

Do We Always Need Query-Level Workflows? Rethinking Agentic Workflow Generation for Multi-Agent Systems

Anonymous ACL submission

Abstract

Multi-Agent Systems (MAS) built on large language models typically solve complex tasks by coordinating multiple agents through workflows. Existing approaches generates workflows either at task level or query level, but their relative costs and benefits remain unclear. After rethinking and empirical analyses, we show that query-level workflow generation is not always necessary, since a small set of top-K best task-level workflows together already covers equivalent or even more queries. We further find that exhaustive execution-based task-level evaluation is both extremely token-costly and frequently unreliable. Inspired by the idea of self-evolution and generative reward modeling, we propose a low-cost task-level generation framework **SCALE**, which means Self prediction of the optimizer with few shot CALibration for Evaluation instead of full validation execution. Extensive experiments demonstrate that **SCALE** maintains competitive performance, with an average degradation of just 0.61% compared to existing approach across multiple datasets, while cutting overall token usage by up to 83%.

1 Introduction

Large Language Model (LLM)-based multi-agent systems (MAS) have recently emerged as a powerful paradigm for solving complex reasoning, coding, and decision-making tasks (Zhang et al., 2024b,a; Zhuge et al., 2024; Niu et al., 2025). By decomposing a task into multiple interacting agents and organizing their collaboration through agentic workflows, MAS can substantially extend the capabilities of a single agent.

Based on the granularity of workflow construction, agentic workflow generation methods fall into two categories: task-level and query-level approaches. Task-level approaches, such as such as search-based Aflow(Zhang et al., 2024b) and learning-based GPTSwarm(Zhuge et al., 2024) and

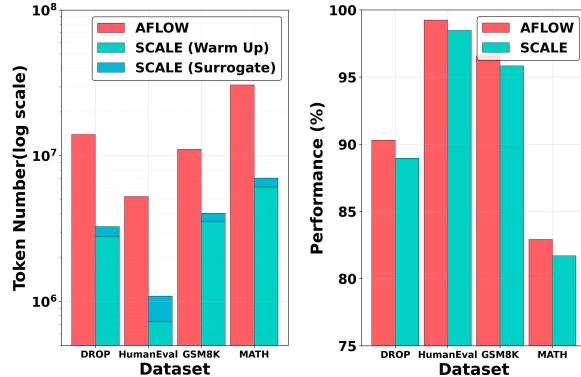


Figure 1: Comparison of Aflow and our rethought task-level workflow generation framework. Left: total token number during workflow generation (log-scale axis). Right: final test performance. Our method **SCALE** achieves comparable performance while significantly reducing token number.

AgentPrune(Zhang et al., 2024a), generate a single workflow intended to perform well across an entire dataset or task distribution. However, this generality comes at a high cost: evaluation dominates computation, as each candidate requires full execution over the validation set. As shown in Figure 1, the whole Aflow’s generation process on four benchmarks consumes approximately 10⁶–10⁸ LLM tokens.

In parallel, query-level approaches generate a separate workflow for each input query(Ye et al., 2025; Wang et al., 2025a; Gao et al., 2025). For every query, the system constructs a customized multi-agent workflow, allowing the agent roles and interaction patterns to adapt to the specific problem. This design aims to better handle heterogeneous queries and can yield strong per-query performance. However, this adaptivity comes with clear costs. A new workflow must be generated for every query, which introduces substantial inference overhead. For many simple or similar inputs, such query-level generation may be unnecessary, providing limited gains relative to its train time computational cost

066 and test time generation cost.
067

068 From the characteristics of these two paradigms,
069 we raise two fundamental questions that have not
070 been systematically examined as shown in Figure 2. First, *Is query-level workflow generation*
071 *always necessary?* Second, *Is high-cost evalua-*
072 *tion in task-level workflow generation necessary?*
073 For the first question, we show that a small set of
074 top-k task-level workflows already achieves strong
075 query coverage comparing to query-level method’s
076 performance. It indicates that query-level work-
077 flow generation is not always necessary in practice.
078 For the second, we find that exhaustive execu-
079 tion-based evaluation of task-level workflows is both
080 extremely expensive and frequently unreliable.

