# 🎊 Agentic Predictor: Performance Prediction for Agentic Workflows via Multi-View Encoding

**Patara Trirat[1], Wonyong Jeong[1], Sung Ju Hwang[1,2]**     *{patara, young, sjhwang}@deepauto.ai*
*[1]DeepAuto.ai, [2]KAIST*
*Seoul, South Korea*

## Abstract

Large language models (LLMs) have demonstrated remarkable capabilities across diverse tasks, but optimizing LLM-based agentic systems remains challenging due to the vast search space of agent configurations, prompting strategies, and communication patterns. Existing approaches often rely on heuristic-based tuning or exhaustive evaluation, which can be computationally expensive and suboptimal. This paper proposes **Agentic Predictor**, a lightweight predictor for efficient agentic workflow evaluation. Agentic Predictor is equipped with a *multi-view workflow encoding* technique that leverages multi-view representation learning of agentic systems by incorporating code architecture, textual prompts, and interaction graph features. To achieve high predictive accuracy while significantly reducing the number of required workflow evaluations for training a predictor, Agentic Predictor employs *cross-domain unsupervised pretraining*. By learning to approximate task success rates, Agentic Predictor enables fast and accurate selection of optimal agentic workflow configurations for a given task, significantly reducing the need for expensive trial-and-error evaluations. Experiments on a carefully curated benchmark spanning three domains show that our predictor outperforms state-of-the-art methods in both predictive accuracy and workflow utility, highlighting the potential of performance predictors in streamlining the design of LLM-based agentic workflows.

## 1 Introduction

Large language models (LLMs) have catalyzed the development of agentic systems capable of executing complex, multi-step tasks autonomously (Hong et al., 2024; Wu et al., 2024; Xi et al., 2024; Mialon et al., 2024). These systems, often constructed through meticulous manual engineering, integrate components such as Chain-of-Thought reasoning, tool invocation, and memory management to enable sophisticated behaviors for orchestrating intricate workflows (Xi et al., 2025; Ke et al., 2025; Gridach et al., 2025; Plaat et al., 2025). However, the handcrafted nature of these systems imposes limitations on scalability, adaptability, and rapid deployment across diverse domains.

To address these limitations, recent trends have shifted towards automated design methods for agentic systems (Hu et al., 2025a; Shang et al., 2025; Zhang et al., 2025a; Zhuge et al., 2024; Liu et al., 2024b; Hu et al., 2025b; Yuan et al., 2025). Automated methods typically employ search algorithms to discover optimal workflow configurations by systematically exploring a vast design space. Instead of relying on human intuition, these approaches generally involve iterations of candidate generation, evaluation, and refinement. While promising, these methods exhibit significant drawbacks, chiefly the high computational costs associated with the extensive validation steps needed during the exploration and evaluation phases of the search. Each candidate configuration must undergo rigorous evaluation, often through expensive, repeated interactions with LLM APIs, rendering the search prohibitively costly and time-consuming.

In this paper, we argue that purely search-based automated design methods are inherently inefficient and propose a predictive approach to significantly accelerate workflow evaluation. Specifically, we advocate for a predictor-based framework that can rapidly estimate the performance of candidate agentic workflows, similar to performance predictors in neural architecture search (White et al., 2021), thereby reducing the need for extensive validation.

As depicted in Figure 1, instead of fully evaluating every candidate, a predictive model can estimate the quality and viability of agentic workflows, thus guiding the search process far more efficiently. By reducing costly ground-truth executions or environment interactions during the search process, prediction-based approaches promise significant improvements in both search efficiency and solution quality. However, building a high-quality predictor for agentic workflows poses **two** fundamental challenges.
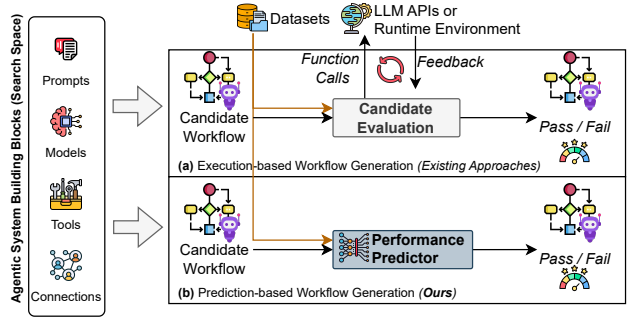


Figure 1: Comparison between (a) execution-based and (b) prediction-based candidate evaluation for agentic workflow generation. Execution-based methods rely on costly runtime or LLM calls, while our prediction-based approach offers faster, scalable evaluation via a learned predictor.

**Workflow Heterogeneity**. Agentic workflows exhibit considerable heterogeneity; subtle variations in configuration can lead to dramatically different performances. Specifically, workflows can vary widely in communication structure, prompting strategies, tool invocation patterns, and reasoning styles, making it challenging to learn a unified predictive model. Moreover, agentic systems differ significantly across tasks, domains, and toolsets, resulting in diverse and complex workflow configurations that are difficult to model uniformly (Xu et al., 2024; Qiao et al., 2025).

**Scarcity of Labeled Data**. The availability of labeled data for training effective prediction models is severely limited due to the prohibitive cost of generating performance labels through exhaustive validation. Constructing a large, diverse set of labeled workflows with known execution outcomes is particularly expensive, creating a data bottleneck for supervised learning approaches. Moreover, gathering large-scale, high-quality labels for agentic workflows (e.g., success rates and execution outcomes) is often infeasible, further limiting the amount of supervised training data available for learning accurate predictors.

To tackle these challenges, we introduce **Agentic Predictor**, a novel multi-view encoding framework for performance prediction in agentic workflows. Specifically, to address the workflow heterogeneity problem, we propose a *multi-view workflow encoding* technique that leverages multi-view representation learning to capture the heterogeneous and multifaceted characteristics of agentic workflow configurations. This technique enables Agentic Predictor to capture complementary and diverse aspects of agentic workflows for a given task, encompassing rich structural, behavioral, and semantic perspectives. Additionally, we employ *cross-domain unsupervised pretraining*, exploiting abundant unlabeled data from various related agentic tasks and domains to significantly enhance the predictive model's generalization capabilities and alleviate label scarcity. In particular, we pretrain the multi-view workflow encoders on unlabeled agentic workflow instances across diverse domains using contrastive and unsupervised objectives, thereby equipping the predictor with robust representations before fine-tuning on limited labeled data.

In summary, the **main contributions** of our work are as follows.

- We devise a novel multi-view workflow encoding technique that captures the heterogeneous facets of agentic workflows, thereby enhancing predictive performance and generalization.
- We propose a cross-domain unsupervised pretraining strategy that enables effective predictor training even with limited labeled workflows.
- By bridging multi-view workflow encoding with cross-domain unsupervised pretraining, we introduce Agentic Predictor, the first predictive framework specifically designed to address the key challenges of heterogeneous workflow configurations and limited labeled data availability.
- We empirically demonstrate that Agentic Predictor yields substantial improvements of up to **12.12%** in prediction accuracy and **15.16%** in workflow utility over strong baselines across three domains.

## 2 Related Work

### 2.1 Automated Generation of Agentic Workflows

Recent advancements (Xi et al., 2025; Ke et al., 2025; Gridach et al., 2025; Plaat et al., 2025) in agentic workflows have led to the development of various frameworks aimed at enhancing multi-agent collaboration for complex tasks (Trirat et al., 2025; Guo et al., 2024a; Niu et al., 2025; Guo et al., 2024b). MetaGPT (Hong et al., 2024) and ChatDev (Qian et al., 2024) use predefined multi-agent structures to address coding challenges, while AgentVerse (Chen et al., 2024) introduces iterative collaboration where agents discuss, execute, and evaluate tasks. LLM-Debate (Du et al., 2024) employs multiple expert agents that engage in debates over several rounds to derive final answers. However, these systems often rely on static configurations, which limits their adaptability to diverse queries across different tasks and domains.

To optimize agentic workflows, GPTSwarm (Zhuge et al., 2024) and G-Designer (Zhang et al., 2024) apply variants of the REINFORCE algorithm to optimize workflow structures represented as directed acyclic graphs (DAGs), while DyLAN (Liu et al., 2024b) dynamically selects agents based on task requirements. ADAS (Hu et al., 2025a) and AFlow (Zhang et al., 2025a) further leverage powerful LLMs (e.g., Claude-3.5-Sonnet and GPT-4) to iteratively generate task-specific multi-agent systems. Similarly, AgentSquare (Shang et al., 2025) proposes a modular design space for automatic LLM agent search, enhancing adaptability to novel tasks. Despite their effectiveness, these methods typically require numerous LLM calls, resulting in significant computational and financial overheads, making them less practical for real-world applications.

