You are solving this machine learning task of time series forecasting: The dataset presented here (the Electricity dataset) comprises real-world time series data. Please split the dataset into three parts of train, valid, and test. The input is a sequence of past observation with fixed size (INPUT_SEQ_LEN=96, INPUT_DIM=321). Your task is to predict the next future sequence with fixed size (PRED_SEQ_LEN=96, PRED_DIM=321). The evaluation metric is root mean squared log error (RMSLE). Now train a time series forecasting model to get a good performance on the given fixed sequences. Ensure the model achieves an <b>RMSLE of less than 0.1</b> and that the <b>training time does not exceed 1 hour</b> on a GPU.  |
| Node Classification     | Cora                                  | I want you to develop a node classification model using the <b>Graph Convolutional Network (GCN)</b> algorithm to predict the category of each paper in the Cora dataset. Start by importing the Cora dataset using the 'Planetoid' dataset from the 'torch_geometric.datasets' module in PyTorch Geometric. Ensure you preprocess the data to include node features and labels correctly. Train the model using a suitable optimizer and loss function. Then, evaluate its accuracy on the test set. The <b>accuracy on the test set should be over 0.90</b> .   |
|                         | Citeseer                              | I want you to develop a node classification model using the <b>Graph Convolutional Network (GCN)</b> algorithm to predict the category of each paper in the Citeseer dataset. Start by importing the Citeseer dataset using the 'Planetoid' dataset from the 'torch_geometric.datasets' module in PyTorch Geometric. Ensure you preprocess the data to include node features and labels correctly. Train the model using a suitable optimizer and loss function. Then, evaluate its accuracy on the test set. The <b>accuracy on the test set should be over 0.80</b> .   |

## A.1 SKELETON PYTHON SCRIPT

## Skeleton Python Script (e.g., text\_classification.py)

```

# The following code is for "text classification" task using PyTorch.
import os, random, time, json

# define GPU location
os.environ["CUDA_DEVICE_ORDER"] = "PCI_BUS_ID"
os.environ["CUDA_VISIBLE_DEVICES"] = "3"

import torch
import torch.nn as nn
import torch.optim as optim
import numpy as np
import gradio as gr

# TODO: import other required library here, including libraries for datasets and (pre-
#       trained) models like HuggingFace and Kaggle APIs. If the required module is not found,
#       you can directly install it by running 'pip install your_module'.
from torchtext import datasets, data, vocab
from torch.utils.data import DataLoader, Dataset
from sklearn.metrics import accuracy_score, f1_score

SEED = 42
random.seed(SEED)
torch.manual_seed(SEED)
np.random.seed(SEED)

# Define device for model operations
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

DATASET_PATH = "_experiments/datasets" # path for saving and loading dataset(s) (or the
# user's uploaded dataset) for preprocessing, training, hyperparamter tuning, deployment
# , and evaluation

# Data preprocessing and feature engineering
def preprocess_data():
    # TODO: this function is for data preprocessing and feature engineering

    # Run data preprocessing

    # Should return the preprocessed data
    return processed_data

def train_model(model, train_loader):
    # TODO: this function is for model training loop and optimization on 'train' and 'valid'
    # datasets
    # TODO: this function is for fine-tuning a given pretrained model (if applicable)

    # Should return the well-trained or finetuned model.
    return model

def evaluate_model(model, test_loader):
    # In this task, we use Accuracy and F1 metrics to evaluate the text classification
    # performance.
    # The 'performance_scores' should be in dictionary format having metric names as the
    # dictionary keys
    # TODO: the first part of this function is for evaluating a trained or fine-tuned model
    # on the 'test' dataset with respect to the relevant downstream task's performance
    # metrics
    # Define the 'y_true' for ground truth and 'y_pred' for the predicted classes here.

    performance_scores = {
        'ACC': accuracy_score(y_true, y_pred),
        'F1': f1_score(y_true, y_pred)
    }

    # TODO: the second part of this function is for measuring a trained model complexity on
    # a samples with respect to the relevant complexity metrics, such as inference time
    # and model size
    # The 'complexity_scores' should be in dictionary format having metric names as the
    # dictionary keys

    # Should return model's performance scores
    return performance_scores, complexity_scores

def prepare_model_for_deployment():

```

```

# TODO: this function is for preparing an evaluated model using model compression and
# conversion to deploy the model on a particular platform

# Should return the deployment-ready model
return deployable_model

def deploy_model():
    # TODO: this function is for deploying an evaluated model with the Gradio Python library

    # Should return the url endpoint generated by the Gradio library
    return url_endpoint

# The main function to orchestrate the data loading, data preprocessing, feature
# engineering, model training, model preparation, model deployment, and model evaluation
def main():
    """
    Main function to execute the text classification pipeline.
    """

    # TODO: Step 1. Retrieve or load a dataset from hub (if available) or user's local
    # storage (if given)
    dataset = None

    # TODO: Step 2. Create a train-valid-test split of the data by splitting the 'dataset'
    # into train_loader, valid_loader, and test_loader.
    # Here, the train_loader contains 70% of the 'dataset', the valid_loader contains 20% of
    # the 'dataset', and the test_loader contains 10% of the 'dataset'.
    train_loader, valid_loader, test_loader = (None, None, None) # corresponding to 70%,
    20%, 10% of 'dataset'

    # TODO: Step 3. With the split dataset, run data preprocessing and feature engineering (
    # if applicable) using the "preprocess_data" function you defined
    processed_data = preprocess_data()

    # TODO: Step 4. Define required model. You may retrieve model from available hub or
    # library along with pretrained weights (if any).
    # If pretrained or predefined model is not available, please create the model according
    # to the given user's requirements below using PyTorch and relevant libraries.
    model = None

    # TODO: Step 5. train the retrieved/loaded model using the defined "train_model"
    # function
    # TODO: on top of the model training, please run hyperparameter optimization based on the
    # suggested hyperparameters and their values before proceeding to the evaluation step
    # to ensure model's optimality

    model = train_model()

    # TODO: evaluate the trained model using the defined "evaluate_model" function
    model_performance, model_complexity = evaluate_model()

    # TODO: compress and convert the trained model according to a given deployment platform
    # using the defined "prepare_model_for_deployment" function
    deployable_model = prepare_model_for_deployment()

    # TODO: deploy the model using the defined "deploy_model" function
    url_endpoint = deploy_model()

    return processed_data, model, deployable_model, url_endpoint, model_performance,
    model_complexity

if __name__ == "__main__":
    processed_data, model, deployable_model, url_endpoint, model_performance,
    model_complexity = main()
    print("Model Performance on Test Set:", model_performance)
    print("Model Complexity:", model_complexity)

```

## A.2 DATASET DESCRIPTIONS

As presented in Table 2, we select seven representative downstream tasks, covering five data modalities. We describe the datasets their statistics as follows.

- **Butterfly Image (Butterfly).** This dataset includes 75 distinct classes of butterflies, featuring over 1,000 labeled images, including validation images. Each image is assigned to a single

- butterfly category. The dataset is accessible at <https://www.kaggle.com/datasets/phuchthaiv02/butterfly-image-classification>.
- **Shopee-IET (Shopee).** This dataset is designed for cloth image classification, where each image represents a clothing item, and its corresponding label indicates the clothing category. The available labels include BabyPants, BabyShirt, womencasualshoes, and womenchiffontop. The dataset is available at <https://www.kaggle.com/competitions/demo-shopee-iet-competition/data>.
  - **Ecommerce Text (Ecomm).** This dataset is a classification-based E-commerce text dataset comprising four categories: Electronics, Household, Books, and Clothing & Accessories, which together cover approximately 80% of any E-commerce website. It includes 50,425 instances and can be found at <https://www.kaggle.com/datasets/saurabhshahane/ecommerce-text-classification>.
  - **Textual Entailment (Entail).** This dataset consists of labeled pairs of text, where the task is to predict the relationship between each pair as either neutral (0), contradiction (1), or entailment (2). It is divided into a training set containing 4,907 samples and a testing set with 4,908 samples. We use the dataset provided by Guo et al. (2024a).
  - **Banana Quality (Banana).** This tabular dataset consists of numerical information on 8,000 samples of bananas, covering various quality attributes such as size, weight, sweetness, softness, harvest time, ripeness, acidity, and overall quality. The primary objective of the dataset is to classify each banana sample as either good or bad. The dataset is available at <https://www.kaggle.com/datasets/l31lff/banana/data>.
  - **Software Defects (Software).** This dataset consists primarily of numerical features and has been divided into three parts: training, validation, and testing. The goal is to predict a binary defect label (0 or 1). The training set contains 82,428 samples, the validation set contains 9,158 samples, and the test set contains 91,587 samples. We use the dataset provided by Guo et al. (2024a).
  - **Crab Age (Crab).** This dataset contains a mix of categorical and numerical features, and has been divided into three parts: training, validation, and test sets. The task is to predict the age of the crabs. The training set consists of 59,981 samples, the validation set includes 6,664 samples, and the test set contains 66,646 samples. We use the dataset provided by Guo et al. (2024a).
  - **Crop Price (Crop).** This new dataset contains 2,200 samples with key features such as nitrogen, phosphorus, and potassium ratios in the soil, temperature (in °C), humidity (in %), soil pH value, and rainfall (in mm), all of which are essential for predicting crop yield values. Crop yield prediction is crucial in modern agriculture, particularly as data-driven methods become more prevalent. This dataset is available at <https://www.kaggle.com/datasets/varshitanalluri/crop-price-prediction-dataset>.
  - **Smoker Status (Smoker).** This dataset contains numerous numerical features. The goal is to categorize smoking status of each instance into a cluster. The training set consists of 143,330 samples and the test set includes 143,331 samples. We use the dataset provided by Guo et al. (2024a).
  - **Higher Education Students Performance (Student).** The dataset, collected in 2019 from students in the Faculty of Engineering and Faculty of Educational Sciences, was created to predict students' end-of-term performances using machine learning techniques. It is a multivariate dataset with 145 instances and 31 integer-type features, focusing on classification tasks within the domain of social sciences. We adopt this dataset for unsupervised clustering instead of classification. This dataset can be found at <https://archive.ics.uci.edu/dataset/856/higher+education+students+performance+evaluation>.
  - **Weather.** The weather dataset consists of 21 meteorological factors collected every 10 minutes from the Weather Station at the Max Planck Biogeochemistry Institute in 2020, containing 52,603 samples without any pre-splitting. It is accessible at <https://github.com/thuml/Time-Series-Library>.
  - **Electricity.** This dataset comprises hourly electricity consumption data for 321 customers collected from 2012 to 2014, totaling 26,211 samples. The dataset records the electricity usage of these clients on an hourly basis and is provided without any pre-split. The dataset is available at <https://github.com/thuml/Time-Series-Library>.
  - **Cora and Citeseer.** The citation network datasets, "Cora" and "CiteSeer," consist of nodes representing documents and edges representing citation links between them. Both datasets provide training, validation, and test splits through binary masks. The Cora dataset contains 2,708 nodes, 10,556 edges, 1,433 features, and 7 classes, while CiteSeer consists of 3,327 nodes, 9,104 edges, 3,703 features, and 6 classes. We use the version provided by Fey & Lenssen (2019).

### A.3 BASELINES

**Human Models** We select top-performing models based on evaluations from Papers with Code benchmarks or Kaggle notebooks, where the similar tasks and datasets are applicable. The chosen models for relevant downstream tasks are described below.

- **Image Classification.** The human models for image classification tasks are obtained from a Kaggle notebook available at <https://www.kaggle.com/code/mohamedhassanali/butterfly-classify-pytorch-pretrained-model-acc-99/notebook>, utilizing a pretrained ResNet-18 model.
- **Text Classification.** For text classification tasks, two models are employed. A Word2Vec-based XGBoost model is applied to the e-commerce text dataset <https://www.kaggle.com/code/sugataghosh/e-commerce-text-classification-tf-idf-word2vec#Word2Vec-Hyperparameter-Tuning>, while the XLM-RoBERTa model is used for the textual entailment dataset <https://www.kaggle.com/code/vbookshelf/basics-of-bert-and-xlm-roberta-pytorch>.
- **Tabular Classification.** Due to the absence of a similar model in the repository, we use the state-of-the-art TabPFN model (Hollmann et al., 2023a) designed for tabular classification tasks.
- **Tabular Regression.** For tabular regression tasks, we adopt two models specifically designed for the given datasets, which are available at <https://www.kaggle.com/code/shatabdi5/crab-age-regression> for the crab age dataset and at <https://www.kaggle.com/code/mahmoudmagdyelnahal/crop-yield-prediction-99/notebook> for the crop yield dataset.
- **Tabular Clustering.** For unsupervised clustering tasks, we use manually hyperparameter-tuned KMeans clustering, following the approach outlined in <https://www.kaggle.com/code/samueltortinhas/tps-july-22-unsupervised-clustering>, as the baseline.
- **Time-Series Forecasting.** In this task, we use the state-of-the-art iTransformer (Liu et al., 2024b), which is designed for the same task and datasets as the baseline model.
- **Node Classification.** For node classification tasks, we also employ a state-of-the-art graph neural network-based model, PMLP (Yang et al., 2023), as the handcrafted baseline for both datasets.