081 Inspired by the idea of self-evolution and gen-
082 erative reward modeling, we propose a low-cost
083 task-level generation framework **SCALE**, which
084 means Self prediction of the optimizer with few
085 shot CALibration for Evaluation instead of full
086 validation execution. By leveraging the inher-
087 ent evaluative ability of LLM-based optimizers,
088 **SCALE** makes self predictions in a generative man-
089 ner and calibrates them using few shot executions,
090 thereby achieving highly reliable predictions with
091 minimal token cost. Experimental analysis fur-
092 ther shows that the calibrated self predictions in
093 **SCALE** closely approximate true execution scores,
094 it achieves a low MAE of 0.16 and maintain con-
095 sistent ranking with a Pearson correlation of 0.52
096 (range: $[-1, 1]$), further validating reliability. Over-
097 all, our contributions are threefold as shown below:

- **Rethinking Insights:** We present a new empirical rethinking of workflow generation in multi-agent systems. Our analysis yields two main findings: (1) Query-level methods is not always necessary in practice. (2) Exhaustive execution-based evaluation in task-level approaches is both costly and unreliable.
- **Improved Framework:** Motivated by these observations and inspired by self-evolution, we develop a low-cost and effective framework **SCALE** for task-level workflow generation. Instead of exhaustively executing candidate workflows on the full validation set, our approach combines the LLM-based optimizer’s self prediction with few shot calibration to evaluate workflows efficiently.
- **Empirical Validation:** Extensive experiments demonstrate that **SCALE** maintains

116 competitive performance, with an average
117 degradation of just 0.61% compared to exist-
118 ing approach across multiple datasets, while
119 cutting overall token usage by up to 83%.

2 Preliminaries

2.1 Agentic MAS Workflow

We consider a task \mathcal{T} given as a dataset of queries $q \in \mathcal{D}$, and an agentic MAS workflow $W \in \mathcal{W}$ that orchestrates a LLM-based MAS to produce an answer y . Formally, we formalizes the workflows space as:

$$\mathcal{W} = \{(P_1, \dots, P_n, E, O_{\theta_1}, \dots, O_{\theta_n}) \mid P_i \in \mathcal{P}, E \in \mathcal{E}, O_{\theta_i} \in \mathcal{O}\} \quad (1)$$

where n is the number of agents, \mathcal{P} is the prompt space, \mathcal{E} is the information control flow that governs the execution order, data dependencies, and communication among these agents, which can be represented in various forms: as executable code (e.g. Hu et al., 2024; Zhang et al., 2024b; Xu et al., 2025) or as directed acyclic graphs(e.g. Zhuge et al., 2024; Zhang et al., 2024a; Wang et al., 2025b; Zhang et al., 2025). \mathcal{O} is a set of predefined LLM-based agents parameterized by O_{θ_i} . The prompt P_i and parameters θ_i enable the agent to adapt its behaviors to the task at hand, such as *Review*(Yao et al., 2022), *Ensemble*(Liang et al., 2024), or *Self-Correction*(Shinn et al., 2023).

At a high level, a workflow is a series of agent calls to answer a certain query $y = W(q)$. Given a task-specific evaluator $s(\cdot, \cdot)$ (e.g., exact match, pass@1), the performance of W on a query q is measured by $s(W(q), q) \in [0, 1]$.

2.2 Agentic MAS Workflow Generation

A number of recent methods have been proposed for agentic workflow generation, which can be grouped into two paradigms according to the granularity at which workflows are generated: task-level approaches and query-level approaches.

2.2.1 Task-level workflow generation

Task-level methods aim to generate a single workflow W^* that performs well on a distribution of queries from the same task as shown in (1)A and (1)B of Figure 2. Given a validation set $\mathcal{D}_{val} = \{q_i\}_{i=1}^N$ sampled from the task \mathcal{T} , the objective is:

$$W_{task}^* = \arg \max_{W \in \mathcal{W}} S^{exec}(W, \mathcal{D}_{val}) \quad (2)$$