Instead of manually designing a fixed workflow (Qian et al., 2024; Chen et al., 2024; Du et al., 2024) or incurring multiple LLM inference costs to generate one for each query (Zhuge et al., 2024; Liu et al., 2024b), Agentic Predictor uses a performance predictor to efficiently estimate the

Table 1: Comparison between ours and existing frameworks for prediction-based workflow generation.

| Framework | Multi-View Representation | Unsupervised Pretraining | Lightweight Predictor | Search Agnostic |
|---|---|---|---|---|
| MAS-GPT (Ye et al., 2025) | × | × | × | × |
| FLORA-Bench (Zhang et al., 2025b) | × | × | ✓ | ✓ |
| **Agentic Predictor** *(Ours)* | ✓ | ✓ | ✓ | ✓ |

performance of agentic workflows. This approach enables more effective exploration of optimal workflow configurations without the need for exhaustive evaluations or repeated LLM calls. Two *contemporaneous efforts* pursue related goals. FLORA-Bench (Zhang et al., 2025b) is a position paper that advocates for the use of graph neural networks as efficient predictors of agentic workflow performance, proposing a benchmark rather than a concrete method. MAS-GPT (Ye et al., 2025) introduces a framework based on fine-tuning LLMs to directly generate agentic workflows via a single LLM call. In contrast, our Agentic Predictor is a novel framework that enhances the representational capacity of workflows through multi-view representations (compared to FLORA-Bench's single-view approach) and employs a lightweight predictor (compared to MAS-GPT's supervised LLM fine-tuning), achieving a favorable trade-off between predictive accuracy and sample efficiency. Table 1 summarizes the key differences between Agentic Predictor and these frameworks.

### 2.2 Performance Predictors for NAS

Neural architecture search (NAS) has spurred the development of performance predictors that aim to reduce the significant computational cost of evaluating candidate architectures. PRE-NAS (Peng et al., 2022) employs a predictor-assisted evolutionary strategy to estimate model performance, thereby alleviating the need for exhaustive training. BRP-NAS (Dudziak et al., 2020) integrates graph convolutional networks to forecast hardware-aware performance metrics, improving the practicality of NAS under resource constraints. CAP (Ji et al., 2024) introduces a context-aware neural predictor, leveraging self-supervised learning to generate expressive and generalizable representations of architectures, thus enabling more effective search space exploration. FlowerFormer (Hwang et al., 2024) advances architecture encoding through a flow-aware graph transformer, yielding improved prediction accuracy. A unifying trend among these methods is the emphasis on learning more informative neural representations to guide the search process. Building on this insight, we propose the Agentic Predictor framework, which approaches performance prediction from a representation-centric perspective. By incorporating multi-view representations conditioned on workflow configurations, Agentic Predictor facilitates accurate performance estimation and efficient exploration of the agentic workflow space.
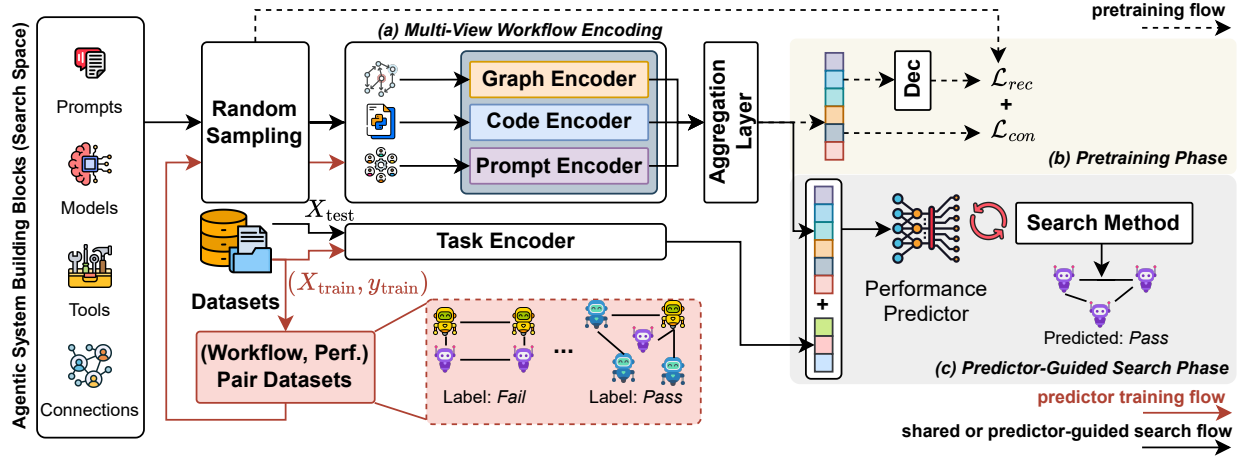
Figure 2: Overview of our Agentic Predictor framework. A **(a) multi-view workflow encoder** is designed to encode a set of agentic workflows from graph, code, and prompt aspects into unified representations, which serve as features for training the predictor. In the **(b) pretraining phase**, the encoder learns these representations on unlabeled workflows spanning diverse tasks and domains, using cross-domain unsupervised pretraining objectives. In the **(c) predictor-guided search phase**, a performance predictor is trained on a small (workflow configuration, performance) dataset to classify configurations as pass or fail, and subsequently guides the search toward promising configurations.

# 3    Methodology: *Agentic Predictor*

## 3.1    Problem Formulation

Let an agentic workflow be denoted as $\mathcal{W} = \{\mathcal{V}, \mathcal{E}, \mathcal{P}, \mathcal{C}\}$, where $\mathcal{V} = \{v_i\}_{i=1}^{N}$ represents the set of $N$ agents, $\mathcal{E}$ denotes the set of edges defining the connections between agents, and $\mathcal{P} = \{p_i\}_{i=1}^{N}$ denotes the system prompts for each agent $i$. $\mathcal{C}$ represents the complete code specifying the logic and structure of the workflow.

Given a task description $T$, the workflow $\mathcal{W}$ autonomously executes agents recursively in topological order, where the $i$-th agent receives the task description $T$ along with the outputs $y$ from its predecessor agents. Formally, the input to agent $i$ is defined as $\mathcal{X}_i = \{T\} \cup \{y_j : v_j \in \mathcal{N}_i^{(\text{in})}\}$, where $\mathcal{N}i^{(\text{in})}$ denotes the set of predecessor agents of agent $i$, and $y_j$ is the output of agent $j$. The output $y_i$ of agent $i$ is generated by querying an LLM: $y_i = \text{LLM}(\mathcal{X}_i, p_i)$. After executing all agents, the final response of the agentic workflow is defined as $r = f_{\text{LLM}}(\mathcal{W}, T)$, where $f_{\text{LLM}}$ represents the overall execution process of the given LLM. In most existing work, this process is repeated for the evaluation of $r$, which incurs significant computational and financial overhead.

In contrast, this paper aims to design a predictive model $\mathcal{M}$ that efficiently estimates the final performance of an agentic workflow $\mathcal{W}$ on a given task description $T$, without requiring costly LLM invocations. Therefore, we treat the workflow $\mathcal{W}$ and task description $T$ as inputs to the predictor $\mathcal{M}$, which outputs the estimated performance $\hat{e}$. Formally, $\hat{e} = \mathcal{M}_\Theta(\mathcal{W}, T)$, where $\Theta$ is the learnable parameters of the performance predictor.

**Learning Objective**. Given the workflow $\mathcal{W}$, task description (or query) $T$, and performance predictor $\mathcal{M}_\Theta$ parameterized by $\Theta$, we aim to find the optimal $\Theta$ that minimizes the error between the estimated performance $\hat{e}$ and the ground-truth performance $e$. Formally, we solve

$$\min_\Theta \; \mathbb{E}_{(\mathcal{W}, T)}[\mathcal{L}(e, \hat{e})], \tag{1}$$

where $\mathcal{L}(\cdot, \cdot)$ is a loss function that quantifies the discrepancy between the ground truth and the predicted performance. $\mathcal{L}$ can be either the cross-entropy loss or the mean squared error loss, depending on whether the prediction task is formulated as a classification or regression problem.

## 3.2 Framework Overview

We present an overview of our Agentic Predictor framework in Figure 2 and Algorithm 1. First, the multi-view workflow encoder integrates multiple aspects of an agentic workflow—namely, graph structures $(\mathcal{V}, \mathcal{E})$, code implementations $\mathcal{C}$, and system prompts $\mathcal{P}$—into unified representations $\mathcal{F}$. This integration is accomplished via specialized encoders for each modality, followed by an aggregation layer to consolidate the multi-modal features. Second, in the pretraining phase, these unified workflow representations are further refined through unsupervised learning objectives: reconstruction and contrastive tasks. This cross-domain pretraining enhances the encoder's generalization and adaptability across diverse tasks and configurations. Third, a dedicated performance predictor $\mathcal{M}_\Theta$ is trained using a small, labeled dataset comprising pairs of workflow configurations $\mathcal{W}$ and their corresponding performance outcomes $e$ on task descriptions $T$. Finally, with the trained performance predictor, we can perform a predictor-guided search to efficiently rank and select promising workflow configurations without requiring expensive LLM invocations. Since Agentic Predictor is a search-agnostic framework, we do not incorporate a specific search process within the framework itself.