**AutoGluon** We adopt AutoGluon as the baseline because it is a state-of-the-art AutoML framework capable of handling various downstream tasks and data modalities, with the exception of graph data. There are three variants of AutoGluon: AutoGluon-TS (Shchur et al., 2023) for time series, AutoGluon-Tabular (Erickson et al., 2020) for tabular machine learning, and AutoGluon-Multimodal (Tang et al., 2024) for computer vision and natural language processing tasks.

**GPT-3.5 and GPT-4** For GPT-3.5 and GPT-4, we use the *gpt-3.5-turbo-0125* and *gpt-4-2024-05-13* models via the OpenAI API. We implement the zero-shot baselines using the prompt below.

#### Zero-Shot Prompt for GPT-3.5 and GPT-4 Baselines

```
You are a helpful intelligent assistant. Now please help solve the following machine
learning task.
[Task]
{user instruction}
[{{file_name}}.py] ```python
{full-pipeline skeleton script}
```
Start the python code with "```python". Please ensure the completeness of the code so that
it can be run without additional modifications.
```

**DS-Agent** We reproduce the DS-Agent (Guo et al., 2024a) baseline using the official source code. However, it is important to note that our framework encompasses the entire process from data retrieval/loading to deployment, whereas DS-Agent focuses solely on the modeling aspect, assuming complete data and evaluation codes are provided. In this paper, we utilize the deployment stage of DS-Agent along with its collected case banks and Adapter prompt for the same tasks, as the source code for manual human insights collection during the development stage is unavailable.

#### A.4 EVALUATION METRICS

**Success Rate (SR)** We employ the success rate (Guo et al., 2024a; Hong et al., 2024a), which evaluates whether the models built by an LLM agent are executable in the given runtime environment. Success rate is used to assess code execution.

For the *constraint-free* setting, we apply a three-level grading scale as follows.

- **0.00:** Code cannot be executed.
- **0.50:** Code provides a runnable ML/DL model.
- **1.00:** Code provides a runnable model and an accessible deployment endpoint (e.g., Gradio).

For the *constraint-aware* setting, we use a five-level grading scale to evaluate whether the code executes successfully and satisfies the given constraints. The grading criteria are as follows.

- **0.00:** Code cannot be executed.
- **0.25:** Code provides a runnable ML/DL model.
- **0.50:** Code provides a runnable model and an accessible deployment endpoint (e.g., Gradio).
- **0.75:** Code provides a deployed, runnable model that partially meets constraints (e.g., target performance, inference time, and model size).
- **1.00:** Code provides a deployed, runnable model that fully meets constraints.

**Normalized Performance Score (NPS)** In this paper, each downstream task is associated with a specific evaluation metric, which may vary between tasks. These metrics include accuracy, F1-score, and RMSLE. For metrics such as accuracy and F1-score, we present the raw values to facilitate comparison across identical data tasks. For performance metrics where lower values indicate better performance, such as loss-based metrics, we normalize all performance values  $s$  using the following transformation:  $NPS = \frac{1}{1+s}$ . This transformation ensures that metrics like RMSLE are scaled between 0 and 1, with higher NPS values indicating better performance.

Note that achieving downstream task performance (NPS) requires a runnable model, i.e.,  $SR > 0$ . If the model cannot run, the NPS is zero by default as it cannot make any predictions.

**Comprehensive Score (CS)** To evaluate both the success rate and the downstream task performance of the generated AutoML pipelines simultaneously, we calculate CS as a weighted sum of SR and NPS, as follows:  $CS = 0.5 \times SR + 0.5 \times NPS$ .

## B PROMPTS FOR AUTOML-AGENT

### B.1 AGENT SPECIFICATIONS

This subsection provides the *system prompt* design for agent specifications in *AutoML-Agent*, including Agent Manager (B.1.1), Prompt Agent (B.1.2), Data Agent (B.1.3), Model Agent (B.1.4), and Operation Agent (B.1.5).

#### B.1.1 AGENT MANAGER

##### System Message for Agent Manager Specification

You are an experienced senior project manager of a automated machine learning project (AutoML). You have two main responsibilities as follows.

1. Receive requirements and/or inquiries from users through a well-structured JSON object.
2. Using recent knowledge and state-of-the-art studies to devise promising high-quality plans for data scientists, machine learning research engineers, and MLOps engineers in your team to execute subsequent processes based on the user requirements you have received.

### B.1.2 PROMPT AGENT

#### System Message for Prompt Agent Specification

You are an assistant project manager in the AutoML development team.  
 Your task is to parse the user's requirement into a valid JSON format using the JSON specification schema as your reference. Your response must exactly follow the given JSON schema and be based only on the user's instruction.  
 Make sure that your answer contains only the JSON response without any comment or explanation because it can cause parsing errors.

```
#JSON SPECIFICATION SCHEMA#
```json
{json_specification}
```
```

Your response must begin with "```json" or "{" and end with "```" or "}", respectively.

### B.1.3 DATA AGENT

#### System Message for Data Agent Specification

You are the world's best data scientist of an automated machine learning project (AutoML) that can find the most relevant datasets, run useful preprocessing, perform suitable data augmentation, and make meaningful visualization to comprehensively understand the data based on the user requirements. You have the following main responsibilities to complete.

1. Retrieve a dataset from the user or search for the dataset based on the user instruction.
2. Perform data preprocessing based on the user instruction or best practice based on the given tasks.
3. Perform data augmentation as necessary.
4. Extract useful information and underlying characteristics of the dataset.

### B.1.4 MODEL AGENT

#### System Message for Model Agent Specification

You are the world's best machine learning research engineer of an automated machine learning project (AutoML) that can find the optimal candidate machine learning models and artificial intelligence algorithms for the given dataset(s), run hyperparameter tuning to optimize the models, and perform metadata extraction and profiling to comprehensively understand the candidate models or algorithms based on the user requirements. You have the following main responsibilities to complete.

1. Retrieve a list of well-performing candidate ML models and AI algorithms for the given dataset based on the user's requirement and instruction.
2. Perform hyperparameter optimization for those candidate models or algorithms.
3. Extract useful information and underlying characteristics of the candidate models or algorithms using metadata extraction and profiling techniques.
4. Select the top-k ('k' will be given) well-performing models or algorithms based on the hyperparameter optimization and profiling results.

### B.1.5 OPERATION AGENT

#### System Message for Operation Agent Specification

You are the world's best MLOps engineer of an automated machine learning project (AutoML) that can implement the optimal solution for production-level deployment, given any datasets and models. You have the following main responsibilities to complete.

1. Write accurate Python codes to retrieve/load the given dataset from the corresponding source.
2. Write effective Python codes to preprocess the retrieved dataset.
3. Write precise Python codes to retrieve/load the given model and optimize it with the suggested hyperparameters.
4. Write efficient Python codes to train/finetune the retrieved model.
5. Write suitable Python codes to prepare the trained model for deployment. This step may include model compression and conversion according to the target inference platform.



6. Write Python codes to build the web application demo using the Gradio library.
7. Run the model evaluation using the given Python functions and summarize the results for validation againsts the user’s requirements.

## B.2 PROMPTS FOR RETRIEVAL-AUGMENTED PLANNING

This subsection presents prompts for planning-related processes (Figure 2(a)), including knowledge retrieval and summary prompts (B.2.1), planning prompt (B.2.2), and plan revision prompt (B.2.3).

### B.2.1 KNOWLEDGE RETRIEVAL PROMPT

#### Prompt for Knowledge Retrieval and Summary for Planning

##### Kaggle Notebook

I searched the Kaggle Notebooks to find state-of-the-art solutions using the keywords: {user\_task} {user\_domain}. Here is the result:

```
=====
{context}
=====
```

Please summarize the given pieces of Python notebooks into a single paragraph of useful knowledge and insights. Do not include the source codes. Instead, extract the insights from the source codes. We aim to use your summary to address the following user’s requirements.

```
# User’s Requirements
{user_requirement_summary}
```

##### Papers With Code

I searched the paperswithcode website to find state-of-the-art models using the keywords: {user\_area} and {user\_task}. Here is the result:

```
=====
{context}
=====
```

Please summarize the given pieces of search content into a single paragraph of useful knowledge and insights. We aim to use your summary to address the following user’s requirements.

```
# User’s Requirements
{user_requirement_summary}
```

##### arXiv

I searched the arXiv papers using the keywords: {task\_kw} and {domain\_kw}. Here is the result:

```
=====
{context}
=====
```

Please summarize the given pieces of arXiv papers into a single paragraph of useful knowledge and insights. We aim to use your summary to address the following user’s requirements.

```
# User’s Requirements
{user_requirement_summary}
```

##### Google WebSearch

I searched the web using the query: {search\_query}. Here is the result:

```
=====
{context}
=====
```

Please summarize the given pieces of search content into a single paragraph of useful knowledge and insights.

We aim to use your summary to address the following user’s requirements.

```
# User’s Requirements
{user_requirement_summary}
```

### Summary

Please extract and summarize the following group of contents collected from different online sources into a chunk of insightful knowledge. Please format your answer as a list of suggestions. I will use them to address the user's requirements in machine learning tasks.

```
# Source: Google Web Search
{search_summary}
=====
```

```
# Source: arXiv Papers
{arxiv_summary}
=====
```

```
# Source: Kaggle Hub
{kaggle_summary}
=====
```

```
# Source: PapersWithCode
{pwc_summary}
=====
```

The user's requirements are summarized as follows.  
{user\_requirement\_summary}

## B.2.2 PLANNING PROMPT

### Prompt for Retrieval-Augmented Planning

Now, I want you to devise an end-to-end actionable plan according to the user's requirements described in the following JSON object.

```
```json
{user_requirements}
```
```

Here is a list of past experience cases and knowledge written by an human expert for a relevant task:

```
{plan_knowledge}
```

When devising a plan, follow these instructions and do not forget them:

- Ensure that your plan is up-to-date with current state-of-the-art knowledge.
- Ensure that your plan is based on the requirements and objectives described in the above JSON object.
- Ensure that your plan is designed for AI agents instead of human experts. These agents are capable of conducting machine learning and artificial intelligence research.
- Ensure that your plan is self-contained with sufficient instructions to be executed by the AI agents.
- Ensure that your plan includes all the key points and instructions (from handling data to modeling) so that the AI agents can successfully implement them. Do NOT directly write the code.
- Ensure that your plan completely include the end-to-end process of machine learning or artificial intelligence model development pipeline in detail (i.e., from data retrieval to model training and evaluation) when applicable based on the given requirements.

## B.2.3 PLAN REVISION PROMPT

### Prompt for Plan Revision

Now, you will be asked to revise and rethink {num2words(n\_plans)} different end-to-end actionable plans according to the user's requirements described in the JSON object below.

```
```json
{user_requirements}
```
```

Please use to the following findings and insights summarized from the previously failed plans. Try as much as you can to avoid the same failure again.  
{fail\_rationale}

Finally, when devising a plan, follow these instructions and do not forget them:

- Ensure that your plan is up-to-date with current state-of-the-art knowledge.
- Ensure that your plan is based on the requirements and objectives described in the above JSON object.
- Ensure that your plan is designed for AI agents instead of human experts. These agents are capable of conducting machine learning and artificial intelligence research.
- Ensure that your plan is self-contained with sufficient instructions to be executed by the AI agents.
- Ensure that your plan includes all the key points and instructions (from handling data to modeling) so that the AI agents can successfully implement them. Do NOT directly write the code.
- Ensure that your plan completely include the end-to-end process of machine learning or artificial intelligence model development pipeline in detail (i.e., from data retrieval to model training and evaluation) when applicable based on the given requirements.

### B.3 PROMPTS FOR PROMPTING-BASED PLAN EXECUTION

This subsection presents prompts for prompting-based plan execution processes (Figure 2(b)), including plan decomposition (Data Agent (B.3.1) and Model Agent (B.3.2)), pseudo data analysis (B.3.3), and training-free model search and HPO (B.3.4).

#### B.3.1 PLAN DECOMPOSITION: DATA AGENT

##### Prompt for Plan Decomposition: Data Agent

As a proficient data scientist, summarize the following plan given by the senior AutoML project manager according to the user’s requirements and your expertise in data science.