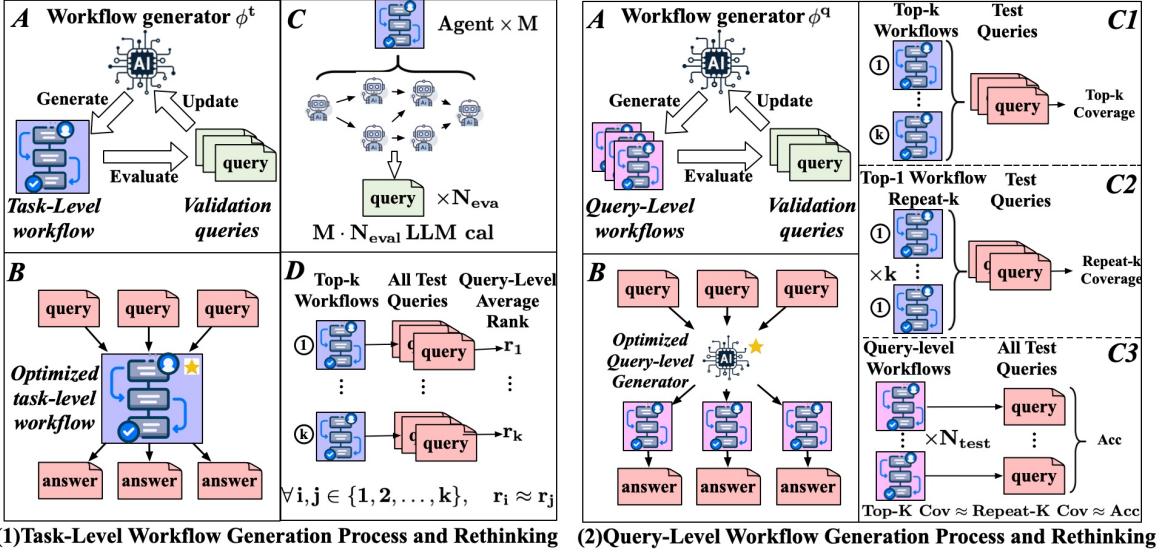


Figure 2: Task-level vs. Query-level workflow generation on their process and rethinking.(1)Task-level generation.(1)A shows searching/training: the generator generates a single workflow using validation queries; (1)B shows inference: the optimized workflow is reused for all test queries. (1)C–D present our rethinking: (1)C shows that repeated full-set evaluations is very token-costly, and (1)D shows top-k workflows have very similar query-level ranks. (2)Query-level generation. (2)A shows training: a workflow is generated per query; (2)B shows inference: producing customized workflows for each input. (2)C1–C3 summarize our rethinking on query-level workflows: top-k task-level workflows, repeat-k runs of the top-1 workflow, and true query-level generation yield comparable coverage/performance.

where \mathcal{W} is the workflow search space, and $S^{exec}(W, \mathcal{D}_{val}) = \frac{1}{|\mathcal{D}_{val}|} \sum_{q \in \mathcal{D}_{val}} s(W(q), q)$ denotes the average execution score of workflow W on dataset \mathcal{D}_{val} .

Despite implementation differences, these methods share the same closed-loop paradigm. First, generating candidate workflows through a task-level workflow generator ϕ^t . Second, evaluating the generated workflow on \mathcal{D}_{val} . Finally, updating the model ϕ^t using evaluation as a feedback. Aflow (Zhang et al., 2024b) uses a LLM as optimizer ϕ^t and improve the initial workflow through an MCTS-style loop. AgentPrune (Zhang et al., 2024a) use a graph model as workflow generator ϕ^t and learn it through reinforcement learning methods to generate better workflows.

Despite good performance, their evaluation is extremely costly, since each candidate’s evaluation requires the MAS workflow to execute on the full validation set. The token number scales with the validation dataset size and agent numbers resulting in a sharp increase as the loop continues.