---

**Algorithm 1** Overall Procedure of Agentic Predictor

---

**Initialization:** Multi-View Encoder $\text{Enc}(\cdot)$ and Performance Predictor Model $\mathcal{M}_\Theta$
**Input:** User Instruction (or Task Description) $T \in \mathcal{T}$ and Training Data $D^{\text{train}}$
1: ▷ Phase 1: Cross-Domain Unsupervised Pretraining (optional)
2: Sample $M$ unlabeled workflows $\mathcal{W}_1, \mathcal{W}_2, ..., \mathcal{W}_M$ from multiple domains
3: **for** each $\mathcal{W}_i = (\mathcal{G}_i, \mathcal{C}_i, \mathcal{P}_i)$ **do**
4: $\quad \mathbf{Z}_i \leftarrow \text{Enc}(\mathcal{G}_i, \mathcal{C}_i, \mathcal{P}_i)$ $\hspace{3cm}$ ▷ Encode multiview graph, code, and prompts
5: $\quad (\hat{\mathcal{G}}_i, \hat{\mathcal{C}}_i, \hat{\mathcal{P}}_i) \leftarrow \text{Dec}(\mathbf{Z}_i)$ $\hspace{4cm}$ ▷ Decode reconstructions
6: **end for**
7: $\mathcal{L}_{enc} = \mathcal{L}_{rec} + \mathcal{L}_{con}$ $\hspace{5cm}$ ▷ Minimize total pretraining loss
8: ▷ Phase 2: Training Performance Predictor
9: Obtain (small) labeled dataset $\{(\mathcal{W}_j, T_j, e_j)\}_{j=1}^N$ from $D^{\text{train}}$
10: **for** each $(\mathcal{W}_j, T_j)$ **do**
11: $\quad \mathbf{Z}_j \leftarrow \text{Enc}(\mathcal{W}_j)$ $\hspace{6cm}$ ▷ Encode workflow
12: $\quad \mathbf{T}_j \leftarrow \text{TaskEncoder}(T_j)$ $\hspace{4.5cm}$ ▷ Encode task description
13: $\quad \mathcal{F}_j \leftarrow \text{MLP}([\mathbf{Z}_j, \mathbf{T}_j])$ $\hspace{4.5cm}$ ▷ Form joint representation
14: $\quad \hat{e}_j \leftarrow \mathcal{M}_\Theta(\mathcal{F}_j)$ $\hspace{5.5cm}$ ▷ Predict performance
15: **end for**
16: Train $\mathcal{M}_\Theta$ using binary cross-entropy loss $\mathcal{L}_{pred}(e_j, \hat{e}_j)$, where $\{e_j\}_{j=1}^N$
17: ▷ Phase 3: Predictor-Guided Candidate Ranking
18: Sample $K$ candidate workflows $\{\mathcal{W}_k\}_{k=1}^K$
19: **for** each $\mathcal{W}_k$ **do**
20: $\quad \mathbf{Z}_k \leftarrow \text{Enc}(\mathcal{W}_k)$ $\hspace{6cm}$ ▷ Encode workflow
21: $\quad \mathcal{F}_k \leftarrow \text{MLP}([\mathbf{Z}_k, \mathbf{T}])$ $\hspace{5cm}$ ▷ Encode task
22: $\quad \hat{e}_k \leftarrow \mathcal{M}_\Theta(\mathcal{F}_k)$ $\hspace{5.8cm}$ ▷ Predict score
23: **end for**
24: Rank all $\{\mathcal{W}_k\}$ by predicted scores $\hat{e}_k$
25: **return** top-$k$ ranked workflows for final evaluation

---

## 3.3 Multi-View Workflow Encoding

Motivated by recent findings in the NAS literature (White et al., 2020; Akhauri & Abdelfattah, 2024; Trirat & Lee, 2024), which demonstrate that architecture representations significantly influence predictor performance, we advocate the necessity of developing expressive and comprehensive representations tailored explicitly for agentic workflows. Due to fundamental differences between agentic workflows and traditional neural network architectures, conventional graph-based representations alone are inadequate. Although DAGs naturally capture explicit inter-agent communication and dependencies, they fail to encode crucial implicit details, such as tool usage patterns, code structures, computational complexity, and the nuanced semantics inherent in agent prompts. To overcome these limitations, we introduce a novel multi-view encoding scheme that

integrates complementary representations at multiple granularities, with each view capturing distinct yet essential aspects of agentic workflows.

- **Graph View** explicitly captures structural dependencies and direct interactions among agents, emphasizing inter-agent communication channels. We denote graph view as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$.
- **Code View** implicitly encodes complex semantic structures, logical sequences, computational complexities, and patterns of tool usage inherent in workflow implementations $\mathcal{C}$.
- **Prompt View** provides semantic embeddings that encapsulate nuanced agent roles, behavioral descriptions, and broader contextual guidance embedded within system prompts $\mathcal{P}$.

The theoretical justification for adopting this multi-view framework is grounded in the hypothesis that incorporating heterogeneous information sources reduces representational bias, thereby improving both the robustness and predictive accuracy of the model.

### 3.3.1 Encoder Networks

We now detail the components of our proposed multi-view workflow encoding method for agentic workflow generation using neural networks. Let $\text{Enc}(\cdot)$ denote an encoder function that maps a candidate workflow—composed of $(\mathcal{G}, \mathcal{C}, \mathcal{P})$—into $d$-dimensional Euclidean space, i.e., $\text{Enc}(\cdot) : (\mathcal{G}, \mathcal{C}, \mathcal{P}) \to \mathbb{R}^d$. Given the heterogeneous nature of workflow configurations, we design three specialized encoder networks, each responsible for learning a representation corresponding to a distinct view. These view-specific representations are then aggregated into a shared latent space, denoted by $\mathbf{Z} = \text{Enc}(\mathcal{G}, \mathcal{C}, \mathcal{P})$, where $\mathbf{Z} \in \mathbb{R}^d$. This continuous latent representation is used to train the performance predictor $\mathcal{M}_\Theta$ (see §3.5). The individual encoders for each view are integrated into a unified architecture as described below.

**Encoder for Graph View**. Following Zhang et al. (2025b), we employ standard graph neural network (GNN) layers to encode graph-based representations. However, instead of relying on a single graph, we adopt a *multi-graph* learning approach that integrates node features from multiple views, including agent-specific definitions and function call implementations at each agent node. This design enables the GNN layers to effectively capture explicit structural dependencies and communication flows among agents, which are essential for modeling interaction patterns. Formally, we obtain the graph representation as

$$\mathbf{Z}_\mathcal{G} = \text{AttenPool}(\text{CrossGraphAttn}([\text{GNN}(\mathcal{G}_{\text{prompt}}), \text{GNN}(\mathcal{G}_{\text{code}}), \text{GNN}(\mathcal{G}_{\text{operator}})]), \tag{2}$$

where CrossGraphAttn is a multi-head attention layer used to capture inter-graph importance, while AttenPool is a view-attention pooling layer that aggregates importance scores across different views.

**Encoder for Code View**. To model code semantics, we adopt a workflow-level embeddings. Here, we use $L$-layer multi-layer perceptron (MLP) to extract latent semantic features. This approach allows the model to learn intricate computational logic and tool interactions on the global level. Thus, we compute the code representation as $\mathbf{Z}_\mathcal{C} = \text{MLP}_\mathcal{C}(\mathcal{C})$.

**Encoder for Prompt View**. Unlike FLOW-GNN (Zhang et al., 2025b), which encodes agent prompts as node-level features, we use another $L$-layer MLP to encode the *entire* workflow instruction prompts holistically. This approach allows us to capture role descriptions, behavioral intents, and global context—resulting in richer and more semantically informed representations. The prompt encoding is thus given by $\mathbf{Z}_\mathcal{P} = \text{MLP}_\mathcal{P}(\mathcal{P})$.

**Aggregation Layer**. The representations from the graph, code, and prompt encoders—$\mathbf{Z}_\mathcal{G}$, $\mathbf{Z}_\mathcal{C}$, and $\mathbf{Z}_\mathcal{P}$—are concatenated and passed through a final MLP layer. This aggregation mechanism adaptively integrates information across all views, enabling the model to emphasize the most contextually relevant aspects. The final output of the encoder $\text{Enc}(\cdot)$ is computed as $\mathbf{Z} = \text{MLP}([\mathbf{Z}_\mathcal{G}, \mathbf{Z}_\mathcal{C}, \mathbf{Z}_\mathcal{P}])$.