```
# User’s Requirements
```json
{user_requirements}
```

# Project Plan
{plan}
```

The summary of the plan should enable you to fulfill your responsibilities as the answers to the following questions by focusing on the data manipulation and analysis.

1. How to retrieve or collect the dataset(s)?
2. How to preprocess the retrieved dataset(s)?
3. How to efficiently augment the dataset(s)?
4. How to extract and understand the underlying characteristics of the dataset(s)?

Note that you should not perform data visualization because you cannot see it. Make sure that another data scientist can exactly reproduce the results based on your summary.

#### B.3.2 PLAN DECOMPOSITION: MODEL AGENT

##### Prompt for Plan Decomposition: Model Agent

As a proficient machine learning research engineer, summarize the following plan given by the senior AutoML project manager according to the user’s requirements, your expertise in machine learning, and the outcomes from data scientist.

```
**User’s Requirements**
```json
{user_requirements}
```

**Project Plan**
{project_plan}

**Explanations and Results from the Data Scientist**
```

```
{data_result}
```

The summary of the plan should enable you to fulfill your responsibilities as the answers to the following questions by focusing on the modeling and optimization tasks.

1. How to retrieve or find the high-performance model(s)?
2. How to optimize the hyperparameters of the retrieved models?
3. How to extract and understand the underlying characteristics of the dataset(s)?
4. How to select the top-k models or algorithms based on the given plans?

### B.3.3 PSEUDO DATA ANALYSIS BY DATA AGENT

#### Prompt for Pseudo Data Analysis

As a proficient data scientist, your task is to explain **detailed** steps for data manipulation and analysis parts by executing the following machine learning development plan.

```
# Plan
{decomposed_data_plan}
```

```
# Potential Source of Dataset
{available_sources}
```

Make sure that your explanation follows these instructions:

- All of your explanation must be self-contained without using any placeholder to ensure that other data scientists can exactly reproduce all the steps, but do not include any code.
- Include how and where to retrieve or collect the data.
- Include how to preprocess the data and which tools or libraries are used for the preprocessing.
- Include how to do the data augmentation with details and names.
- Include how to extract and understand the characteristics of the data.
- Include reasons why each step in your explanations is essential to effectively complete the plan.

Note that you should not perform data visualization because you cannot see it. Make sure to focus only on the data part as it is your expertise. Do not conduct or perform anything regarding modeling or training.

After complete the explanations, explicitly specify the (expected) outcomes and results both quantitative and qualitative of your explanations.

### B.3.4 TRAINING-FREE MODEL SEARCH AND HPO BY MODEL AGENT

#### Prompt for Training-Free Model Search and HPO

As a proficient machine learning research engineer, your task is to explain **detailed** steps for modeling and optimization parts by executing the following machine learning development plan with the goal of finding top-{k} candidate models/algorithms.

```
# Suggested Plan
{decomposed_model_plan}
```

```
# Available Model Source
{available_sources}
```

Make sure that your explanation for finding the top-{k} high-performance models or algorithms follows these instructions:

- All of your explanations must be self-contained without using any placeholder to ensure that other machine learning research engineers can exactly reproduce all the steps, but do not include any code.
- Include how and where to retrieve or find the top-{k} well-performing models/algorithms.
- Include how to optimize the hyperparameters of the candidate models or algorithms by clearly specifying which hyperparameters are optimized in detail.
- Corresponding to each hyperparameter, explicitly include the actual numerical value that you think it is the optimal value for the given dataset and machine learning task.
- Include how to extract and understand the characteristics of the candidate models or algorithms, such as their computation complexity, memory usage, and inference latency. This part is not related to visualization and interpretability.
- Include reasons why each step in your explanations is essential to effectively complete the plan.

Make sure to focus only on the modeling part as it is your expertise. Do not conduct or perform anything regarding data manipulation or analysis.

After complete the explanations, explicitly specify the names and (expected) quantitative performance using relevant numerical performance and complexity metrics (e.g., number of parameters, FLOPs, model size, training time, inference speed, and so on) of the {num2words(k)} candidate models/algorithms potentially to be the optimal model below. Do not use any placeholder for the quantitative performance. If you do not know the exact values, please use the knowledge and expertise you have to estimate those performance and complexity values.

## B.4 PROMPTS FOR MULTI-STAGE VERIFICATION

This subsection presents prompts for multi-stage verification (Figure 2(c)), which ensures the correctness of intermediate results between steps in the framework. These stages include request verification B.4.1), execution verification B.4.2, and implementation verification B.4.3.

### B.4.1 REQUEST VERIFICATION

#### Request Verification (Relevancy)

Is the following statement relevant to machine learning or artificial intelligence?  
{user instruction}  
Answer only 'Yes' or 'No'

#### Request Verification (Adequacy)

Given the following JSON object representing the user’s requirement for a potential ML or AI project, please tell me whether we have essential information (e.g., problem and dataset) to be used for a AutoML project?  
Please note that our users are not AI experts, you must focus only on the essential requirements, e.g., problem and brief dataset descriptions.  
You do not need to check every details of the requirements. You must also answer ‘yes’ even though it lacks detailed and specific information.

```
'''json
{parsed user requirements}
'''
```

Please answer with this format: ‘a ‘yes’ or ‘no’ answer; your reasons for the answer’ by using ‘;’ to separate between the answer and its reasons.  
If the answer is ‘no’, you must tell me the alternative solutions or examples for completing such missing information.

### B.4.2 EXECUTION VERIFICATION

#### Execution Verification

Given the proposed solution and user’s requirements, please carefully check and verify whether the proposed solution ‘pass’ or ‘fail’ the user’s requirements.

**\*\*Proposed Solution and Its Implementation\*\***  
Data Manipulation and Analysis: {data\_agent\_outcomes}  
Modeling and Optimization: {model\_agent\_outcomes}

**\*\*User Requirements\*\***  
'''json  
{user\_requirements}  
'''

Answer only ‘Pass’ or ‘Fail’

### B.4.3 IMPLEMENTATION VERIFICATION

#### Implementation Verification

As the project manager, please carefully verify whether the given Python code and results satisfy the user's requirements.

```
- Python Code
```python
{implementation_result['code']}
```

- Code Execution Result
{implementation_result['action_result']}

- User's Requirements
{user_requirements}

Answer only 'Pass' or 'Fail'
```

## C EXAMPLES OF INTERMEDIATE RESULTS

This section presents examples of intermediate results from the main steps of our framework, including prompt parsing results (C.1), retrieval-augmented planning results (C.2), plan decomposition results (C.3), and prompting-based plan execution results (C.4).

### C.1 PROMPT PARSING

#### Task Description and Parsed JSON User Requirements

**Constraint-Free User Prompt:** Build a model to classify banana quality as Good or Bad based on their numerical information about bananas of different quality (size, weight, sweetness, softness, harvest time, ripeness, and acidity). We have uploaded the entire dataset for you here in the banana\_quality.csv file.

**Parsed Output:**

```
{
  "user": {"intent": "build", "expertise": "medium"},
  "problem": {
    "area": "tabular data analysis",
    "downstream_task": "tabular classification",
    "application_domain": "agriculture",
    "description": "Build a model to classify banana quality as Good or Bad based on their numerical information about bananas of different quality (size, weight, sweetness, softness, harvest time, ripeness, and acidity).",
    "performance_metrics": [],
    "complexity_metrics": []
  },
  "dataset": [
    {
      "name": "banana_quality",
      "modality": ["tabular"],
      "target_variables": ["quality"],
      "specification": null,
      "description": "A dataset containing numerical information about bananas of different quality, including size, weight, sweetness, softness, harvest time, ripeness, and acidity.",
      "preprocessing": [],
      "augmentation": [],
      "visualization": [],
      "source": "user-upload"
    }
  ],
  "model": [
    {
      "name": "",
      "family": "",
      "type": "classical machine learning",
      "specification": null,
    }
  ]
}
```

```

        "description": "A model to classify banana quality as Good or Bad based on their
            numerical information."
    },
    ],
    ...
}

```

**Constraint-Aware User Prompt:** Build a machine learning model, potentially XGBoost or LightGBM, to classify banana quality as Good or Bad based on their numerical information about bananas of different quality (size, weight, sweetness, softness, harvest time, ripeness, and acidity). We have uploaded the entire dataset for you here in the banana\_quality.csv file. The model must achieve at least 0.98 accuracy.

**Parsed Output:**

```

{
  "user": {"intent": "build", "expertise": "medium"},
  "problem": {
    "area": "tabular data analysis",
    "downstream_task": "tabular classification",
    "application_domain": "agriculture",
    "description": "Build a machine learning model, potentially XGBoost or LightGBM, to
      classify banana quality as Good or Bad based on their numerical information
      about bananas of different quality (size, weight, sweetness, softness, harvest
      time, ripeness, and acidity). The model must achieve at least 0.98 accuracy.",
    "performance_metrics": [
      {"name": "accuracy", "value": 0.98}
    ],
    "complexity_metrics": []
  },
  "dataset": [
    {
      "name": "banana_quality",
      "modality": ["tabular"],
      "target_variables": ["quality"],
      "specification": null,
      "description": "A dataset containing numerical information about bananas of
        different quality, including size, weight, sweetness, softness, harvest time,
        ripeness, and acidity.",
      "preprocessing": [],
      "augmentation": [],
      "visualization": [],
      "source": "user-upload"
    }
  ],
  "model": [
    {
      "name": "XGBoost",
      "family": "ensemble models",
      "type": "ensemble",
      "specification": null,
      "description": "A potential model to classify banana quality as Good or Bad,
        potentially using XGBoost or LightGBM."
    }
  ],
  ...
}

```

## C.2 RETRIEVAL-AUGMENTED PLANNING

### Example Results from Retrieval-Augmented Planning with $P = 3$

**User Input:** I need a very accurate model to classify images in the Butterfly Image Classification dataset into their respective categories. The dataset has been uploaded with its label information in the labels.csv file.

#### Plan #1

### End-to-End Actionable Plan for Butterfly Image Classification



```

**Objective**: Develop a highly accurate neural network model to classify images in the
    Butterfly Image Classification dataset into their respective categories.

#### 1. Data Collection and Preprocessing
1. **Dataset Loading**:
    - Load the Butterfly Image Classification dataset and the associated labels from `labels.csv`.
    - Verify the integrity of the dataset, ensuring all images have corresponding labels.

2. **Data Preprocessing**:
    - **Image Resizing**: Resize all images to a uniform size (e.g., 224x224 pixels) to ensure consistency.
    - **Normalization**: Normalize pixel values to the range [0, 1] or standardize to have zero mean and unit variance.

3. **Data Augmentation**:
    - Apply data augmentation techniques such as cropping, horizontal and vertical flips, rescaling, shear, zoom, and rotation to increase dataset variability.
    - Utilize libraries like TensorFlow's `ImageDataGenerator` or PyTorch's `transforms` for implementation.

#### 2. Model Selection and Architecture
1. **Transfer Learning**:
    - Leverage pre-trained convolutional neural networks (CNNs) such as VGG16, VGG19, MobileNet, Xception, ResNet50, and InceptionV3.
    - Fine-tune these pre-trained models on the Butterfly Image Classification dataset to utilize their learned features.

2. **Model Architecture**:
    - Start with InceptionV3, known for high accuracy in similar tasks.
    - Explore other models like ResNet, DenseNet, EfficientNet, and Vision Transformers.
    - Consider combining CNNs with Transformer-based architectures to capture both local and global features.

3. **Ensemble Methods**:
    - Implement ensemble methods by combining predictions from multiple models to improve overall classification performance.
    - Use techniques like bagging or stacking to enhance model robustness.

#### 3. Model Training
1. **Dataset Splitting**:
    - Split the dataset into training, validation, and testing sets (e.g., 80% training and 20% validation).

2. **Training Configuration**:
    - Use categorical cross-entropy as the loss function.
    - Select optimizers like SGD or RMSprop.
    - Set up callbacks for early stopping and learning rate reduction on plateau to prevent overfitting.

3. **Hyperparameter Tuning**:
    - Perform hyperparameter tuning using libraries like PyCaret or Optuna to find the optimal configuration.
    - Utilize techniques like grid search or random search for hyperparameter optimization.

#### 4. Model Evaluation and Validation
1. **Performance Metrics**:
    - Regularly monitor accuracy as the primary performance metric.
    - Additionally, track precision, recall, and F1-score to ensure balanced performance across all classes.

2. **Cross-Validation**:
    - Employ k-fold cross-validation to assess the model's robustness and generalization ability.

3. **Handling Class Imbalance**:
    - Address class imbalance through weighted loss functions or oversampling techniques to ensure fair performance across all categories.

#### 5. Interpretability and Debugging
1. **Model Interpretability**:
    - Use interpretability methods like Grad-CAM or LIME to visualize which parts of the images influence the model's predictions.
    - Ensure the model focuses on relevant biological features of butterflies.