2.2.2 Query-level workflow generation

As shown in (2)A and (2)B of Figure 2, query-level methods instead learn a query-level workflow generator ϕ^q that maps each query q to its own workflow $W_q = \phi^q(q)$ and the optimization objec-

tive is to maximize expected performance over the validation set:

$$\phi^{q,*} = \arg \max_{\phi^q} \frac{1}{|\mathcal{D}_{val}|} \sum_{q \in \mathcal{D}_{val}} [s(W_q(q), q)] \quad (3)$$

Existing approaches differ mainly in how ϕ^q is trained. MAS-GPT (Ye et al., 2025) adopts supervised fine-tuning on a curated dataset of query–workflow pairs. ScoreFlow (Wang et al., 2025a) optimizes the workflow generator ϕ^q using a preference-based optimization approach which enhances the original DPO (Rafailov et al., 2023). For many simple or structurally similar inputs, such query-level workflow generation may be unnecessary, providing limited gains relative to its test-time generation cost.

3 Rethinking Agentic Workflow Generation for Multi-Agent Systems

Our rethinking is twofold. First, as shown in (2)C1–C3 of Figure 2, we rethink the necessity of query-level workflow generation. Second, as shown in (1)C and (1)D of Figure 2, we rethink task-level methods, arguing that execution on validation set for evaluation is both token-costly and unreliable.

Dataset	Aflow			S.Flow		
	Top-1 Perf	Top-5 Perf	Cov	Repeat-5 Perf	Cov	All Perf
DROP	90.30	89.81	93.87	90.04	92.45	91.48
HumanEval	99.24	98.17	100.00	97.82	99.24	98.91
GSM8K	96.58	95.89	97.35	96.17	96.87	97.79
MATH	82.92	79.84	87.04	82.81	83.39	84.35

Table 1: Comparison of task-level and query-level workflow effectiveness. **Perf** denotes average test performance (%). **Cov** denotes coverage (%) of test queries. **S.Flow** denotes the query-level method ScoreFlow.

3.1 Is Query-level Workflow Generation Always Necessary?

Query-level methods offers fine-grained adaptivity, but it introduces considerable test-time generation cost that task-level methods don’t have. This raises a fundamental question:

Is query-level workflow generation always necessary to achieve strong performance?

For each dataset, we evaluate the following settings.(1)Task-level Top-1: We report the test performance of Aflow’s (Zhang et al., 2024b) best generated workflow. (2)Task-level Top-5: We report average test performance and coverage of Aflow’s top-5 workflows. *Coverage* is defined as the fraction of test queries that are covered by these workflows. A query is counted as covered if at least one of these workflows answers it correctly. (3)Repeat-5 of Top-1: We execute the Aflow’s top-1 workflow on test set 5 times. We report the average performance and coverage on test queries. This isolates the effect of the test-time execution stochasticity. (4)Query-level workflows: A dedicated workflow is generated for each query using a representative query-level method ScoreFlow (Wang et al., 2025a), and we report its average test performance.

As shown in Table 1, we obtain three main observations. First, a single Top-1 task-level workflow already performs strongly, indicating substantial structural sharing across queries. Second, Top-5 gains a clear increase in query coverage even more than query-level method’s performance, meaning that improvements can come from gathering few candidate task-level workflows rather than generating single strictly better workflows. Third, Repeat-5 achieves coverage comparable to query-level method, showing that much of the gain by query-level method compared to task-level can be covered by stochastic execution rather than diverse workflow structures.

Taken together, these results suggest that most benefits attributed to query-level method can already be achieved by a small pool of reusable task-

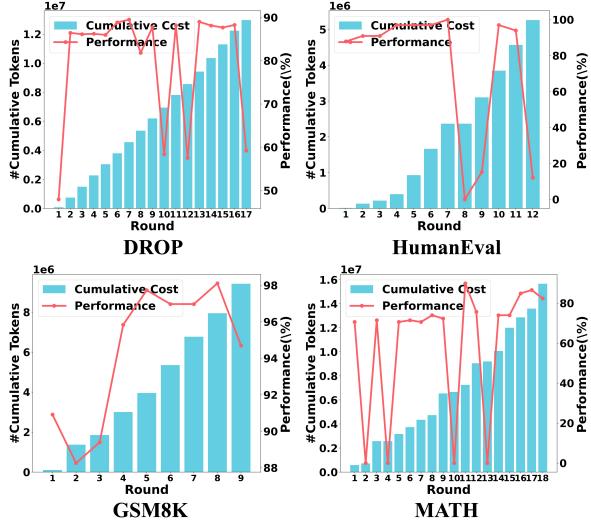


Figure 3: Cumulative Token Number v.s. Performance during Aflow’s task-level workflow generation process.

level workflows or even repeated execution of a single strong task-level workflow. The main advantage can come from coverage and stochasticity, not necessarily in query-level workflows.