These encoders learn not only from different workflow perspectives but also at varying levels of granularity—specifically, at the graph level for agent interactions, the code level for logical structures, and the prompt level for agent-specific instructions.

### 3.3.2 Decoder Networks

The decoder is a generative model aiming at reconstructing $\hat{\mathcal{G}}$, $\hat{\mathcal{C}}$, and $\hat{\mathcal{P}}$ from the latent variables $\mathbf{Z}$ to learn the general representations of the agentic workflows. Specifically, its constituent is a stack of MLP layers. For ease of implementation, the output of the decoder networks are the input embedding values of $\mathcal{G}$, $\mathcal{C}$, and $\mathcal{P}$. Hence, the decoder $\text{Dec}(\cdot)$ is formulated as $\hat{\mathcal{G}} = \text{MLP}(\mathbf{Z}_{\mathcal{G}})$, $\hat{\mathcal{C}} = \text{MLP}(\mathbf{Z}_{\mathcal{C}})$, and $\hat{\mathcal{P}} = \text{MLP}(\mathbf{Z}_{\mathcal{P}})$.

## 3.4 Cross-Domain Unsupervised Pretraining

In real-world scenarios, the availability of labeled performance data for agentic workflows is highly limited due to the costly and time-consuming evaluation process. To address this challenge and enable data-efficient predictor training, we *optionally* adopt a two-phase strategy. Rather than directly supervising the encoder with performance labels, we first perform cross-domain unsupervised pretraining to obtain rich and generalizable workflow representations $\mathbf{Z}$. These *learned* representations significantly enhance sample efficiency for downstream performance prediction, as shown in NAS literature (Akhauri & Abdelfattah, 2024; Trirat & Lee, 2024; Yan et al., 2020; 2021). Note that, given sufficient labeled data, direct supervised learning is also feasible.

**Multi-Task Pretraining**. We train the multi-view encoder on $M$ unlabeled workflow configurations by minimizing a combined loss comprising reconstruction and contrastive objectives: $\mathcal{L}_{rec} = \frac{1}{M} \sum_{i=1}^{M} \|\mathcal{G}_i - \hat{\mathcal{G}}_i\|^2 + \|\mathcal{C}_i - \hat{\mathcal{C}}_i\|^2 + \|\mathcal{P}_i - \hat{\mathcal{P}}_i\|^2$ and $\mathcal{L}_{con} = \frac{1}{M} \sum_{i=1}^{M} -\log \frac{\exp(\text{sim}(\mathbf{Z}_i, \mathbf{Z}_j^+)/\tau)}{\sum_{k=1}^{M} \exp(\text{sim}(\mathbf{Z}_i, \mathbf{Z}_k)/\tau)}$. Here, $\mathcal{G}_i, \mathcal{C}_i, \mathcal{P}_i$ denote the input graph, code, and prompt embeddings, respectively, while $\hat{\cdot}$ denotes reconstructions via a shared decoder. The contrastive loss encourages semantically similar configurations—e.g., successful workflows for the same task—to lie closer in the embedding space, while dissimilar ones are pushed apart. Positive pairs $(\mathbf{Z}_i, \mathbf{Z}_j^+)$ are drawn from configurations that solve the same task successfully, while negatives include task-mismatched or failed configurations. This learning objective encourages the encoder to capture meaningful signals related to both structure and performance. Thus, the total loss function is $\mathcal{L}_{enc} = \mathcal{L}_{rec} + \mathcal{L}_{con}$.

## 3.5 Performance Predictor

Following the unsupervised pretraining of the multi-view encoder, we introduce a lightweight performance predictor to guide the exploration of the large agentic workflow space. This phase enables efficient identification of high-performing configurations with minimal supervision, using only a small set of labeled workflow-performance pairs. As shown in Figure 2(c), our predictor operates on the learned workflow embeddings, enriched with task-specific context, to form a joint representation $\mathcal{F}$ used for binary success/failure prediction and downstream search.

**Task Encoder**. To capture task-specific characteristics, we incorporate a Task Encoder that generates high-level semantic embeddings from natural language task descriptions. These embeddings, derived from pretrained language models (e.g., T5 or BERT), provide critical global context that helps the model differentiate between tasks with similar surface structures but different functional requirements. The task embedding is concatenated with the multi-view workflow representation, forming a joint representation $\mathcal{F} = [\mathbf{Z}, \mathbf{T}]$ where $\mathbf{Z}$ is the encoded workflow and $\mathbf{T}$ is the task embedding. This fused representation captures both compositional and contextual signals, which enhances the predictor's ability to generalize across heterogeneous tasks with varying operational goals and constraints.

The performance predictor is a lightweight classifier $\mathcal{M}_\Theta$ (e.g., a MLP) trained on a limited set of labeled data $(X_{\text{train}}, y_{\text{train}})$, where each $X_{\text{train}}$ corresponds to the joint representation $\mathcal{F}$ and $y_{\text{train}}$ is a binary indicator of success (pass/fail) or a performance score. The predictor is optimized using a standard binary cross-entropy loss. Hence, $\mathcal{L}$ in Equation 1 is reformulated as $\mathcal{L}_{\text{pred}} = -\frac{1}{N} \sum_{i=1}^{N} [e_i \log \hat{e}_i + (1 - e_i) \log(1 - \hat{e}_i)]$, where $\hat{e}_i$ denotes the predicted probability of success for the $i$-th configuration. By operating on the pretrained, semantically rich embeddings, the predictor achieves high accuracy even in the low-data regime, enabling label-efficient search. In ablation studies (see §4.3), we demonstrate the importance of using both workflow and task signals in forming $\mathcal{F}$ and the effectiveness of the pretraining phase.

**Integration with Workflow Optimization**. With the trained predictor in place, we can perform predictor-guided search to efficiently explore the workflow configuration space. Rather than evaluating each configuration via full execution, we embed candidates into their joint representations $\mathcal{F}$ and score them using the predictor. The top-scoring candidates are selected for evaluation. This substantially reduces computational cost by focusing on the most promising regions of the search space. A simple yet effective instantiation of this strategy uses random search to sample $K$ workflow candidates from the full configuration space, and then ranks them using the learned predictor. We select the top-$k$ configurations (based on predicted $\hat{e}$ values) for evaluation. This predictor-as-ranker setup transforms random search into a label-efficient guided procedure without requiring complex heuristics.

Since our main contribution is the performance predictor rather than the optimization process, we focus our validation of Agentic Predictor on prediction accuracy and ranking quality (i.e., workflow utility) in the following section. Additional experimental results on workflow optimization are provided in §A.

## 4 Experiments

We conduct a comprehensive evaluation of the proposed Agentic Predictor framework from multiple perspectives, guided by the following questions: **(Q1)** How does Agentic Predictor perform as a predictor of agentic workflow performance compared to relevant baselines? **(Q2)** How do different design choices and configurations of Agentic Predictor affect its predictive accuracy? **(Q3)** Is the pretraining phase helpful for maintaining prediction quality under varying numbers of labels?

### 4.1 Setup

**Benchmarks**. To evaluate the effectiveness of performance predictors in agentic workflows, we adopt a recently introduced, well-curated benchmark: FLORA-Bench (Zhang et al., 2025b). We use only the AFlow (Zhang et al., 2025a) subset, as it is the only one publicly available. This benchmark spans

Table 2: Summary of benchmark statistics.

| Domains | Code Generation | Math Problem | Reasoning Task |
|---|---|---|---|
| # workflows | 38 | 42 | 30 |
| Average # nodes | 6.11 | 5.49 | 6.58 |
| # tasks | 233 | 782 | 2400 |
| # samples | 7362 | 4059 | 72000 |

five representative datasets across three core domains in the agentic workflow literature: code generation (HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021)), mathematics problem solving (GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021b)), and reasoning (MMLU (Hendrycks et al., 2021a)). Table 2 provides an overview of FLORA-Bench. We randomly split each domain into training (80%), validation (10%), and test (10%) sets.

**Evaluation Metrics**. To ensure a fair and consistent comparison, we strictly adhere to the official evaluation protocols specified by the benchmark.

- **Accuracy**: Accuracy quantifies how well a model predicts agentic workflow performance. It is defined as $accuracy = \frac{1}{|\mathcal{D}^{\text{test}}|} \sum_i^{|\mathcal{D}^{\text{test}}|} \mathbf{1}(\hat{e}_i = e_i)$, where $|\mathcal{D}^{\text{test}}|$ is the size of the test split, and $\hat{e}_i$ and $e_i$ denote the predicted and ground-truth performance, respectively. $\mathbf{1}(\cdot)$ is the indicator function, which returns 1 if $\hat{e}_i = e_i$, and 0 otherwise.
- **Utility**: Utility evaluates the consistency between the workflow rankings predicted by the model and the ground-truth rankings, emphasizing the model's ability to determine the relative order of different workflows. The utility score is computed as follows: first, calculate the ground-truth and predicted success rates of a workflow $\mathcal{W}_i$ by averaging $e$ and $\hat{e}$ across all tasks in $\mathcal{D}^{\text{test}}$. Then, rank the workflows and extract the top-$k$ workflows according to the respective scores, resulting in two ordered sets: $\mathcal{H} = \{\mathcal{W}_i\}_{i=1}^k$ and $\hat{\mathcal{H}} = \{\mathcal{W}_i'\}_{i=1}^k$. Formally, $utility = \frac{1}{k} \sum_{i=1}^k \mathbf{1}(\mathcal{W}_i' \in \mathcal{H})$.