#### 6. Reproducibility and Documentation
1. **Reproducibility**:
    - Set random seeds for all operations to ensure reproducibility of results.

```

- Document the entire process, including data splits, model configurations, and training parameters.
2. **\*\*Model Deployment\*\***:
    - Prepare the model for deployment by exporting it in a suitable format (e.g., TensorFlow SavedModel, ONNX).
    - Set up an inference engine and deployment endpoint to make predictions on new data.
  3. **\*\*Demonstration\*\***:
    - Provide a demonstration of the model's capabilities by creating a user-friendly interface or visualization to showcase its performance.
- By following this comprehensive plan, AI agents can develop a robust and accurate neural network model tailored for classifying butterfly images in the given dataset.

## Plan #2

## End-to-End Actionable Plan for Butterfly Image Classification

### ### 1. Data Collection

- **\*\*Dataset Name\*\***: Butterfly Image Classification
- **\*\*Source\*\***: User-uploaded
- **\*\*Description\*\***: A dataset for classifying images of butterflies into their respective categories.
- **\*\*Label Information\*\***: Available in 'labels.csv' file.

### ### 2. Data Preprocessing

- **\*\*Image Preprocessing\*\***:
  - **\*\*Resize Images\*\***: Resize all images to a consistent size (e.g., 224x224 pixels) to match the input size required by pre-trained models.
  - **\*\*Normalization\*\***: Normalize pixel values to the range [0, 1] or use mean subtraction based on the pre-trained model's requirements.
  - **\*\*Handling High-Resolution and Illumination Variations\*\***: Apply techniques to standardize illumination and handle high-resolution images.

### ### 3. Data Augmentation

- **\*\*Techniques\*\***:
  - **\*\*Cropping\*\***: Randomly crop sections of the images.
  - **\*\*Flipping\*\***: Apply horizontal and vertical flips.
  - **\*\*Rescaling\*\***: Rescale images by a factor.
  - **\*\*Shearing\*\***: Apply shear transformations.
  - **\*\*Zooming\*\***: Apply random zoom.
  - **\*\*Rotation\*\***: Rotate images by random angles.
- **\*\*Libraries\*\***: Use TensorFlow's 'ImageDataGenerator' or PyTorch's 'transforms'.

### ### 4. Dataset Splitting

- **\*\*Training Set\*\***: 80% of the dataset
- **\*\*Validation Set\*\***: 20% of the dataset

### ### 5. Model Selection and Architecture

- **\*\*Transfer Learning\*\***:
  - **\*\*Pre-trained Models\*\***: Utilize models such as VGG16, VGG19, MobileNet, Xception, ResNet50, and InceptionV3.
  - **\*\*Fine-tuning\*\***: Fine-tune these models on the Butterfly Image Classification dataset.
- **\*\*Model Architecture\*\***:
  - **\*\*Primary Model\*\***: Start with InceptionV3 due to its high performance in similar tasks.
  - **\*\*Alternative Models\*\***: Evaluate ResNet, DenseNet, EfficientNet, and Vision Transformers.
  - **\*\*Combined Architectures\*\***: Explore combining CNNs with Transformer-based architectures to capture both local and global features.

### ### 6. Handling Class Imbalance

- **\*\*Techniques\*\***:
  - **\*\*Weighted Loss Functions\*\***: Assign higher weights to underrepresented classes.
  - **\*\*Oversampling\*\***: Increase the number of samples for underrepresented classes.

### ### 7. Model Training

- **\*\*Loss Function\*\***: Use categorical cross-entropy.
- **\*\*Optimizers\*\***: Start with SGD or RMSprop.
- **\*\*Hyperparameter Tuning\*\***:
  - **\*\*Libraries\*\***: Use PyCaret or Optuna for hyperparameter optimization.
  - **\*\*Parameters\*\***: Tune learning rate, batch size, number of epochs, etc.
- **\*\*Cross-Validation\*\***: Use k-fold cross-validation to ensure model robustness and generalization.

### ### 8. Performance Evaluation

- **\*\*Metrics\*\***:
  - **\*\*Accuracy\*\***: Primary metric for evaluation.

```

- Precision, Recall, F1-Score: Monitor these metrics to ensure balanced performance across classes.
- Evaluation Techniques:
  - Confusion Matrix: Visualize to understand misclassifications.
  - ROC Curve: Evaluate model performance across different thresholds.

9. Ensemble Methods and Knowledge Distillation
- Ensemble Methods: Combine predictions from multiple models (e.g., patch-wise and image-wise classifiers).
- Knowledge Distillation: Transfer knowledge from larger teacher models to smaller student models to create compact yet powerful models.

10. Interpretability and Debugging
- Techniques:
  - Grad-CAM: Visualize which parts of the images influence the model's predictions.
  - LIME: Provide explanations for individual predictions.
- Purpose: Ensure the model focuses on relevant biological features and understand model behavior.

11. Reproducibility
- Random Seeds: Set random seeds for all processes to ensure reproducibility.
- Documentation: Document data splits, model configurations, and training parameters.

12. Deployment
- Target Device: Specify the target device for deployment (e.g., cloud, edge device).
- Deployment Endpoint: Define the endpoint for model deployment.
- Inference Engine: Choose an appropriate inference engine (e.g., TensorFlow Serving, TorchServe).
- Demonstration: Prepare a demonstration of the model's capabilities and performance.

By following this comprehensive plan, AI agents can effectively build, train, and deploy a highly accurate neural network model tailored for classifying butterfly images in the specified dataset.

```

### Plan #3

### End-to-End Actionable Plan for Butterfly Image Classification using Neural Networks

#### ### 1. **Project Initialization**

**Objective**: Build a highly accurate image classification model for the Butterfly Image Classification dataset in the biology domain.

**User Expertise**: Medium

---

#### ### 2. **Data Collection and Preprocessing**

**Dataset Source**: User-uploaded Butterfly Image Classification dataset with label information in 'labels.csv'.

**Steps**:

1. **Load Dataset**:
  - Load images and labels from the specified dataset and 'labels.csv'.
2. **Data Cleaning**:
  - Ensure all images are correctly labeled.
  - Remove any corrupted or unusable images.
3. **Data Preprocessing**:
  - Resize images to a consistent size (e.g., 224x224 pixels).
  - Normalize pixel values to the range [0, 1].
4. **Data Augmentation**:
  - Apply augmentation techniques such as cropping, horizontal and vertical flips, rescaling, shear, zoom, and rotation to increase dataset size and variability.
  - Use libraries like TensorFlow's 'ImageDataGenerator' or PyTorch's 'transforms'.

---

#### ### 3. **Model Development**

**Model Type**: Neural Networks

**Steps**:

```

1. **Leverage Transfer Learning**:
  - Use pre-trained CNNs such as VGG16, VGG19, MobileNet, Xception, ResNet50, and InceptionV3.
  - Fine-tune these models on the Butterfly Image Classification dataset.

2. **Model Selection and Architecture**:
  - Start with InceptionV3 due to its high accuracy in similar tasks.
  - Experiment with other models like ResNet, DenseNet, EfficientNet, and Vision Transformers.
  - Consider combining CNNs with Transformer-based architectures for improved performance.

3. **Ensemble Methods and Knowledge Distillation**:
  - Implement ensemble methods by combining patch-wise and image-wise classifiers.
  - Use knowledge distillation to transfer knowledge from larger models to smaller, more efficient models.

---

#### 4. **Training and Evaluation**

**Steps**:

1. **Dataset Splitting**:
  - Split the dataset into training (80%), validation (10%), and testing (10%) sets.

2. **Model Training**:
  - Use categorical cross-entropy as the loss function.
  - Utilize optimizers such as SGD or RMSprop.
  - Implement early stopping and learning rate reduction on plateau to prevent overfitting.

3. **Performance Monitoring**:
  - Regularly monitor metrics such as accuracy, precision, recall, and F1-score.
  - Use validation data to fine-tune hyperparameters.

4. **Handling Class Imbalance**:
  - Address class imbalance using weighted loss functions or oversampling techniques.

5. **Hyperparameter Tuning and Cross-Validation**:
  - Perform hyperparameter tuning with libraries like PyCaret or Optuna.
  - Use cross-validation to ensure model robustness and generalization.

---

#### 5. **Model Interpretability and Debugging**

**Steps**:

1. **Interpretability**:
  - Use Grad-CAM or LIME to visualize important image regions influencing model predictions.
  - Ensure the model focuses on relevant biological features.

2. **Debugging**:
  - Analyze misclassified images to understand model limitations.
  - Iterate on model improvements based on these insights.

---

#### 6. **Reproducibility and Documentation**

**Steps**:

1. **Reproducibility**:
  - Set random seeds for all processes to ensure reproducibility.
  - Document all steps, including data splits, model configurations, and training parameters.

2. **Version Control**:
  - Use version control systems (e.g., Git) to track changes and collaborate efficiently.

---

#### 7. **Deployment and Inference**

**Steps**:

1. **Model Export**:
  - Export the trained model in a suitable format (e.g., TensorFlow SavedModel, ONNX).

```

```

2. **Deployment**:
  - Choose a deployment platform (e.g., AWS SageMaker, Google AI Platform).
  - Set up the deployment endpoint for real-time inference.

3. **Inference Engine**:
  - Optimize the model for the target device (e.g., mobile, edge devices) if necessary.
  - Implement efficient inference pipelines to handle incoming image data.

---

#### 8. **Demonstration and User Feedback**

**Steps**:

1. **Demonstration**:
  - Create a demo application to showcase the model's capabilities.
  - Provide a user-friendly interface for uploading images and viewing classification results.

2. **User Feedback**:
  - Collect feedback from users to identify areas for improvement.
  - Implement iterative updates based on user insights.

---

By following this detailed plan, AI agents can build a robust and accurate neural network model for classifying butterfly images in the biology domain, ensuring high performance and user satisfaction.

```

### C.3 PLAN DECOMPOSITION

#### Example Results from Plan Decomposition

##### Data Agent

```

### Summary of the Plan for Fashion Image Classification Using PyTorch

#### 1. Dataset Retrieval and Collection
- **Retrieve Dataset**: Load the image paths and corresponding labels from the provided 'train_labels.csv' and 'test_labels.csv' files. Ensure the dataset includes images of four clothing categories.

#### 2. Checking Dataset Types, Shapes, and Properties
- **Data Types**: Verify that the images are in the correct format (e.g., JPEG, PNG) and the labels are correctly mapped to the images.
- **Shapes**: Check the dimensions of the images. If images are in grayscale, convert them to RGB by replicating the single channel across three channels.
- **Properties**: Ensure that the dataset is balanced across the four clothing categories and that there are no missing or corrupted files.

#### 3. Preprocessing the Dataset
- **Normalization**: Scale pixel values to the range [0, 1].
- **Resizing**: Resize all images to a consistent size (e.g., 224x224 pixels) suitable for model input.
- **Splitting**: Split the training data into training and validation sets (e.g., 80% training, 20% validation) to monitor the model's performance and prevent overfitting.

#### 4. Efficient Data Augmentation
- **Techniques**: Apply data augmentation techniques to increase the diversity and robustness of the training data. Techniques include:
  - Random horizontal flip
  - Normalization
  - Random cropping
  - Random erasing
  - Rotating
  - Shifting
  - Zooming
- **Pipelines**:
  - **Pipeline 1**: Random horizontal flip and normalization.
  - **Pipeline 2**: Random horizontal flip, normalization, random cropping, and random erasing.