3.2 Is High-Cost Evaluation in Task-Level Workflow Generation Necessary?

In this subsection, we revisit the evaluation in task-level workflow generation from two complementary perspectives: (1) how many tokens actually incurs in exhaustive execution-based evaluation, and (2) whether such costly evaluation truly leads to meaningfully different workflows.

3.2.1 How many evaluation tokens are incurred during task-level workflow generation?

We analyze the cumulative evaluation token number incurred during Aflow’s workflow generation. Figure 3 illustrates the relationship between evaluation token number and performance across each generated candidate workflow. For each benchmark, we plot the test performance achieved by candidate workflows, together with the cumulative token number which is defined as the total tokens used to evaluate all candidate workflows generated up to and including the current round.

Figure 3 shows that cumulative evaluation token number grows fast continuously with search rounds, while test performance quickly saturates and yields only marginal or negative gains afterward. This mismatch is evident: the cost is exploding but the corresponding performance improvement is minimal or even negative. These results indicate that the prevailing task-level evaluation paradigm is both

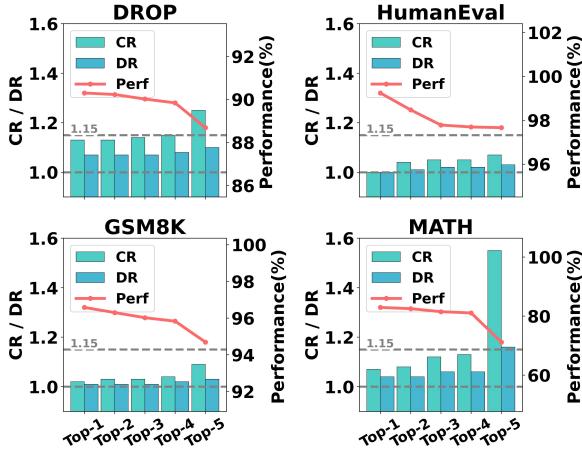


Figure 4: Performance and ranking statistics of the top-5 task-level workflows generated by Aflow across four benchmarks. **Perf** denotes average test performance. **CR** and **DR** denote average competition rank and dense rank, respectively, computed over test queries.

expensive and unreliable in the high-performance regime, motivating the need for cheaper and more reliable task-level workflow evaluation.

3.2.2 Do high-cost task-level evaluations actually distinguish better workflows?

To further examine the efficacy of Aflow’s costly evaluation process, we analyze the Top-5 workflows, defined as the five candidates with the highest validation performance across all candidate workflows produced in one complete run.

Beyond test performance, we also analyze query-level ranking. Concretely, all Top-5 workflows are executed and ranked for each test query. Then Top-5 workflows’ competition rank (**CR**) which allows ties and dense rank (**DR**) are averaged across queries. This reveals how consistently one workflow beats another. If the evaluation are strongly discriminative, these ranks would show clear separation among Top-5 workflows.

Figure 4 shows that the Top-5 (especially Top-4) task-level workflows obtained by Aflow exhibit very similar performance. Their performances vary only slightly across all benchmarks, while both CR and DR remain close to 1 with minimal variation, which indicates that expensive full validation provides limited benefit in identifying substantially better workflows.

3.3 Rethinking Results

In Section 3.1, we observed that query-level workflow generation is not always necessary, because a small set of task-level workflows or even repeated executions of a single workflow already covers

more queries. In Section 3.2, we further showed that exhaustive execution-based task-level evaluation is extremely costly while providing limited discriminative benefit.

Together, these findings show that current agentic workflow methods waste a lot of computation either generating unnecessary query-level workflows or evaluating task-level workflows that have almost the same performances. This motivates a new framework for task-level workflow generation that avoids both query-level generation and high-cost full-execution-based evaluation.

4 Methodology

4.1 Motivation

We propose a low-cost task-level generation framework **SCALE**, which means Self prediction with few shot CALibration for Evaluation. Inspired by the self-evolution paradigm and generative reward modeling in agentic systems, we treat the workflow optimizer as a self-predictor. In other words, the same LLM that generates workflows is also prompted to estimate the candidate workflow’s expected performance.