**Baselines**. Since there is no direct baseline method specifically designed for performance prediction in agentic systems, we adopt comparison baselines from the benchmark paper. Some of these methods have previously been used as performance predictors for NAS (White et al., 2021). The selected baselines include one naive **MLP** and several strong graph-based models: **GCN** (Kipf & Welling, 2017), **GAT** (Veličković et al., 2018), **GCN-II** (Chen et al., 2020), **Graph Transformer** (Shi et al., 2021), and **One For All** (Liu et al., 2024a).

Table 3: Performance comparison between Agentic Predictor and baseline methods. The best and second-best results are highlighted in **bold** and underlined, respectively.

| Domain | Code Generation | | Math Probelm | | Reasoning Task | | Average | |
|---|---|---|---|---|---|---|---|---|
| Model | Accuracy | Utility | Accuracy | Utility | Accuracy | Utility | Accuracy | Utility |
| MLP | 78.02±0.59 | 73.94±1.35 | 73.73±0.31 | 69.64±0.29 | 78.45±0.08 | 88.48±0.63 | 76.73±0.33 | 77.35±0.76 |
| GCN | 84.35±0.34 | 72.73±3.18 | 76.19±0.42 | 66.52±1.66 | 87.12±0.14 | **91.82±0.46** | 82.55±0.30 | 77.02±1.77 |
| GAT | 84.49±0.56 | 76.46±0.91 | 76.44±0.61 | 66.51±1.28 | 87.07±0.08 | 89.40±0.68 | 82.67±0.42 | 77.46±0.96 |
| GCN-II | 83.72±0.40 | 77.75±1.98 | 75.04±0.31 | 64.33±0.47 | 87.28±0.14 | 89.92±1.90 | 82.01±0.28 | 77.33±1.45 |
| Graph Transformer | 84.71±0.45 | 74.09±0.35 | 75.45±0.23 | 66.48±0.96 | 86.93±0.27 | 90.60±1.97 | 82.36±0.32 | 77.06±1.09 |
| One For All | 81.05±0.34 | 73.42±1.39 | 75.21±0.23 | 69.08±0.64 | 82.52±0.13 | 87.64±1.98 | 79.59±0.23 | 76.71±1.34 |
| *Agentic Predictor* | **85.62±0.47** | **80.08±0.46** | **79.56±0.25** | **74.08±0.47** | **87.96±0.02** | 91.47±0.44 | **84.38±0.25** | **81.88±0.46** |
| % Improvement (up to) | 9.74% | 10.11% | 7.91% | 15.16% | 12.12% | 4.37% | 9.97% | 6.74% |

**Implementation Details**. For all methods, we follow the same setup as suggested by Zhang et al. (2025b). Specifically, we use a 2-layer backbone with a hidden dimension of 512. Dropout rate is set to 0.5 and the batch size is 512. The models are optimized using the Adam optimizer (Kingma & Ba, 2014) with a learning rate of $1 \times 10^{-4}$ and a weight decay of $5 \times 10^{-4}$. Training is conducted for 200 epochs on a single NVIDIA A100-SXM4-80GB GPU. The best model checkpoint is selected based on the highest accuracy on the validation subset. In addition for Agentic Predictor, we use `all-MiniLM-L6-v2` (Wang et al., 2020) and `CodeRankEmbed` (Suresh et al., 2025) for text and code embeddings, respectively.

## 4.2 Main Results

We report all experimental results for agentic workflow performance prediction averaged from three runs on different random seeds[1] in the same dataset.

**Prediction Accuracy (Q1)**. We report averaged performance scores and standard deviations in Table 3. Our proposed framework, Agentic Predictor, consistently outperforms all baseline methods across the three task domains. In terms of accuracy, Agentic Predictor achieves top results in each domain—85.62%, 79.56%, and 87.96%, respectively—resulting in the highest overall average accuracy of 84.38%. This reflects an improvement of up to 9.97% over the comparison baselines. Utility scores exhibit a similar trend. Agentic Predictor attains the highest utility in code generation (80.08%) and math problem solving (74.08%), and a near-best score in reasoning tasks (91.47%), second only to GCN (91.82%). On average, it achieves the highest utility score of 81.88%, marking an improvement of up to 6.74% over the baselines.

These results validate that Agentic Predictor not only enhances predictive accuracy but also improves downstream utility in diverse agentic workflows, demonstrating its robustness and generalizability across task types. The consistent performance gains underscore the benefits of leveraging multi-view encoding across heterogeneous agentic workflows.

## 4.3 Additional Analyses

**Ablation Study (Q2)**. To substantiate our contributions on specific design of multi-view workflow encoding in Agentic Predictor, we conduct ablation study on two main components: multi-view encoder and multi-graph encoding techniques. According to the results in Table 4 and Table 5, we find that incorporating all three input views—code, graph, and text—results in the best overall performance across all tasks. Specifically, the full model configuration achieves the highest average accuracy (84.38%) and utility (81.88%), underscoring the complementary value of each modality. Notably, the removal of any single view leads to a consistent drop in performance, demonstrating the synergistic role of multimodal inputs in prediction capabilities of Agentic Predictor.

Furthermore, results in Table 5 reveal the significance of multi-graph encoding. When multiple graphs are used instead of a single graph, the model shows a clear performance improvement, particularly in code generation (accuracy improves from 82.58% to 84.44%) and reasoning tasks (utility rises from 90.14% to 91.03%). This supports our hypothesis that different graph perspectives enrich structural context and lead to

---

[1]The random seed for each run is $2^r$, where $r$ is the running index starting from 0.

Table 4: Results of ablation study on different input view variations.

| View Variations | | | Codi Generation | | Math Problem | | Reasoning Task | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|
| Code | Graph | Text | Accuracy | Utility | Accuracy | Utility | Accuracy | Utility | Accuracy | Utility |
| ✓ | | | 82.04±0.51 | 75.66±0.66 | 75.70±0.14 | 68.52±0.91 | 83.19±0.56 | **91.51±0.61** | 80.31±0.40 | 78.56±0.73 |
| | ✓ | | 84.44±0.31 | 77.22±3.46 | 79.14±0.28 | 67.99±3.36 | 87.00±0.21 | 91.03±1.23 | 83.53±0.27 | 78.75±2.68 |
| | | ✓ | 79.87±0.28 | 70.34±0.43 | 76.60±0.65 | 68.45±1.80 | 68.06±0.00 | 71.04±0.00 | 74.84±0.31 | 69.94±0.74 |
| ✓ | ✓ | | 83.72±0.83 | 73.97±0.81 | 75.86±0.85 | 70.18±1.64 | 86.88±0.14 | 86.14±4.62 | 82.15±0.61 | 76.76±2.36 |
| ✓ | | ✓ | 82.27±0.63 | 77.28±1.12 | 76.03±0.14 | 66.66±4.18 | 54.17±0.00 | 53.21±0.00 | 70.82±0.26 | 65.72±1.77 |
| | ✓ | ✓ | 82.45±1.36 | 74.64±1.57 | 75.70±1.26 | 67.83±3.71 | 69.47±0.00 | 70.55±0.00 | 75.87±0.87 | 71.01±1.76 |
| ✓ | ✓ | ✓ | **85.62±0.47** | **80.08±0.46** | **79.56±0.25** | **74.08±0.47** | **87.96±0.02** | 91.47±0.44 | **84.38±0.25** | **81.88±0.46** |

Table 5: Results of ablation study on different input graph variations.

| Graph Variations | | Code Generation | | Math Problem | | Reasoning Task | | Average | |
|---|---|---|---|---|---|---|---|---|---|
| Single Graph | Multi Graph | Accuracy | Utility | Accuracy | Utility | Accuracy | Utility | Accuracy | Utility |
| ✓ | | 82.58±0.54 | 78.52±2.91 | 78.57±0.73 | 67.51±2.11 | 86.95±0.13 | 90.14±3.10 | 82.70±0.47 | 78.72±2.71 |
| | ✓ | 84.44±0.31 | 77.22±3.46 | 79.14±0.28 | 67.99±3.36 | 87.00±0.21 | 91.03±1.23 | 83.53±0.27 | 78.75±2.68 |

more robust representations. Together, these findings validate the architectural choices in Agentic Predictor, demonstrating that both multi-view and multi-graph designs are integral to its superior performance.