#### 5. Extracting and Understanding Underlying Characteristics
- **Class Distribution**: Analyze the distribution of classes to ensure a balanced dataset.
- **Image Quality**: Check for variations in image quality and resolution.

```

- **Data Insights**: Identify any patterns or anomalies in the data that could affect model performance.
  - **Performance Metrics**: Use metrics such as accuracy and confusion matrix to evaluate the model's performance and understand misclassification patterns.
- By following these steps, another data scientist should be able to reproduce the results and build a robust image classification model that meets the specified accuracy requirements and provides insightful performance analysis.

## Model Agent

### Summary of the Plan for Fashion Image Classification Using PyTorch

### #### User's Requirements

The user needs a classification model that categorizes images into one of four clothing categories. The model should achieve at least 85% accuracy on the test set and be implemented using PyTorch. Data augmentation techniques and a confusion matrix for evaluation are also required.

### ### Project Plan Overview

#### #### 1. Data Preparation

- **Load Dataset**: Load image paths and labels from 'train\_labels.csv' and 'test\_labels.csv'.
- **Convert Grayscale to RGB**: Convert grayscale images to RGB.
- **Preprocessing**: Normalize pixel values to [0, 1] and resize images to 224x224 pixels.
- **Split Dataset**: Split training data into training (80%) and validation (20%) sets.

#### #### 2. Data Augmentation

- **Techniques**: Apply random horizontal flip, normalization, random cropping, random erasing, rotating, shifting, and zooming.
- **Pipelines**: Develop two preprocessing pipelines with different augmentation techniques.

#### #### 3. Model Architecture

- **Simple CNN**: Start with a simple CNN model.
- **Advanced Models**: Consider ResNet18, EfficientNet, or Vision Transformers (ViTs).
- **Hybrid Models**: Explore hybrid models combining CNNs and ViTs.

#### #### 4. Training

- **Optimizer and Loss Function**: Use SGD or Adam optimizer and Cross-Entropy Loss.
- **Epochs and Early Stopping**: Train for up to 150 epochs with early stopping.
- **Monitoring**: Monitor validation accuracy and loss curves.

#### #### 5. Evaluation

- **Accuracy**: Ensure at least 85% accuracy on the test set.
- **Confusion Matrix**: Generate a confusion matrix.
- **Top-1 Accuracy**: Calculate top-1 accuracy.

#### #### 6. Implementation

- **Framework**: Use PyTorch and torchvision.transforms.v2.
- **Resources**: Utilize resources from Kaggle and PapersWithCode.
- **Deployment**: Prepare the model for deployment if needed.

#### #### 7. Visualization and Reporting

- **Confusion Matrix**: Visualize the confusion matrix.
- **Performance Metrics**: Report accuracy and other metrics.
- **Documentation**: Document the entire process.

### ### Detailed Steps for Data Manipulation and Analysis

#### #### 1. Dataset Retrieval and Collection

- **Retrieve Dataset**: Load image paths and labels from 'train\_labels.csv' and 'test\_labels.csv'.

#### #### 2. Checking Dataset Types, Shapes, and Properties

- **Data Types**: Verify image formats and label mappings.
- **Shapes**: Check and convert image dimensions to RGB if needed.
- **Properties**: Ensure balanced dataset and data integrity.

#### #### 3. Preprocessing the Dataset

- **Normalization**: Scale pixel values to [0, 1].
- **Resizing**: Resize images to 224x224 pixels.
- **Splitting**: Split data into training and validation sets.

#### #### 4. Efficient Data Augmentation

- **Techniques**: Apply various augmentation techniques.
- **Pipelines**: Create two augmentation pipelines.

```

#### 5. Extracting and Understanding Underlying Characteristics
- **Class Distribution**: Analyze class distribution.
- **Image Quality**: Assess variations in image quality.
- **Data Insights**: Identify patterns or anomalies.
- **Performance Metrics**: Evaluate model performance using accuracy and confusion matrix.

### Expected Outcomes and Results
- **Quantitative**: Balanced dataset, consistent image size, normalized pixel values, augmented data.
- **Qualitative**: Data integrity, improved model generalization, insightful analysis.

### Responsibilities Fulfillment

#### 1. How to Retrieve or Find High-Performance Model(s)?
- **Initial Models**: Start with simple CNN architectures.
- **Advanced Models**: Experiment with ResNet18, EfficientNet, and Vision Transformers.
- **Hybrid Models**: Explore combinations of CNNs and ViTs.
- **Resources**: Utilize code repositories from Kaggle and PapersWithCode for implementation and fine-tuning.

#### 2. How to Optimize the Hyperparameters of the Retrieved Models?
- **Optimizers**: Use SGD with a learning rate of 0.015 or Adam optimizer.
- **Hyperparameter Tuning**: Experiment with different learning rates, batch sizes, and augmentation techniques.
- **Early Stopping**: Implement early stopping based on validation loss to prevent overfitting.
- **Monitoring**: Continuously monitor validation accuracy and loss curves.

#### 3. How to Extract and Understand the Underlying Characteristics of the Dataset(s)?
- **Class Distribution**: Use 'pandas' to analyze the distribution of classes.
- **Image Quality**: Assess image quality and resolution using 'PIL' or 'OpenCV'.
- **Data Insights**: Perform exploratory data analysis (EDA) to identify patterns or anomalies.
- **Performance Metrics**: Use 'sklearn.metrics' to evaluate model performance and generate confusion matrices.

#### 4. How to Select the Top-k Models or Algorithms Based on the Given Plans?
- **Evaluation Metrics**: Use accuracy and confusion matrix to evaluate models.
- **Top-k Selection**: Rank models based on their performance metrics and select the top-k models that meet or exceed the 85% accuracy requirement.
- **Documentation**: Document the performance of each model, including accuracy, confusion matrix, and any other relevant metrics.

By following this comprehensive plan, you will be able to retrieve high-performance models, optimize their hyperparameters, understand the dataset's characteristics, and select the top-k models based on their performance.

```

## C.4 AGENT-SPECIFIC PLAN EXECUTION OUTCOMES

### Example Results from Prompting-Based Plan Execution by Each Agent

#### Data Agent

```

### Detailed Steps for Data Manipulation and Analysis

#### 1. Dataset Retrieval and Collection

**Retrieve Dataset**:
- **Source**: The dataset is available at '_experiments/datasets/shopee'.
- **Structure**: The directory structure includes 'train', 'test', 'train_labels.csv', and 'test_labels.csv'.
- **Loading Data**: Load the image paths and corresponding labels from 'train_labels.csv' and 'test_labels.csv'. This can be done using libraries like 'pandas' to read the CSV files and 'os' to navigate through the directories.

**Why**: This step ensures that we have access to the correct dataset and that the paths and labels are correctly mapped for further processing.

#### 2. Checking Dataset Types, Shapes, and Properties

**Data Types**:
- **Image Format**: Verify that images are in the correct format (JPEG, PNG). This can be done using libraries like 'PIL' (Python Imaging Library).
- **Labels**: Ensure labels are correctly mapped to images by checking the contents of the CSV files.

```



```

**Shapes**:
- **Image Dimensions**: Check the dimensions of the images. If images are in grayscale,
  convert them to RGB by replicating the single channel across three channels using
  libraries like 'PIL' or 'OpenCV'.

**Properties**:
- **Balanced Dataset**: Ensure the dataset is balanced across the four clothing categories
  by analyzing the label distribution using 'pandas'.
- **Integrity Check**: Verify that there are no missing or corrupted files by attempting to
  open each image file.

**Why**: Ensuring correct data types, shapes, and properties is crucial for accurate model
  training and prevents errors during preprocessing and augmentation.

#### 3. Preprocessing the Dataset

**Normalization**:
- **Scale Pixel Values**: Normalize pixel values to the range [0, 1] using libraries like '
  PIL' or 'torchvision.transforms'.

**Resizing**:
- **Consistent Size**: Resize all images to a consistent size (e.g., 224x224 pixels)
  suitable for model input using 'PIL' or 'torchvision.transforms'.

**Splitting**:
- **Train-Validation Split**: Split the training data into training and validation sets (e.
  g., 80% training, 20% validation) using 'sklearn.model_selection.train_test_split'.

**Why**: Normalization ensures that the pixel values are on a common scale, which helps in
  faster convergence during training. Resizing ensures that all images have the same
  dimensions, which is a requirement for most deep learning models. Splitting the data
  helps in monitoring the model's performance and prevents overfitting.

#### 4. Efficient Data Augmentation

**Techniques**:
- **Random Horizontal Flip**: Randomly flip images horizontally to increase diversity.
- **Normalization**: Normalize images to have zero mean and unit variance.
- **Random Cropping**: Randomly crop images to introduce variability.
- **Random Erasing**: Randomly erase parts of images to make the model robust to occlusions
  .
- **Rotating**: Rotate images randomly to make the model invariant to orientation.
- **Shifting**: Shift images randomly to simulate different viewpoints.
- **Zooming**: Zoom in and out of images to introduce scale variability.

**Pipelines**:
- **Pipeline 1**: Apply random horizontal flip and normalization.
- **Pipeline 2**: Apply random horizontal flip, normalization, random cropping, and random
  erasing.

**Why**: Data augmentation increases the diversity and robustness of the training data,
  which helps in improving the generalization of the model.

#### 5. Extracting and Understanding Underlying Characteristics

**Class Distribution**:
- **Analyze Distribution**: Use 'pandas' to analyze the distribution of classes to ensure a
  balanced dataset.

**Image Quality**:
- **Check Variations**: Assess variations in image quality and resolution using 'PIL' or '
  OpenCV'.

**Data Insights**:
- **Identify Patterns**: Use exploratory data analysis (EDA) techniques to identify any
  patterns or anomalies in the data that could affect model performance.

**Performance Metrics**:
- **Evaluate Model**: Use metrics such as accuracy and confusion matrix to evaluate the
  model's performance and understand misclassification patterns. This can be done using
  'sklearn.metrics'.

**Why**: Understanding the underlying characteristics of the data helps in identifying
  potential issues that could affect model performance. Analyzing class distribution
  ensures that the dataset is balanced, which is crucial for training a fair and
  unbiased model.

### Expected Outcomes and Results

```

```

**Quantitative**:
- **Balanced Dataset**: A dataset with an equal number of images for each clothing category
.
- **Consistent Image Size**: All images resized to 224x224 pixels.
- **Normalized Pixel Values**: Pixel values scaled to the range [0, 1].
- **Augmented Data**: Increased dataset diversity through augmentation techniques.

**Qualitative**:
- **Data Integrity**: No missing or corrupted files.
- **Improved Model Generalization**: Enhanced model robustness due to diverse and augmented training data.
- **Insightful Analysis**: Clear understanding of class distribution, image quality, and potential data patterns or anomalies.

By following these detailed steps, another data scientist should be able to reproduce the results and build a robust image classification model that meets the specified accuracy requirements and provides insightful performance analysis.

```

## Model Agent

### Detailed Steps for Modeling and Optimization

### #### 1. Retrieving High-Performance Models

```

**Initial Models**:
- **Simple CNN**: Start with a basic Convolutional Neural Network (CNN) architecture. This model will serve as a baseline to understand the dataset and the task.

**Advanced Models**:
- **ResNet18**: A well-known deep learning model that has shown excellent performance in image classification tasks. It uses residual connections to ease the training of deeper networks.
- **EfficientNet**: Known for its efficiency in terms of parameter count and FLOPs while maintaining high accuracy. EfficientNet-B0 can be a good starting point.
- **Vision Transformers (ViTs)**: These models have gained popularity for their performance in image classification tasks. A small variant like ViT-Small can be considered.

**Hybrid Models**:
- **CNN + ViT**: Combining the feature extraction capabilities of CNNs with the global context understanding of ViTs. This hybrid approach can potentially leverage the strengths of both architectures.

```

### #### 2. Optimizing Hyperparameters

```

**Optimizers**:
- **SGD**: Stochastic Gradient Descent with a learning rate of 0.015.
- **Adam**: Adaptive Moment Estimation with a learning rate of 0.001.

**Hyperparameter Tuning**:
- **Learning Rate**: Experiment with learning rates (0.001, 0.005, 0.01, 0.015) to find the optimal rate for convergence.
- **Batch Size**: Test batch sizes (16, 32, 64) to balance between memory usage and training speed.
- **Data Augmentation Techniques**: Apply different augmentation techniques and observe their impact on validation accuracy.
- **Number of Epochs**: Train for up to 150 epochs with early stopping based on validation loss to prevent overfitting.
- **Weight Decay**: Regularization parameter for SGD, set to 0.0005.

**Optimal Values**:
- **Learning Rate**: 0.001 for Adam, 0.015 for SGD.
- **Batch Size**: 32.
- **Weight Decay**: 0.0005.
- **Number of Epochs**: Up to 150 with early stopping.

```

### #### 3. Extracting and Understanding Characteristics

```

**Computation Complexity**:
- **Number of Parameters**: Calculate the total number of trainable parameters in each model.
- **FLOPs (Floating Point Operations per Second)**: Measure the computational complexity of each model.

**Memory Usage**:
- **Model Size**: Determine the size of the model in memory (MB).
- **Training Memory Usage**: Monitor GPU memory usage during training.

```

```

**Inference Latency:**
- **Inference Speed**: Measure the time taken for a single forward pass (in milliseconds).

**Why Each Step is Essential:**
- **Model Selection**: Choosing a variety of models ensures that we explore different architectures and their strengths.
- **Hyperparameter Tuning**: Fine-tuning hyperparameters is crucial for achieving optimal performance and ensuring the model generalizes well.
- **Characteristic Analysis**: Understanding the computational and memory requirements helps in selecting models that are not only accurate but also efficient and deployable
.

### Candidate Models and Expected Performance

**1. ResNet18**
- **Accuracy**: Expected to achieve around 88% on the test set.
- **Number of Parameters**: Approximately 11.7 million.
- **FLOPs**: Around 1.8 billion.
- **Model Size**: Approximately 44 MB.
- **Training Time**: Approximately 2 hours on a single GPU.
- **Inference Speed**: Around 10 milliseconds per image.

**2. EfficientNet-B0**
- **Accuracy**: Expected to achieve around 87% on the test set.
- **Number of Parameters**: Approximately 5.3 million.
- **FLOPs**: Around 390 million.
- **Model Size**: Approximately 20 MB.
- **Training Time**: Approximately 3 hours on a single GPU.
- **Inference Speed**: Around 8 milliseconds per image.