To reduce overconfidence, we separate generation and evaluation prompts and obtain scores in a dedicated evaluation context. We calibrate self predictions using execution results from few shot queries, typically 1–3% of the validation set. As a result, **SCALE** enables task-level evaluation without full validation execution, while maintaining reliable workflow scoring and ranking.

4.2 Overall framework

Our method **SCALE** operates in two stages: a short warm-up stage with full execution for evaluation, followed by a surrogate evaluation stage. Specifically, we use Aflow (Zhang et al., 2024b) for a few steps as the warm-up stage, and then instead of repeatedly computing $S^{\text{exec}}(W, \mathcal{D}_{\text{val}})$, we estimate workflow quality using a self prediction score calibrated by few shot execution as the surrogate evaluation stage.

4.2.1 Warm-up Stage

We begin with M warm-up rounds. Concretely, starting from a single question-answering LLM agent call as the initial workflow W_1 with execution score S_1^{exec} , we run an MCTS-style loop for $t = \{1, \dots, M\}$ consisting of four steps:

1. Selection. Given the existing workflows and scores $\{S_i^{\text{exec}}\}_{i=1}^t$, we select a parent workflow

365 using a soft mixed policy:

$$P_i = \lambda \cdot \frac{1}{t} + (1-\lambda) \cdot \frac{\exp(\alpha(S_i^{\text{exec}} - S_{\max}^{\text{exec}}))}{\sum_{j=1}^t \exp(\alpha(S_j^{\text{exec}} - S_{\max}^{\text{exec}}))}, \quad (4)$$

366 where $S_{\max}^{\text{exec}} = \max_j S_j^{\text{exec}}$, and λ, α control the
367 exploration-exploitation trade-off.

368 **2. Expansion.** A new workflow is generated by
369 editing W_i using the LLM-based optimizer:
370

$$371 \quad W_{t+1} = \phi(W_i; P_{t+1}^{\text{optimizer}}) \quad (5)$$

372 The optimizer prompt $P_{t+1}^{\text{optimizer}}$ is dynamically
373 built from local experience E_i^{local} and global expe-
374 rience E^{global} introduced in backpropagation step.

375 **3. Evaluation.** We execute W_{t+1} on the valida-
376 tion set $S_{t+1}^{\text{exec}} = \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{q \in \mathcal{D}_{\text{val}}} s(W_{t+1}(q), q)$.

377 **4. Backpropagation** The evaluation result up-
378 dates both local and global experience. For W_i
379 the local experience is $E_i^{\text{local}} \leftarrow E_i^{\text{local}} \cup (e_i^{t+1})$
380 where $e_i^{t+1} = ((W_i, S_i), \Delta_i^{t+1}, (W_{t+1}, S_{t+1}))$
381 and Δ_i^{t+1} denotes the optimizer's natural language
382 description of the edit from W_i to W_{t+1} . We
383 also maintain a global experience: $E^{\text{global}} \leftarrow$
384 $E^{\text{global}} \cup \{(W_{t+1}, S_{t+1}^{\text{exec}})\}$. These experience is
385 reused to update future optimizer prompts and the
386 selection policy.

387 4.2.2 Surrogate Evaluation Stage

388 After warm-up stage, we continue the loop but the
389 subsequent workflows are evaluated through self
390 prediction with few shot execution calibration. At
391 iteration $t > M$, we select and expand using Equa-
392 tion 4 and Equation 5 to get the newly generated
393 workflow W_{t+1} . We evaluate it with our method:

394 **Self prediction** Firstly, we query the optimizer
395 itself under a dynamically dedicated evalua-
396 tion prompt $P_{t+1}^{\text{optimizer}}$ to obtain the prediction:

$$397 \quad S_{t+1}^{\text{pred}} = S^{\text{pred}}(W_{t+1}) = \phi(W_{t+1}; P_{t+1}^{\text{eval}}) \quad (6)$$

398 where ϕ is the same LLM-based optimizer used for
399 workflow expansion in Equation 5. The evaluation
400 prompt P_{t+1}^{eval} is separated from the optimization
401 prompt $P_{t+1}^{\text{optimizer}}$ to reduce overconfidence and
402 prompt entanglement. The full template of P^{eval} is
403 provided in Appendix A.1.