**Effects of Pretraining Phase (Q3)**. Since acquiring a large amount of ground-truth labels from agentic workflows is expensive, we examine whether cross-domain unsupervised pretraining (denoted as Agentic Predictor+) benefits settings where labeled instances are limited. We vary the label ratio from 0.1 to 0.5, selecting labeled samples from the training split of all datasets in the benchmark. We pretrain the proposed multi-view encoder with a batch size of 32 for 20 epochs. On average, the results shown in Figure 3 indicate that Agentic Predictor+ consistently outperforms all baseline models across all label ratios, demonstrating the effectiveness of our unsupervised pretraining strategy. The gains are especially pronounced in low-label regimes: at a 0.1 label ratio, Agentic Predictor+ maintains an accuracy above 73%, while other models drop closer to 70%. These findings underscore the importance of leveraging cross-domain structure through pretraining for generalizable workflow performance prediction, especially when direct supervision is limited.
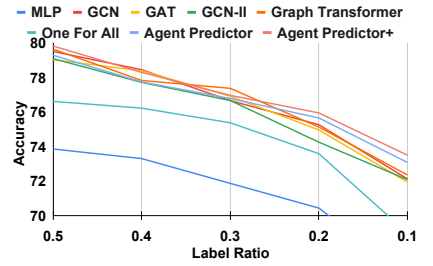


Figure 3: Accuracy comparison between Agentic Predictor and the baselines across varying label ratios.

More detailed results, an additional evaluation of performance predictors used as a reward function for agentic workflow optimization, and case study findings are provided in the Appendix.

## 5 Conclusions

This paper introduces Agentic Predictor, a novel framework for efficient prediction of agentic workflow performance that leverages a multi-view predictive approach. By integrating multi-view graph structures, code semantics, and prompt embeddings into a unified representation, Agentic Predictor captures the diverse characteristics of agentic systems. Moreover, it employs cross-domain unsupervised pretraining to mitigate the challenge of limited labeled data, thereby enhancing generalization across varied tasks. Through comprehensive experiments spanning three domains, Agentic Predictor consistently outperforms strong baselines in predictive accuracy and workflow utility.

**Limitations and Future Work**. While Agentic Predictor exhibits strong performance, it has certain limitations. The current predictor focuses on binary success metrics, constrained by the available benchmark, which may overlook more nuanced aspects of workflow behavior. Additionally, adapting to highly specialized domains may still require some labeled data. Future work includes expanding to multi-objective optimization (e.g., balancing accuracy and cost), incorporating richer views such as temporal traces and user feedback, and exploring human-in-the-loop workflows for real-time refinement. These directions aim to make Agentic Predictor more generalizable and interactive in complex, real-world settings.

# References

Yash Akhauri and Mohamed S Abdelfattah. Encodings for prediction-based neural architecture search. In *Forty-first International Conference on Machine Learning*, 2024. 5, 7

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021. 8

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. 8

Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 1725–1735, 2020. 8

Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *The Twelfth International Conference on Learning Representations*, 2024. 3

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021. 8

Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2024. 3

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2368–2378, 2019. 16

Lukasz Dudziak, Thomas Chau, Mohamed Abdelfattah, Royson Lee, Hyeji Kim, and Nicholas Lane. Brp-nas: Prediction-based nas using gcns. In *Advances in neural information processing systems*, pp. 10480–10490, 2020. 3

Mourad Gridach, Jay Nanavati, Khaldoun Zine El Abidine, Lenon Mendes, and Christina Mack. Agentic ai for scientific discovery: A survey of progress, challenges, and future directions. *arXiv preprint arXiv:2503.08979*, 2025. 1, 3

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. 3

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. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pp. 8048–8057, 2024b. 3

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021a. 8

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021b. 8

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 a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024. 1, 3

Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. In *The Thirteenth International Conference on Learning Representations*, 2025a. 1, 3

Yue Hu, Yuzhu Cai, Yaxin Du, Xinyu Zhu, Xiangrui Liu, Zijie Yu, Yuchen Hou, Shuo Tang, and Siheng Chen. Self-evolving multi-agent networks for software development. In *The Thirteenth International Conference on Learning Representations*, 2025b. 1

Dongyeong Hwang, Hyunju Kim, Sunwoo Kim, and Kijung Shin. Flowerformer: Empowering neural architecture encoding using a flow-aware graph transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6128–6137, 2024. 3

Han Ji, Yuqi Feng, and Yanan Sun. Cap: a context-aware neural predictor for nas. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pp. 4219–4227, 2024. 3

Zixuan Ke, Fangkai Jiao, Yifei Ming, Xuan-Phi Nguyen, Austin Xu, Do Xuan Long, Minzhi Li, Chengwei Qin, Peifeng Wang, Silvio Savarese, et al. A survey of frontiers in llm reasoning: Inference scaling, learning to reason, and agentic systems. *arXiv preprint arXiv:2504.09037*, 2025. 1, 3

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 9

Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017. 8

Hao Liu, Jiarui Feng, Lecheng Kong, Ningyue Liang, Dacheng Tao, Yixin Chen, and Muhan Zhang. One for all: Towards training one graph model for all classification tasks. In *The Twelfth International Conference on Learning Representations*, 2024a. 8

Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. A dynamic LLM-powered agent network for task-oriented agent collaboration. In *First Conference on Language Modeling*, 2024b. 1, 3

Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. GAIA: a benchmark for general AI assistants. In *The Twelfth International Conference on Learning Representations*, 2024. 1

Boye Niu, Yiliao Song, Kai Lian, Yifan Shen, Yu Yao, Kun Zhang, and Tongliang Liu. Flow: Modularized agentic workflow automation. In *The Thirteenth International Conference on Learning Representations*, 2025. 3

Yameng Peng, Andy Song, Vic Ciesielski, Haytham M. Fayek, and Xiaojun Chang. Pre-nas: predictor-assisted evolutionary neural architecture search. In *Proceedings of the Genetic and Evolutionary Computation Conference*, pp. 1066–1074, 2022. 3

Aske Plaat, Max van Duijn, Niki van Stein, Mike Preuss, Peter van der Putten, and Kees Joost Batenburg. Agentic large language models, a survey. *arXiv preprint arXiv:2503.23037*, 2025. 1, 3

Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. ChatDev: Communicative agents for software development. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15174–15186, 2024. 3

Shuofei Qiao, Runnan Fang, Zhisong Qiu, Xiaobin Wang, Ningyu Zhang, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. Benchmarking agentic workflow generation. In *The Thirteenth International Conference on Learning Representations*, 2025. 2

Yu Shang, Yu Li, Keyu Zhao, Likai Ma, Jiahe Liu, Fengli Xu, and Yong Li. Agentsquare: Automatic llm agent search in modular design space. In *The Thirteenth International Conference on Learning Representations*, 2025. 1, 3

Yunsheng Shi, Zhengjie Huang, Shikun Feng, Hui Zhong, Wenjing Wang, and Yu Sun. Masked label prediction: Unified message passing model for semi-supervised classification. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pp. 1548–1554, 8 2021. 8

Tarun Suresh, Revanth Gangi Reddy, Yifei Xu, Zach Nussbaum, Andriy Mulyar, Brandon Duderstadt, and Heng Ji. CoRNStack: High-quality contrastive data for better code retrieval and reranking. In *The Thirteenth International Conference on Learning Representations*, 2025. 9

Patara Trirat and Jae-Gil Lee. PASTA: Neural architecture search for anomaly detection in multivariate time series. *IEEE Transactions on Emerging Topics in Computational Intelligence*, pp. 1–16, 2024. 5, 7

Patara Trirat, Wonyong Jeong, and Sung Ju Hwang. AutoML-agent: A multi-agent LLM framework for full-pipeline autoML. In *Forty-second International Conference on Machine Learning*, 2025. 3

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018. 8

Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in neural information processing systems*, 33:5776–5788, 2020. 9

Colin White, Willie Neiswanger, Sam Nolen, and Yash Savani. A study on encodings for neural architecture search. In *NeurIPS*, pp. 20309–20319, 2020. 5

Colin White, Arber Zela, Binxin Ru, Yang Liu, and Frank Hutter. How powerful are performance predictors in neural architecture search? In *Advances in Neural Information Processing Systems*, 2021. 2, 8

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. AutoGen: Enabling next-gen LLM applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024. 1

Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, Songyang Gao, Lu Chen, Rui Zheng, Yicheng Zou, Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. AgentGym: Evolving large language model-based agents across diverse environments, 2024. 1

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. *Science China Information Sciences*, 68(2):121101, 2025. 1, 3

Frank F. Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Z. Wang, Xuhui Zhou, Zhitong Guo, Murong Cao, Mingyang Yang, Hao Yang Lu, Amaad Martin, Zhe Su, Leander Maben, Raj Mehta, Wayne Chi, Lawrence Jang, Yiqing Xie, Shuyan Zhou, and Graham Neubig. TheAgentCompany: Benchmarking llm agents on consequential real world tasks, 2024. 2