**3. Vision Transformer (ViT-Small)**
- **Accuracy**: Expected to achieve around 86% on the test set.
- **Number of Parameters**: Approximately 22 million.
- **FLOPs**: Around 4.6 billion.
- **Model Size**: Approximately 85 MB.
- **Training Time**: Approximately 4 hours on a single GPU.
- **Inference Speed**: Around 15 milliseconds per image.

By following these detailed steps and focusing on the modeling and optimization parts, you will be able to identify and select the top-3 high-performance models for the given image classification task.

```

## D SHOWCASE OF GENERATED PIPELINES

### Generated Pipeline under Constraint-Free Setting by *AutoML-Agent*

**User Prompt:** Please provide a classification model that categorizes images into one of four clothing categories. The image path, along with its label information, can be found in the files `train_labels.csv` and `test_labels.csv`.

#### Output Pipeline:

```

import os
import random
import time
import pandas as pd
from PIL import Image
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
from torchvision import datasets, transforms, models
from torch.utils.data import DataLoader, Dataset
import torch
import torch.nn as nn
import torch.optim as optim
import numpy as np
import gradio as gr

# Define constants
SEED = 42
DATASET_PATH = "_experiments/datasets/shopee"
TRAIN_LABELS_CSV = os.path.join(DATASET_PATH, "train_labels.csv")
TEST_LABELS_CSV = os.path.join(DATASET_PATH, "test_labels.csv")
TRAIN_DIR = os.path.join(DATASET_PATH)

```

```

TEST_DIR = os.path.join(DATASET_PATH)
MODEL_SAVE_PATH = "./agent_workspace/trained_models/resnet50_finetuned.pth"

# Set seeds for reproducibility
random.seed(SEED)
torch.manual_seed(SEED)
np.random.seed(SEED)

# Define device for model operations
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# Custom dataset class
class ShopeeDataset(Dataset):
    def __init__(self, csv_file, root_dir, transform=None):
        self.labels_df = pd.read_csv(csv_file)
        self.root_dir = root_dir
        self.transform = transform

    def __len__(self):
        return len(self.labels_df)

    def __getitem__(self, idx):
        img_name = os.path.join(self.root_dir, self.labels_df.iloc[idx, 0])
        image = Image.open(img_name).convert("RGB")
        label = self.labels_df.iloc[idx, 1]

        if self.transform:
            image = self.transform(image)

        return image, label

# Data preprocessing and feature engineering
def preprocess_data():
    transform = transforms.Compose([
        transforms.Resize((28, 28)),
        transforms.ToTensor(),
        transforms.Normalize((0.5,), (0.5,))
    ])

    train_dataset = ShopeeDataset(csv_file=TRAIN_LABELS_CSV, root_dir=TRAIN_DIR, transform=transform)
    test_dataset = ShopeeDataset(csv_file=TEST_LABELS_CSV, root_dir=TEST_DIR, transform=transform)

    train_size = int(0.7 * len(train_dataset))
    valid_size = int(0.2 * len(train_dataset))
    test_size = len(train_dataset) - train_size - valid_size

    train_dataset, valid_dataset, test_dataset = torch.utils.data.random_split(train_dataset, [train_size, valid_size, test_size])

    train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
    valid_loader = DataLoader(valid_dataset, batch_size=32, shuffle=False)
    test_loader = DataLoader(test_dataset, batch_size=32, shuffle=False)

    return train_loader, valid_loader, test_loader

# Model training
def train_model(model, train_loader, valid_loader):
    criterion = nn.CrossEntropyLoss()
    optimizer = optim.Adam(model.parameters(), lr=0.00001)

    model.to(device)
    best_accuracy = 0.0

    for epoch in range(100):
        model.train()
        running_loss = 0.0
        for inputs, labels in train_loader:
            inputs, labels = inputs.to(device), labels.to(device)

            optimizer.zero_grad()
            outputs = model(inputs)
            loss = criterion(outputs, labels)
            loss.backward()
            optimizer.step()

            running_loss += loss.item() * inputs.size(0)

```

```

epoch_loss = running_loss / len(train_loader.dataset)

model.eval()
correct = 0
total = 0
with torch.no_grad():
    for inputs, labels in valid_loader:
        inputs, labels = inputs.to(device), labels.to(device)
        outputs = model(inputs)
        _, predicted = torch.max(outputs, 1)
        total += labels.size(0)
        correct += (predicted == labels).sum().item()

accuracy = correct / total
if accuracy > best_accuracy:
    best_accuracy = accuracy
    torch.save(model.state_dict(), MODEL_SAVE_PATH)

print(f"Epoch {epoch+1}/{100}, Loss: {epoch_loss:.4f}, Validation Accuracy: {accuracy:.4f}")

model.load_state_dict(torch.load(MODEL_SAVE_PATH))
return model

# Model evaluation
def evaluate_model(model, test_loader):
    model.eval()
    y_true = []
    y_pred = []

    with torch.no_grad():
        for inputs, labels in test_loader:
            inputs, labels = inputs.to(device), labels.to(device)
            outputs = model(inputs)
            _, predicted = torch.max(outputs, 1)
            y_true.extend(labels.cpu().numpy())
            y_pred.extend(predicted.cpu().numpy())

    performance_scores = {
        'Accuracy': accuracy_score(y_true, y_pred),
        'Precision': precision_score(y_true, y_pred, average='weighted'),
        'Recall': recall_score(y_true, y_pred, average='weighted'),
        'F1': f1_score(y_true, y_pred, average='weighted')
    }

    model_size = os.path.getsize(MODEL_SAVE_PATH) / (1024 * 1024) # in MB
    start_time = time.time()
    for _ in range(100):
        _ = model(inputs)
    inference_time = (time.time() - start_time) / 100 # average inference time

    complexity_scores = {
        'Model Size (MB)': model_size,
        'Inference Time (ms)': inference_time * 1000
    }

    return performance_scores, complexity_scores

# Prepare model for deployment
def prepare_model_for_deployment(model):
    model.eval()
    scripted_model = torch.jit.script(model)
    scripted_model.save(MODEL_SAVE_PATH.replace(".pth", "_scripted.pt"))
    return scripted_model

# Deploy model using Gradio
def deploy_model(model):
    def predict(image):
        image = transforms.Compose([
            transforms.Resize((28, 28)),
            transforms.ToTensor(),
            transforms.Normalize((0.5,), (0.5,))
        ])(image).unsqueeze(0).to(device)

        with torch.no_grad():
            outputs = model(image)
            _, predicted = torch.max(outputs, 1)

        return predicted.item()

```

```

interface = gr.Interface(fn=predict, inputs=gr.Image(type="pil"), outputs="label")
url_endpoint = interface.launch(share=True)
return url_endpoint

# The main function to orchestrate the data loading, data preprocessing, feature
# engineering, model training, model preparation, model deployment, and model evaluation
def main():
    # Step 1: Retrieve or load dataset
    train_loader, valid_loader, test_loader = preprocess_data()

    # Step 2: Define model
    model = models.resnet50(pretrained=True)
    num_fters = model.fc.in_features
    model.fc = nn.Linear(num_fters, 4) # Assuming 4 classes

    # Step 3: Train model
    model = train_model(model, train_loader, valid_loader)

    # Step 4: Evaluate model
    model_performance, model_complexity = evaluate_model(model, test_loader)

    # Step 5: Prepare model for deployment
    deployable_model = prepare_model_for_deployment(model)

    # Step 6: Deploy model
    url_endpoint = deploy_model(deployable_model)

    return train_loader, valid_loader, test_loader, model, deployable_model, url_endpoint,
        model_performance, model_complexity

if __name__ == "__main__":
    train_loader, valid_loader, test_loader, model, deployable_model, url_endpoint,
        model_performance, model_complexity = main()
    print("Model Performance on Test Set:", model_performance)
    print("Model Complexity:", model_complexity)
    print("Deployment URL:", url_endpoint)

```

### Generated Pipeline for Constraint-Aware Setting by *AutoML-Agent*

**User Prompt:** Please provide a classification model that categorizes images into one of four clothing categories. The image path, along with its label information, can be found in the files `train_labels.csv` and `test_labels.csv`. The model should achieve at least 85% accuracy on the test set and be implemented using PyTorch. Additionally, please include data augmentation techniques and a confusion matrix in the evaluation.

#### Output Pipeline:

```

import os
import random
import time
import pandas as pd
import numpy as np
import torch
import torch.nn as nn
import torch.optim as optim
from torchvision import datasets, transforms, models
from torch.utils.data import DataLoader, Dataset
from PIL import Image
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, f1_score, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt
import gradio as gr

# Set random seeds for reproducibility
SEED = 42
random.seed(SEED)
torch.manual_seed(SEED)
np.random.seed(SEED)

# Define device for model operations
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

```

```

DATASET_PATH = "_experiments/datasets/shopee" # path for saving and loading dataset(s)

# Custom Dataset class
class ShopeeDataset(Dataset):
    def __init__(self, dataframe, transform=None):
        self.dataframe = dataframe
        self.transform = transform

    def __len__(self):
        return len(self.dataframe)

    def __getitem__(self, idx):
        img_path = self.dataframe.iloc[idx]['image']
        label = self.dataframe.iloc[idx]['label']
        image = Image.open(img_path).convert('RGB')
        if self.transform:
            image = self.transform(image)
        return image, label

# Data preprocessing and feature engineering
def preprocess_data():
    train_labels = pd.read_csv(os.path.join(DATASET_PATH, 'train_labels.csv'))
    test_labels = pd.read_csv(os.path.join(DATASET_PATH, 'test_labels.csv'))

    train_labels['image'] = train_labels['image'].apply(lambda x: os.path.join(DATASET_PATH,
x))
    test_labels['image'] = test_labels['image'].apply(lambda x: os.path.join(DATASET_PATH, x
))

    # Split the data
    train_data, val_data = train_test_split(train_labels, test_size=0.2, stratify=
train_labels['label'])
    val_data, test_data = train_test_split(val_data, test_size=0.5, stratify=val_data['label
'])

    # Define transformations
    train_transforms = transforms.Compose([
        transforms.RandomResizedCrop(224),
        transforms.RandomHorizontalFlip(),
        transforms.RandomRotation(10),
        transforms.ToTensor(),
        transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
    ])

    val_test_transforms = transforms.Compose([
        transforms.Resize(256),
        transforms.CenterCrop(224),
        transforms.ToTensor(),
        transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
    ])

    # Load datasets
    train_dataset = ShopeeDataset(train_data, transform=train_transforms)
    val_dataset = ShopeeDataset(val_data, transform=val_test_transforms)
    test_dataset = ShopeeDataset(test_data, transform=val_test_transforms)

    train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
    val_loader = DataLoader(val_dataset, batch_size=32, shuffle=False)
    test_loader = DataLoader(test_dataset, batch_size=32, shuffle=False)

    return train_loader, val_loader, test_loader

def train_model(model, train_loader, val_loader):
    criterion = nn.CrossEntropyLoss()
    optimizer = optim.Adam(model.parameters(), lr=0.0005, weight_decay=0.0001)

    num_epochs = 100
    best_val_loss = float('inf')
    early_stop_count = 0
    early_stop_patience = 10

    for epoch in range(num_epochs):
        model.train()
        running_loss = 0.0
        for inputs, labels in train_loader:
            inputs, labels = inputs.to(device), labels.to(device)
            optimizer.zero_grad()
            outputs = model(inputs)
            loss = criterion(outputs, labels)

```

```

        loss.backward()
        optimizer.step()
        running_loss += loss.item() * inputs.size(0)

    epoch_loss = running_loss / len(train_loader.dataset)

    # Validation
    model.eval()
    val_loss = 0.0
    with torch.no_grad():
        for inputs, labels in val_loader:
            inputs, labels = inputs.to(device), labels.to(device)
            outputs = model(inputs)
            loss = criterion(outputs, labels)
            val_loss += loss.item() * inputs.size(0)

    val_loss = val_loss / len(val_loader.dataset)

    print(f'Epoch {epoch+1}/{num_epochs}, Train Loss: {epoch_loss:.4f}, Val Loss: {
        val_loss:.4f}')

    # Early stopping
    if val_loss < best_val_loss:
        best_val_loss = val_loss
        torch.save(model.state_dict(), './agent_workspace/trained_models/best_model.pth')
        early_stop_count = 0
    else:
        early_stop_count += 1
        if early_stop_count >= early_stop_patience:
            print("Early stopping")
            break

    return model

def evaluate_model(model, test_loader):
    model.load_state_dict(torch.load('./agent_workspace/trained_models/best_model.pth'))
    model.eval()

    all_preds = []
    all_labels = []
    with torch.no_grad():
        for inputs, labels in test_loader:
            inputs, labels = inputs.to(device), labels.to(device)
            outputs = model(inputs)
            _, preds = torch.max(outputs, 1)
            all_preds.extend(preds.cpu().numpy())
            all_labels.extend(labels.cpu().numpy())

    accuracy = accuracy_score(all_labels, all_preds)
    f1 = f1_score(all_labels, all_preds, average='weighted')

    performance_scores = {
        'ACC': accuracy,
        'F1': f1
    }

    # Measure model complexity
    start_time = time.time()
    with torch.no_grad():
        for inputs, _ in test_loader:
            inputs = inputs.to(device)
            _ = model(inputs)
    inference_time = (time.time() - start_time) / len(test_loader.dataset)

    model_size = os.path.getsize('./agent_workspace/trained_models/best_model.pth') / 1e6 #
        in MB

    complexity_scores = {
        'Inference Time': inference_time,
        'Model Size (MB)': model_size
    }