404 **Few shot execution calibration** To reduce bias,
405 we execute W_{t+1} only on a small subset $\mathcal{D}_{\text{few}} \subset$

406 \mathcal{D}_{val} with $|\mathcal{D}_{\text{few}}| \ll |\mathcal{D}_{\text{val}}|$:

$$407 \quad S_{t+1}^{\text{few}} = \frac{1}{|\mathcal{D}_{\text{few}}|} \sum_{q \in \mathcal{D}_{\text{few}}} s(W_{t+1}(q), q) \quad (7)$$

408 The sampling of \mathcal{D}_{few} is guided by the full-
409 execution statistics collected during warm-up. With
410 warming up workflows $\{W_m\}_{t=1}^M$, for each vali-
411 dation query $q \in \mathcal{D}_{\text{val}}$, we compute its empiri-
412 cal warm-up score $\bar{s}(q) = \frac{1}{M} \sum_{t=1}^M s(W_m(q), q)$
413 which reflects how well warm-up workflows al-
414 ready solve q . We then partition the range of
415 $\bar{s}(q)$ into K bins $\{B_k\}_{k=1}^K$ (e.g., $K = 10$), where
416 $B_k = \{q \in \mathcal{D}_{\text{val}} \mid \bar{s}(q) \in I_k\}$ and $\{I_k\}_{k=1}^K$ are dis-
417 joint score intervals covering $[0, 1]$. Let $n_k = |B_k|$
418 be the number of queries in bin k . We define a bin-
419 level sampling distribution by a softmax over bin
420 counts $p_k = \frac{\exp(\gamma n_k)}{\sum_{j=1}^K \exp(\gamma n_j)}$, $k = 1, \dots, K$, where
421 $\gamma > 0$ is the sampling temperature.

422 Given a target budget $|\mathcal{D}_{\text{few}}| = \rho |\mathcal{D}_{\text{val}}|$ with
423 $\rho \in [0.01, 0.03]$, we sample queries without re-
424 placement by first sampling a bin index k from the
425 categorical distribution $\{p_k\}$, and then sampling a
426 query q uniformly from B_k . Repeating this pro-
427 cedure until $|\mathcal{D}_{\text{few}}|$ is reached yields a few shot subset
428 that preserves the warm-up difficulty distribution
429 while covering both easy and hard queries.

430 **Calibrated surrogate score** The final score used
431 to evaluate the workflow is:

$$432 \quad \hat{S}_{t+1} = (1 - \alpha_{t+1}) S_{t+1}^{\text{pred}} + \alpha_{t+1} S_{t+1}^{\text{few}} \quad (8)$$

433 where $\alpha_{t+1} \in [0, 1]$ controls how strongly we trust
434 the few shot execution relative to self prediction.

435 In practice, α_{t+1} is set adaptively based on the
436 discrepancy between S_{t+1}^{pred} and S_{t+1}^{few} and the few
437 shot sampling ratio. Let $\epsilon_{t+1} = |S_{t+1}^{\text{pred}} - S_{t+1}^{\text{few}}|$,
438 and let $\tau > 0$ be a calibration tolerance. We also
439 define the few shot ratio $\psi = \frac{|\mathcal{D}_{\text{few}}|}{|\mathcal{D}_{\text{val}}|}$, and an upper
440 bound $\alpha_{\max} \in (0, 1]$ on the calibration strength.
441 We set α_{t+1} as:

$$442 \quad \alpha_{t+1} = \begin{cases} 0, & \text{if } \epsilon_{t+1} \leq \tau, \\ \min\left(\frac{\epsilon_{t+1}}{\tau} \psi, \alpha_{\max}\right), & \text{otherwise.} \end{cases} \quad (9)$$

443 Intuitively, when S_{t+1}^{pred} and S_{t+1}^{few} agree within the
444 tolerance τ , we keep \hat{S}_{t+1} equal to the self prediction
445 ($\alpha_{t+1} = 0$). When the discrepancy exceeds
446 τ , we increase α_{t+1} proportionally to how many
447 tolerance units ϵ_{t+1} spans and to the few shot ratio