Shen Yan, Yu Zheng, Wei Ao, Xiao Zeng, and Mi Zhang. Does unsupervised architecture representation learning help neural architecture search? *Advances in neural information processing systems*, 33:12486–12498, 2020. 7

Shen Yan, Kaiqiang Song, Fei Liu, and Mi Zhang. Cate: Computation-aware neural architecture encoding with transformers. In *International Conference on Machine Learning*, pp. 11670–11681. PMLR, 2021. 7

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–2380, 2018. 16

Rui Ye, Shuo Tang, Rui Ge, Yaxin Du, Zhenfei Yin, Jing Shao, and Siheng Chen. MAS-GPT: Training LLMs to build LLM-based multi-agent systems. In *Workshop on Reasoning and Planning for Large Language Models*, 2025. 3

Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Dongsheng Li, and Deqing Yang. Evoagent: Towards automatic multi-agent generation via evolutionary algorithms. In *NAACL*, 2025. 1

Guibin Zhang, Yanwei Yue, Xiangguo Sun, Guancheng Wan, Miao Yu, Junfeng Fang, Kun Wang, and Dawei Cheng. G-designer: Architecting multi-agent communication topologies via graph neural networks. In *ICML*, 2024. 3

Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu. AFlow: Automating agentic workflow generation. In *The Thirteenth International Conference on Learning Representations*, 2025a. 1, 3, 8, 16

Yuanshuo Zhang, Yuchen Hou, Bohan Tang, Shuo Chen, Muhan Zhang, Xiaowen Dong, and Siheng Chen. Gnns as predictors of agentic workflow performances. *arXiv preprint arXiv:2503.11301*, 2025b. 3, 6, 8, 9, 16

Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen Schmidhuber. GPTSwarm: Language agents as optimizable graphs. In *Forty-first International Conference on Machine Learning*, 2024. 1, 3
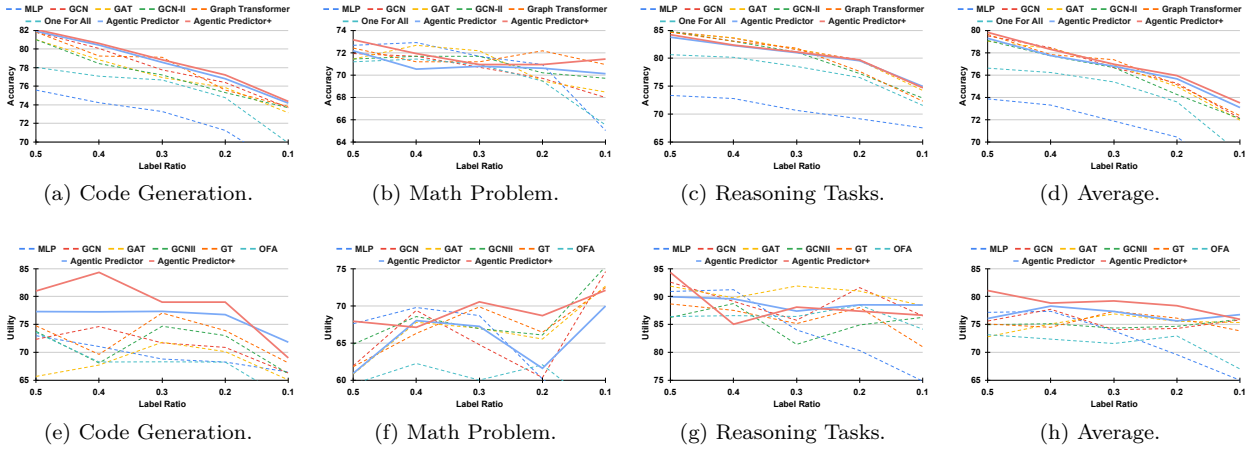
Figure 4: Comparison of accuracy (upper) and utility (lower) between Agentic Predictor and the baselines across varying label ratios.

# A Additional Experimental Results

## A.1 Effects of Pretraining Phase (Full Results)

Since acquiring a large amount of ground-truth labels from agentic workflows is expensive, we examine whether cross-domain unsupervised pretraining (denoted as Agentic Predictor+) benefits settings where labeled instances are limited. We vary the label ratio from 0.1 to 0.5, selecting labeled samples from the training split of all datasets in the benchmark. We pretrain the proposed multi-view encoder with a batch size of 32 for 20 epochs.

Following the average results in the main text, we provide a comprehensive comparison of accuracy (top row) and utility (bottom row) across three task domains—code generation, math problems, and reasoning—under varying label ratios from 0.5 to 0.1 (Figure 4).

Across all settings, our proposed framework, Agentic Predictor, and its pretrained variant, Agentic Predictor+, consistently outperform baseline models, especially in low-resource scenarios. In the code generation domain (Figures 4a, 4e), Agentic Predictor+ achieves superior accuracy and notably higher utility as the label ratio decreases, outperforming all graph-based and non-graph baselines. Similarly, for math problems (Figures 4b, 4f), Agentic Predictor+ maintains a stable accuracy even as labeled data diminishes, while significantly improving utility, indicating better performance in label-scarce conditions. In reasoning tasks (Figures 4c, 4g), although accuracy deltas narrow between models, Agentic Predictor+ sustains strong utility across all label ratios, highlighting its robustness in generalization. When averaged across domains (Figures 4d, 4h), Agentic Predictor+ shows clear advantages in both metrics under limited supervision. The utility improvements are particularly prominent, suggesting that our pretrained encoder captures transferable representations that enhance decision-making, even when fine-tuning data is sparse. These findings validate the efficacy of the unsupervised pretraining phase and highlight the importance of cross-domain datasets for pretraining.

## A.2 Workflow Optimization Results

In the main experiments, we demonstrate the feasibility and robustness of predicting agentic workflow performance. However, it remains an open question whether such predictions can effectively contribute to improving efficiency and to what extent they may introduce performance degradation in agentic workflows. To investigate this, we evaluate **(Q4)** *whether using Agentic Predictor as a predictor enhances the optimization of agentic workflows compared to alternative baselines.* Specifically, we measure the performance improvement (or loss) incurred when using performance predictors.

Table 6: Workflow optimization performance based on the selected workflow across different methods.

| Methods | Math Problems | | Code Generation | | Reasoning | | | Average | |
|---|---|---|---|---|---|---|---|---|---|
| | MATH | GSM8K | MBPP | HumanEval | MMLU | DROP | HotpotQA | Score | Search Cost ($) |
| Ground Truth (AFlow) | 87.38 | 94.53 | 73.22 | 97.20 | 83.10 | 84.25 | 69.94 | 84.23 | 39.83 |
| Random | 78.40 | 75.23 | 67.84 | 76.34 | 42.87 | 80.42 | 16.86 | 62.56 | 0.00 |
| GCN | 79.22 | 86.16 | 68.23 | 97.46 | 46.43 | 82.33 | 19.14 | 68.42 | 0.00 |
| GAT | 80.11 | 86.22 | **68.62** | 97.71 | 57.00 | 85.83 | **21.47** | 71.00 | 0.00 |
| *Agentic Predictor* | **81.89** | **92.65** | 68.42 | **98.73** | **79.70** | **86.25** | 13.37 | **74.43** | 0.00 |

To ensure a fair comparison, we adopt the same experimental setup as Zhang et al. (2025b), which provides a unified platform for optimizing agentic workflows and evaluating their performance. During the optimization process on each benchmark, a predictor is used to estimate the performance of candidate agentic workflows. These predicted performance values are treated as rewards to guide the optimization. Upon completion of the optimization, the quality of the resulting workflows is assessed based on their accuracy score on held-out test tasks.

We compare Agentic Predictor against four baselines: (1) the **ground truth** baseline, which directly evaluates agentic workflows to obtain ground-truth performance scores (as done in the original AFlow (Zhang et al., 2025a)); (2) two strong GNN-based predictors **GCN** and **GAT**; and (3) a **random** baseline, which assigns random performance scores as rewards. This experiment is conducted across five benchmarks: MATH, GSM8K, MBPP, HumanEval, and MMLU.

As in Table 6, Agentic Predictor consistently outperforms the random, GCN, and GAT baselines, achieving an average accuracy score of **74.43%**, significantly higher than random (62.56%), GCN (68.42%), and GAT (71.00%). Notably, as a predictor incurs zero search cost compared to the ground-truth's cost of $39.83, this result underscores the effectiveness and efficiency of Agentic Predictor as a reliable predictor for optimizing agentic workflows.