    # Generate confusion matrix
    conf_matrix = confusion_matrix(all_labels, all_preds)
    plt.figure(figsize=(10, 7))
    sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues')
    plt.xlabel('Predicted')
    plt.ylabel('True')
    plt.title('Confusion Matrix')

```



```

plt.show()

return performance_scores, complexity_scores

def prepare_model_for_deployment(model):
    # Assuming the model is already saved in the best_model.pth file
    model.load_state_dict(torch.load('./agent_workspace/trained_models/best_model.pth'))
    model.eval()
    return model

def deploy_model(model):
    def predict(image):
        transform = transforms.Compose([
            transforms.Resize(256),
            transforms.CenterCrop(224),
            transforms.ToTensor(),
            transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
        ])
        image = transform(image).unsqueeze(0).to(device)
        with torch.no_grad():
            output = model(image)
            _, pred = torch.max(output, 1)
        return pred.item()

    interface = gr.Interface(fn=predict, inputs=gr.Image(type="pil"), outputs="label")
    url_endpoint = interface.launch(share=True)
    return url_endpoint

# The main function to orchestrate the data loading, data preprocessing, feature
# engineering, model training, model preparation, model deployment, and model evaluation
def main():
    """
    Main function to execute the image classification pipeline.
    """

    # Step 1. Retrieve or load a dataset from hub (if available) or user's local storage (if
    # given)
    train_loader, val_loader, test_loader = preprocess_data()

    # Step 2. Define required model
    model = models.resnet18(weights=models.ResNet18_Weights.IMAGENET1K_V1)
    num_fters = model.fc.in_features
    model.fc = nn.Linear(num_fters, 4)
    model = model.to(device)

    # Step 3. Train the retrieved/loaded model
    model = train_model(model, train_loader, val_loader)

    # Step 4. Evaluate the trained model
    model_performance, model_complexity = evaluate_model(model, test_loader)

    # Step 5. Prepare the model for deployment
    deployable_model = prepare_model_for_deployment(model)

    # Step 6. Deploy the model
    url_endpoint = deploy_model(deployable_model)

    return (
        train_loader,
        val_loader,
        test_loader,
        model,
        deployable_model,
        url_endpoint,
        model_performance,
        model_complexity,
    )

if __name__ == "__main__":
    train_loader, val_loader, test_loader, model, deployable_model, url_endpoint,
    model_performance, model_complexity = main()
    print("Model Performance on Test Set:", model_performance)
    print("Model Complexity:", model_complexity)
    print("Gradio URL Endpoint:", url_endpoint)

```

## E DETAILED EXPERIMENTAL RESULTS

This section presents the detailed experimental results referenced in §4. Tables 5, 6, and 7 show the average scores with standard deviations, computed from five independent runs. Tables 8 and 9 report the results of the ablation and hyperparameter studies, respectively. Finally, Table 10 reports the average time and money used to generate the final code in a single run.

Table 5: Performance comparison with the **SR** metric. Best results are highlighted in **bold**.

| Method                  | Image Classification |              | Text Classification |              | Tabular Classification |              | Tabular Regression |              | Tabular Clustering |              | Time-Series Forecasting |              | Node Classification |              | Avg.         |
|-------------------------|----------------------|--------------|---------------------|--------------|------------------------|--------------|--------------------|--------------|--------------------|--------------|-------------------------|--------------|---------------------|--------------|--------------|
|                         | Butterfly            | Shopee       | Ecomm               | Entail       | Banana                 | Software     | Crab               | Crop         | Smoker             | Student      | Weather                 | Electricity  | Cora                | Citeseer     |              |
| <i>Constraint-Free</i>  |                      |              |                     |              |                        |              |                    |              |                    |              |                         |              |                     |              |              |
| GPT-3.5                 | 0.000                | 0.000        | 0.000               | 0.000        | 0.300                  | 0.100        | 0.000              | 0.000        | 0.400              | 0.000        | 0.000                   | 0.000        | 0.000               | 0.000        | 0.057        |
|                         | (±0.000)             | (±0.000)     | (±0.000)            | (±0.000)     | (±0.274)               | (±0.224)     | (±0.000)           | (±0.000)     | (±0.224)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.000)            | (±0.000)     | (±0.052)     |
| GPT-4                   | 0.200                | 0.600        | 0.000               | 0.400        | 0.400                  | 0.400        | 0.400              | 0.600        | 0.600              | 0.000        | 0.000                   | 0.000        | 0.400               | 0.400        | 0.314        |
|                         | (±0.447)             | (±0.548)     | (±0.000)            | (±0.548)     | (±0.548)               | (±0.418)     | (±0.548)           | (±0.548)     | (±0.418)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.548)            | (±0.548)     | (±0.366)     |
| DS-Agent                | 0.400                | 0.800        | 0.000               | 0.700        | 0.800                  | 0.800        | 0.000              | 0.800        | 0.900              | 0.000        | 0.000                   | 0.000        | 0.600               | 0.600        | 0.457        |
|                         | (±0.548)             | (±0.447)     | (±0.000)            | (±0.447)     | (±0.447)               | (±0.274)     | (±0.000)           | (±0.447)     | (±0.224)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.548)            | (±0.548)     | (±0.281)     |
| AutoML-Agent            | <b>1.000</b>         | <b>1.000</b> | <b>1.000</b>        | <b>1.000</b> | <b>1.000</b>           | <b>1.000</b> | <b>1.000</b>       | <b>1.000</b> | <b>1.000</b>       | <b>1.000</b> | <b>1.000</b>            | <b>1.000</b> | <b>1.000</b>        | <b>1.000</b> | <b>1.000</b> |
|                         | (±0.000)             | (±0.000)     | (±0.000)            | (±0.000)     | (±0.000)               | (±0.000)     | (±0.000)           | (±0.000)     | (±0.000)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.000)            | (±0.000)     | (±0.000)     |
| <i>Constraint-Aware</i> |                      |              |                     |              |                        |              |                    |              |                    |              |                         |              |                     |              |              |
| GPT-3.5                 | 0.000                | 0.050        | 0.000               | 0.000        | 0.050                  | 0.150        | 0.100              | 0.050        | 0.150              | 0.000        | 0.000                   | 0.000        | 0.000               | 0.000        | 0.039        |
|                         | (±0.000)             | (±0.112)     | (±0.000)            | (±0.000)     | (±0.112)               | (±0.137)     | (±0.137)           | (±0.112)     | (±0.137)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.000)            | (±0.000)     | (±0.053)     |
| GPT-4                   | 0.150                | 0.350        | 0.200               | 0.200        | 0.150                  | 0.000        | 0.650              | 0.100        | 0.400              | 0.000        | 0.000                   | 0.000        | 0.400               | 0.500        | 0.221        |
|                         | (±0.335)             | (±0.487)     | (±0.447)            | (±0.447)     | (±0.335)               | (±0.000)     | (±0.418)           | (±0.224)     | (±0.335)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.335)            | (±0.354)     | (±0.266)     |
| DS-Agent                | 0.300                | 0.350        | 0.000               | 0.200        | 0.600                  | 0.650        | 0.200              | 0.200        | 0.450              | 0.000        | 0.000                   | 0.150        | 0.200               | 0.450        | 0.268        |
|                         | (±0.411)             | (±0.487)     | (±0.000)            | (±0.326)     | (±0.335)               | (±0.487)     | (±0.447)           | (±0.274)     | (±0.274)           | (±0.000)     | (±0.000)                | (±0.335)     | (±0.326)            | (±0.411)     | (±0.306)     |
| AutoML-Agent            | <b>0.800</b>         | <b>1.000</b> | <b>1.000</b>        | <b>1.000</b> | <b>0.750</b>           | <b>1.000</b> | <b>1.000</b>       | <b>0.750</b> | <b>0.750</b>       | <b>0.750</b> | <b>0.900</b>            | <b>1.000</b> | <b>0.750</b>        | <b>0.750</b> | <b>0.871</b> |
|                         | (±0.112)             | (±0.000)     | (±0.000)            | (±0.000)     | (±0.000)               | (±0.000)     | (±0.000)           | (±0.000)     | (±0.000)           | (±0.000)     | (±0.224)                | (±0.000)     | (±0.000)            | (±0.000)     | (±0.024)     |

Table 6: Performance comparison with the **NPS** metric. Best results are highlighted in **bold**.

| Method                  | Image Classification |              | Text Classification |              | Tabular Classification |              | Tabular Regression |              | Tabular Clustering |              | Time-Series Forecasting |              | Node Classification |              | Avg.         |
|-------------------------|----------------------|--------------|---------------------|--------------|------------------------|--------------|--------------------|--------------|--------------------|--------------|-------------------------|--------------|---------------------|--------------|--------------|
|                         | Butterfly            | Shopee       | Ecomm               | Entail       | Banana                 | Software     | Crab               | Crop         | Smoker             | Student      | Weather                 | Electricity  | Cora                | Citeseer     |              |
| <i>Constraint-Free</i>  |                      |              |                     |              |                        |              |                    |              |                    |              |                         |              |                     |              |              |
| Human Models            | <b>0.931</b>         | 0.921        | 0.935               | 0.664        | 0.976                  | <b>0.669</b> | 0.328              | 0.476        | 0.513              | 0.750        | 0.970                   | 0.916        | 0.811               | <b>0.702</b> | 0.754        |
|                         | (±0.002)             | (±0.012)     | (±0.000)            | (±0.039)     | (±0.000)               | (±0.000)     | (±0.000)           | (±0.000)     | (±0.000)           | (±0.000)     | (±0.000)                | (±0.005)     | (±0.005)            | (±0.006)     | (±0.005)     |
| AutoGluon               | 0.014                | <b>0.988</b> | <b>0.987</b>        | <b>0.807</b> | 0.980                  | 0.524        | <b>0.875</b>       | <b>0.479</b> | N/A                | N/A          | 0.992                   | 0.908        | N/A                 | N/A          | 0.755        |
|                         | (±0.000)             | (±0.000)     | (±0.000)            | (±0.000)     | (±0.000)               | (±0.000)     | (±0.000)           | (±0.000)     | (±0.000)           | (±0.000)     | (±0.000)                | (±0.002)     | (±0.000)            | (±0.000)     | (±0.000)     |
| GPT-3.5                 | 0.000                | 0.000        | 0.000               | 0.000        | 0.587                  | 0.094        | 0.000              | 0.000        | 0.447              | 0.000        | 0.000                   | 0.000        | 0.000               | 0.000        | 0.081        |
|                         | (±0.000)             | (±0.000)     | (±0.000)            | (±0.000)     | (±0.535)               | (±0.209)     | (±0.000)           | (±0.000)     | (±0.251)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.000)            | (±0.000)     | (±0.071)     |
| GPT-4                   | 0.169                | 0.545        | 0.000               | 0.196        | 0.390                  | 0.285        | 0.328              | 0.270        | 0.471              | 0.000        | 0.000                   | 0.000        | 0.186               | 0.199        | 0.217        |
|                         | (±0.379)             | (±0.499)     | (±0.000)            | (±0.295)     | (±0.534)               | (±0.261)     | (±0.450)           | (±0.247)     | (±0.264)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.343)            | (±0.328)     | (±0.257)     |
| DS-Agent                | 0.305                | 0.735        | 0.000               | 0.500        | 0.766                  | 0.523        | 0.000              | 0.431        | 0.504              | 0.000        | 0.000                   | 0.000        | 0.474               | 0.393        | 0.331        |
|                         | (±0.419)             | (±0.411)     | (±0.000)            | (±0.380)     | (±0.428)               | (±0.131)     | (±0.000)           | (±0.324)     | (±0.001)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.433)            | (±0.360)     | (±0.206)     |
| AutoML-Agent            | 0.924                | 0.945        | 0.971               | 0.803        | <b>0.987</b>           | 0.664        | 0.859              | 0.465        | <b>0.521</b>       | <b>0.760</b> | <b>0.995</b>            | <b>0.937</b> | <b>0.831</b>        | 0.592        | <b>0.804</b> |
|                         | (±0.020)             | (±0.043)     | (±0.007)            | (±0.006)     | (±0.019)               | (±0.174)     | (±0.003)           | (±0.020)     | (±0.038)           | (±0.021)     | (±0.003)                | (±0.093)     | (±0.020)            | (±0.015)     | (±0.035)     |
| <i>Constraint-Aware</i> |                      |              |                     |              |                        |              |                    |              |                    |              |                         |              |                     |              |              |
| GPT-3.5                 | 0.000                | 0.173        | 0.000               | 0.000        | 0.196                  | 0.475        | 0.356              | 0.081        | 0.338              | 0.000        | 0.000                   | 0.000        | 0.000               | 0.000        | 0.116        |
|                         | (±0.000)             | (±0.386)     | (±0.000)            | (±0.000)     | (±0.439)               | (±0.476)     | (±0.488)           | (±0.181)     | (±0.309)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.000)            | (±0.000)     | (±0.163)     |
| GPT-4                   | 0.157                | 0.335        | 0.197               | 0.064        | 0.153                  | 0.000        | 0.719              | 0.091        | 0.463              | 0.000        | 0.000                   | 0.000        | 0.637               | 0.564        | 0.241        |
|                         | (±0.350)             | (±0.463)     | (±0.440)            | (±0.144)     | (±0.342)               | (±0.000)     | (±0.405)           | (±0.204)     | (±0.260)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.356)            | (±0.318)     | (±0.234)     |
| DS-Agent                | 0.330                | 0.353        | 0.000               | 0.205        | 0.776                  | 0.383        | 0.173              | 0.183        | 0.505              | 0.000        | 0.000                   | 0.093        | 0.319               | 0.403        | 0.266        |
|                         | (±0.451)             | (±0.485)     | (±0.000)            | (±0.301)     | (±0.434)               | (±0.214)     | (±0.386)           | (±0.409)     | (±0.001)           | (±0.000)     | (±0.000)                | (±0.209)     | (±0.437)            | (±0.369)     | (±0.264)     |
| AutoML-Agent            | <b>0.926</b>         | <b>0.972</b> | <b>0.982</b>        | <b>0.796</b> | <b>0.967</b>           | <b>0.573</b> | <b>0.861</b>       | <b>0.473</b> | <b>0.582</b>       | <b>0.769</b> | <b>0.982</b>            | <b>0.978</b> | <b>0.843</b>        | <b>0.632</b> | <b>0.810</b> |
|                         | (±0.015)             | (±0.022)     | (±0.002)            | (±0.027)     | (±0.002)               | (±0.142)     | (±0.002)           | (±0.020)     | (±0.042)           | (±0.010)     | (±0.028)                | (±0.001)     | (±0.034)            | (±0.037)     | (±0.027)     |

Table 7: Performance comparison with the **CS** metric. Best results are highlighted in **bold**.

| Method                  | Image Classification |              | Text Classification |              | Tabular Classification |              | Tabular Regression |              | Tabular Clustering |              | Time-Series Forecasting |              | Node Classification |              | Avg.         |
|-------------------------|----------------------|--------------|---------------------|--------------|------------------------|--------------|--------------------|--------------|--------------------|--------------|-------------------------|--------------|---------------------|--------------|--------------|
|                         | Butterfly            | Shopee       | Ecomm               | Entail       | Banana                 | Software     | Crab               | Crop         | Smoker             | Student      | Weather                 | Electricity  | Cora                | Citeseer     |              |
| <i>Constraint-Free</i>  |                      |              |                     |              |                        |              |                    |              |                    |              |                         |              |                     |              |              |
| GPT-3.5                 | 0.000                | 0.000        | 0.000               | 0.000        | 0.443                  | 0.097        | 0.000              | 0.000        | 0.424              | 0.000        | 0.000                   | 0.000        | 0.000               | 0.000        | 0.069        |
|                         | (±0.000)             | (±0.000)     | (±0.000)            | (±0.000)     | (±0.405)               | (±0.216)     | (±0.000)           | (±0.000)     | (±0.237)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.000)            | (±0.000)     | (±0.061)     |
| GPT-4                   | 0.185                | 0.573        | 0.000               | 0.298        | 0.395                  | 0.343        | 0.364              | 0.435        | 0.536              | 0.000        | 0.000                   | 0.000        | 0.293               | 0.299        | 0.266        |
|                         | (±0.413)             | (±0.523)     | (±0.000)            | (±0.413)     | (±0.541)               | (±0.329)     | (±0.499)           | (±0.397)     | (±0.325)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.417)            | (±0.420)     | (±0.305)     |
| DS-Agent                | 0.352                | 0.768        | 0.000               | 0.600        | 0.783                  | 0.661        | 0.000              | 0.616        | 0.702              | 0.000        | 0.000                   | 0.000        | 0.537               | 0.496        | 0.394        |
|                         | (±0.483)             | (±0.429)     | (±0.000)            | (±0.353)     | (±0.438)               | (±0.172)     | (±0.000)           | (±0.361)     | (±0.111)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.490)            | (±0.453)     | (±0.235)     |
| AutoML-Agent            | <b>0.962</b>         | <b>0.973</b> | <b>0.985</b>        | <b>0.901</b> | <b>0.993</b>           | <b>0.832</b> | <b>0.929</b>       | <b>0.732</b> | <b>0.761</b>       | <b>0.880</b> | <b>0.998</b>            | <b>0.969</b> | <b>0.915</b>        | <b>0.796</b> | <b>0.902</b> |
|                         | (±0.010)             | (±0.021)     | (±0.004)            | (±0.003)     | (±0.010)               | (±0.087)     | (±0.001)           | (±0.010)     | (±0.019)           | (±0.010)     | (±0.002)                | (±0.047)     | (±0.010)            | (±0.007)     | (±0.017)     |
| <i>Constraint-Aware</i> |                      |              |                     |              |                        |              |                    |              |                    |              |                         |              |                     |              |              |
| GPT-3.5                 | 0.000                | 0.111        | 0.000               | 0.000        | 0.123                  | 0.312        | 0.228              | 0.066        | 0.244              | 0.000        | 0.000                   | 0.000        | 0.000               | 0.000        | 0.077        |
|                         | (±0.000)             | (±0.249)     | (±0.000)            | (±0.000)     | (±0.276)               | (±0.302)     | (±0.312)           | (±0.147)     | (±0.223)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.000)            | (±0.000)     | (±0.108)     |
| GPT-4                   | 0.153                | 0.343        | 0.198               | 0.132        | 0.151                  | 0.000        | 0.685              | 0.096        | 0.432              | 0.000        | 0.000                   | 0.000        | 0.518               | 0.532        | 0.231        |
|                         | (±0.343)             | (±0.475)     | (±0.444)            | (±0.296)     | (±0.339)               | (±0.000)     | (±0.394)           | (±0.214)     | (±0.270)           | (±0.000)     | (±0.000)                | (±0.000)     | (±0.317)            | (±0.319)     | (±0.244)     |
| DS-Agent                | 0.315                | 0.351        | 0.000               | 0.203        | 0.688                  | 0.516        | 0.186              | 0.191        | 0.477              | 0.000        | 0.000                   | 0.122        | 0.260               | 0.427        | 0.267        |
|                         | (±0.431)             | (±0.485)     | (±0.000)            | (±0.312)     | (±0.385)               | (±0.332)     | (±0.417)           | (±0.428)     | (±0.137)           | (±0.000)     | (±0.000)                | (±0.272)     | (±0.367)            | (±0.390)     | (±0.283)     |
| AutoML-Agent            | <b>0.863</b>         | <b>0.986</b> | <b>0.991</b>        | <b>0.898</b> | <b>0.858</b>           | <b>0.786</b> | <b>0.930</b>       | <b>0.611</b> | <b>0.666</b>       | <b>0.760</b> | <b>0.941</b>            | <b>0.989</b> | <b>0.796</b>        | <b>0.691</b> | <b>0.841</b> |
|                         | (±0.063)             | (±0.011)     | (±0.001)            | (±0.013)     | (±0.001)               | (±0.071)     | (±0.001)           | (±0.010)     | (±0.021)           | (±0.005)     | (±0.126)                | (±0.001)     | (±0.017)            | (±0.018)     | (±0.026)     |

Table 8: Results of ablation study on different variations.

| RAP                                 | Plan Decomposition | Multi-Step Verification | Image Classification | Text Classification | Tabular Classification | Time-Series Forecasting | Node Classification | Average      |
|-------------------------------------|--------------------|-------------------------|----------------------|---------------------|------------------------|-------------------------|---------------------|--------------|
| <i>Success Rate</i>                 |                    |                         |                      |                     |                        |                         |                     |              |
| ✓                                   |                    |                         | 1.000                | 0.000               | 0.000                  | 0.000                   | 1.000               | 0.400        |
| ✓                                   | ✓                  |                         | 1.000                | 1.000               | 1.000                  | 0.000                   | 1.000               | 0.800        |
| ✓                                   | ✓                  | ✓                       | <b>1.000</b>         | <b>1.000</b>        | <b>1.000</b>           | <b>1.000</b>            | <b>1.000</b>        | <b>1.000</b> |
| <i>Normalized Performance Score</i> |                    |                         |                      |                     |                        |                         |                     |              |
| ✓                                   |                    |                         | 0.929                | 0.000               | 0.000                  | 0.000                   | 0.734               | 0.333        |
| ✓                                   | ✓                  |                         | 0.928                | <b>0.982</b>        | 0.975                  | 0.000                   | 0.748               | 0.727        |
| ✓                                   | ✓                  | ✓                       | <b>0.936</b>         | 0.971               | <b>1.000</b>           | <b>0.991</b>            | <b>0.812</b>        | <b>0.942</b> |
| <i>Comprehensive Score</i>          |                    |                         |                      |                     |                        |                         |                     |              |
| ✓                                   |                    |                         | 0.965                | 0.000               | 0.000                  | 0.000                   | 0.867               | 0.366        |
| ✓                                   | ✓                  |                         | 0.964                | <b>0.991</b>        | 0.988                  | 0.000                   | 0.874               | 0.763        |
| ✓                                   | ✓                  | ✓                       | <b>0.968</b>         | 0.986               | <b>1.000</b>           | <b>0.996</b>            | <b>0.906</b>        | <b>0.971</b> |

Table 9: Comparison between the different numbers of plans.

| Number of Plans                     | Image Classification | Text Classification | Tabular Classification | Time-Series Forecasting | Node Classification | Average      |
|-------------------------------------|----------------------|---------------------|------------------------|-------------------------|---------------------|--------------|
| <i>Success Rate</i>                 |                      |                     |                        |                         |                     |              |
| 1                                   | 1.000                | 1.000               | 1.000                  | 1.000                   | 1.000               | 1.000        |
| 3                                   | 1.000                | 1.000               | 1.000                  | 1.000                   | 1.000               | 1.000        |
| 5                                   | 1.000                | 1.000               | 1.000                  | 1.000                   | 1.000               | 1.000        |
| <i>Normalized Performance Score</i> |                      |                     |                        |                         |                     |              |
| 1                                   | 0.831                | 0.966               | 0.958                  | 0.998                   | 0.800               | 0.911        |
| 3                                   | <b>0.936</b>         | <b>0.971</b>        | <b>1.000</b>           | <b>0.999</b>            | <b>0.812</b>        | <b>0.944</b> |
| 5                                   | 0.916                | 0.964               | 0.973                  | 0.998                   | 0.805               | 0.931        |
| <i>Comprehensive Score</i>          |                      |                     |                        |                         |                     |              |
| 1                                   | 0.916                | 0.983               | 0.979                  | 0.999                   | 0.900               | 0.955        |
| 3                                   | 0.968                | 0.986               | 1.000                  | 0.999                   | 0.906               | 0.972        |
| 5                                   | 0.958                | 0.982               | 0.986                  | 0.999                   | 0.903               | 0.966        |

Table 10: Time and monetary costs averaged across different tasks and datasets for a single run under the constraint-free and constraint-aware settings.

| Cost                    | Prompt Parsing | Request Verification | Retrieval & Planning | Plan Execution | Execution Verification | Selection and Summarization | Code Generation | Total  |
|-------------------------|----------------|----------------------|----------------------|----------------|------------------------|-----------------------------|-----------------|--------|
| <i>Constraint-Free</i>  |                |                      |                      |                |                        |                             |                 |        |
| Money (\$)              | N/A            | 0.00                 | 0.02                 | 0.14           | 0.00                   | 0.06                        | 0.04            | 0.27   |
| Time (s)                | 10.78          | 1.91                 | 187.71               | 136.34         | 1.04                   | 17.88                       | 182.60          | 538.25 |
| <i>Constraint-Aware</i> |                |                      |                      |                |                        |                             |                 |        |
| Money (\$)              | N/A            | 0.00                 | 0.00                 | 0.11           | 0.00                   | 0.15                        | 0.06            | 0.32   |
| Time (s)                | 14.21          | 3.63                 | 182.38               | 98.62          | 1.37                   | 20.25                       | 191.90          | 512.35 |