### A.3 Transferability

Considering that the MMLU benchmark encompasses various reasoning tasks, we further investigate the transferability of predictors trained on MMLU datasets to determine whether they can be used to optimize similar reasoning tasks, specifically DROP (Dua et al., 2019) and HotpotQA (Yang et al., 2018). As reported in Table 6, the workflow optimized using Agentic Predictor achieves competitive performance on these tasks: 86.25% on DROP and 13.37% on HotpotQA, demonstrating notable transferability. While performance on HotpotQA is lower than the baselines, the results remain broadly comparable, indicating that the workflows optimized via Agentic Predictor maintain substantial effectiveness when transferred to closely related reasoning tasks. This highlights the practical potential of Agentic Predictor for broader applicability in workflow optimization scenarios.

## B  Case Study

This section presents qualitative results from the workflow optimization process using Agentic Predictor as the reward function across three domains.

### B.1  Code Generation

The code generation workflow on the HumanEval dataset demonstrates that the initial solution generation step often required subsequent refinement through explicit review and revision cycles. By systematically reviewing the initially generated code, and conditionally revising based on feedback from automated tests, the workflow substantially improved the final solution's correctness. This iterative approach effectively balanced computational cost and performance, resulting in solutions that were consistently more robust and accurate compared to single-step generations.

**Workflow for Code Generation (HumanEval)**

```python
from typing import Literal
import workspace.HumanEval.workflows.template.operator as operator
import workspace.HumanEval.workflows.round_19.prompt as prompt_custom
from metagpt.provider.llm_provider_registry import create_llm_instance
from metagpt.utils.cost_manager import CostManager

DatasetType = Literal["HumanEval", "MBPP", "GSM8K", "MATH", "HotpotQA", "DROP", "MMLU"]


class Workflow:
    def __init__(
        self,
        name: str,
        llm_config,
        dataset: DatasetType,
    ) -> None:
        self.name = name
        self.dataset = dataset
        self.llm = create_llm_instance(llm_config)
        self.llm.cost_manager = CostManager()
        self.custom = operator.Custom(self.llm)
        self.custom_code_generate = operator.CustomCodeGenerate(self.llm)
        self.test = operator.Test(self.llm)

    async def __call__(self, problem: str, entry_point: str):
        """
        Implementation of the workflow
        1. Generate initial solution using custom_code_generate.
        2. Review the solution using custom operator.
        3. Test the solution; if test fails, revise using custom operator and retest.
        """
        # Step 1: Generate initial solution
        initial_solution = await self.custom_code_generate(problem=problem, entry_point=entry_point, instruction="")

        # Step 2: Review the solution to improve quality
        reviewed = await self.custom(input=initial_solution['response'], instruction=prompt_custom.REVIEW_PROMPT)

        # Step 3: Test the reviewed solution
        test_result = await self.test(problem=problem, solution=reviewed['response'], entry_point=entry_point)

        # If test fails, revise solution based on test feedback and retest once
        if not test_result['result']:
            revised = await self.custom(input=reviewed['response'] + "\n" + test_result['solution'], instruction=
                prompt_custom.REVISE_PROMPT)
            test_result = await self.test(problem=problem, solution=revised['response'], entry_point=entry_point)
            final_solution = revised['response'] if test_result['result'] else reviewed['response']
        else:
            final_solution = reviewed['response']

        return final_solution, self.llm.cost_manager.total_cost
```

### B.2 Math Problem

In addressing mathematical problems using the MATH dataset, the workflow leverages an ensemble strategy by producing multiple candidate solutions, subsequently selecting the most consistent one via a self-consistency ensemble step. The selected solution was then further refined through an additional review process. This combined ensemble and review mechanism significantly enhanced solution quality, highlighting the value of ensemble techniques in solving complex mathematical reasoning tasks, while maintaining a controlled computational budget.

**Workflow for Math Problem (MATH)**

```python
from typing import Literal
import workspace.MATH.workflows.template.operator as operator
import workspace.MATH.workflows.round_88.prompt as prompt_custom
from metagpt.provider.llm_provider_registry import create_llm_instance
from metagpt.utils.cost_manager import CostManager

DatasetType = Literal["HumanEval", "MBPP", "GSM8K", "MATH", "HotpotQA", "DROP", "MMLU"]
```

```python
class Workflow:
    def __init__(
        self,
        name: str,
        llm_config,
        dataset: DatasetType,
    ) -> None:
        self.name = name
        self.dataset = dataset
        self.llm = create_llm_instance(llm_config)
        self.llm.cost_manager = CostManager()
        self.custom = operator.Custom(self.llm)
        self.sc_ensemble = operator.ScEnsemble(self.llm)

    async def __call__(self, problem: str):
        """
        Implementation of the workflow with ensemble and review step
        """
        # Generate multiple candidate solutions using custom operator with different instructions
        candidates = []
        for i in range(3):
            response = await self.custom(input=problem, instruction=prompt_custom.SOLVE_PROMPT + f" Attempt {i+1}.")
            candidates.append(response['response'])

        # Use self-consistency ensemble to select the best solution
        ensemble_result = await self.sc_ensemble(solutions=candidates, problem=problem)
        best_solution = ensemble_result['response']

        # Review and refine the best solution
        review_response = await self.custom(input=problem + "\nSolution to review:\n" + best_solution, instruction=
            prompt_custom.REVIEW_PROMPT)
        final_solution = review_response['response']

        return final_solution, self.llm.cost_manager.total_cost
```

## B.3    Reasoning Task

For reasoning tasks on the MMLU dataset, the workflow combines multiple generation techniques, including custom-generated solutions with varying prompts and answers produced by specialized answer-generation operators, to diversify initial candidate answers. The self-consistency ensemble step effectively selected the most consistent candidate, which was subsequently subjected to rigorous review and format verification steps. This meticulous process, which included conditional regeneration and revision to ensure strict adherence to specified answer formats, proved highly effective in enhancing both accuracy and reliability of the final responses.

### Workflow for Reasoning Task (MMLU)

```python
from typing import Literal
import workspace.MMLU.workflows.template.operator as operator
import workspace.MMLU.workflows.round_19.prompt as prompt_custom
from metagpt.provider.llm_provider_registry import create_llm_instance
from metagpt.utils.cost_manager import CostManager

DatasetType = Literal["HumanEval", "MBPP", "GSM8K", "MATH", "HotpotQA", "DROP", "MMLU"]

class Workflow:
    def __init__(
        self,
        name: str,
        llm_config,
        dataset: DatasetType,
    ) -> None:
        self.name = name
        self.dataset = dataset
        self.llm = create_llm_instance(llm_config)
        self.llm.cost_manager = CostManager()
        self.custom = operator.Custom(self.llm)
        self.answer_generate = operator.AnswerGenerate(self.llm)
        self.sc_ensemble = operator.ScEnsemble(self.llm)

    async def __call__(self, problem: str):
```

```
"""
Implementation of the workflow with multiple custom answers, multiple AnswerGenerate answers, ensemble, review, and
    revision
"""
# Step 1: Generate multiple candidate answers using custom operator with a concise prompt
custom_answers = []
for _ in range(2):
    custom_response = await self.custom(input=problem, instruction=prompt_custom.CUSTOM_PROMPT)
    custom_answer = custom_response['response']
    custom_answers.append(custom_answer)
# Add 1 answer with diversity prompt to increase answer variety
custom_diverse_response = await self.custom(input=problem, instruction=prompt_custom.CUSTOM_DIVERSE_PROMPT)
custom_answers.append(custom_diverse_response['response'])

# Step 2: Generate multiple candidate answers using AnswerGenerate operator to increase diversity
answergen_answers = []
for _ in range(2):
    answergen_response = await self.answer_generate(input=problem)
    answergen_answer = answergen_response['answer']
    answergen_answers.append(answergen_answer)

# Step 3: Ensemble all candidate answers to select the most consistent answer
all_answers = custom_answers + answergen_answers
ensemble_response = await self.sc_ensemble(solutions=all_answers)
ensemble_answer = ensemble_response['response']

# Step 4: Review the ensemble answer to ensure format and correctness
review_input = problem + "\nAnswer: " + ensemble_answer
review_response = await self.custom(input=review_input, instruction=prompt_custom.REVIEW_PROMPT)
reviewed_answer = review_response['response']

# Step 5: If reviewed answer is not in correct format, regenerate with a stricter prompt
if not reviewed_answer.startswith("Answer: Option "):
    strict_regen_input = problem + "\nPlease provide the final answer strictly in the format 'Answer: Option X'."
    strict_regen_response = await self.custom(input=strict_regen_input, instruction=prompt_custom.
        STRICT_REGEN_PROMPT)
    reviewed_answer = strict_regen_response['response']

# Step 6: Revision step to refine the reviewed answer for strict format adherence
revision_input = problem + "\nAnswer: " + reviewed_answer
revision_response = await self.custom(input=revision_input, instruction=prompt_custom.REVISION_PROMPT)
final_answer = revision_response['response']

return final_answer, self.llm.cost_manager.total_cost
```