Method	DROP		HotpotQA		GSM8K		MATH		HumanEval		MBPP	
	Perf	Cost	Perf	Cost	Perf	Cost	Perf	Cost	Perf	Cost	Perf	Cost
ScoreFlow	91.48%	2.06e7	76.85%	2.72e7	97.79%	2.83e7	84.35%	4.03e7	98.91%	3.86e6	89.63%	3.93e6
AgentPrune	89.22%	6.65e6	76.73%	2.81e7	95.42%	1.07e7	81.53%	2.82e7	98.03%	1.65e6	88.37%	1.34e6
Aflow	90.30%	1.40e7	77.57%	3.11e7	96.58%	1.11e7	82.92%	3.06e7	99.24%	5.26e6	89.74%	3.92e6
SCALE	88.96%	3.27e6	77.41%	1.43e7	95.83%	4.04e6	81.70%	7.04e6	98.47%	1.09e6	90.32%	6.83e5
Δ	-1.34%	$\downarrow 76\%$	-0.16%	$\downarrow 54\%$	-0.75%	$\downarrow 63\%$	-1.22%	$\downarrow 77\%$	-0.77%	$\downarrow 79\%$	+0.58%	$\downarrow 83\%$
SCALE _{S^{pred}}	78.61%	4.45e6	75.40%	1.50e7	92.51%	7.69e6	76.33%	7.55e6	09.16%	2.30e5	90.62%	5.75e5
SCALE _{S^{few}}	86.06%	2.85e6	73.05%	1.55e7	94.22%	6.28e6	78.10%	5.49e6	96.95%	4.13e5	91.20%	6.17e5
SCALE _{S^{conf}}	00.00%	1.78e6	00.00%	1.04e7	00.00%	4.43e6	57.61%	8.48e6	97.71%	3.15e5	88.86%	1.24e6

Table 2: Test performance **Perf** and token number **Cost** comparison across six benchmarks. Δ reports the performance change and cost reduction of **SCALE** relative to Aflow.

ψ , but cap it by α_{\max} . This realizes the heuristic that large disagreements are more likely due to prediction error and should be corrected more aggressively toward the few shot estimate, while still respecting the limited size of \mathcal{D}_{few} .

To this end, we replace full validation execution with the calibrated surrogate score \hat{S}_i in Equation 8. The surrogate scores of low-cost stage $\{\hat{S}_i\}_{i>M}$ are used for selection, expansion, and experience update together with the full-execution scores in warm up stage $\{S_i^{\text{exec}}\}_{i=1}^M$.

Overall, **SCALE** eliminates full validation runs in the main search phase: only few shot execution is performed to calibrate the optimizer’s self prediction. As a result, token number scales with $|\mathcal{D}_{\text{few}}|$ instead of $|\mathcal{D}_{\text{val}}|$, cutting the token cost substantially while maintaining the performance.

5 Experiments

We empirically evaluate our framework on multiple benchmarks, aiming to answer these two questions: **Q1:** Can **SCALE** reduce cost while maintaining performance? **Q2:** Why does calibrated prediction approximate full execution score well?

5.1 Experimental Setup

Agent and Optimizer In all experiments, we adopt Qwen-Plus (Hui et al., 2024) as base models $\{O_{\theta_i}\}_{i=1}^n$ of executor agents and Qwen3-8B (Yang et al., 2025) as the workflow optimizer ϕ .

Benchmarks We evaluate on six benchmarks spanning diverse domains: DROP (Dua et al., 2019) and HotpotQA (Yang et al., 2018) (multi-hop reasoning), GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) (mathematical reasoning), HumanEval (Chen, 2021) and MBPP (Austin et al., 2021) (program synthesis). The dataset splits follow Zhang et al. (2024b).

Baselines We compare **SCALE** with both task-level methods Aflow (Zhang et al., 2024b), Agent-Prune (Zhang et al., 2024a) and query-level method ScoreFlow (Wang et al., 2025a). We also include internal ablations that vary only in the surrogate score: **SCALE**_{S^{pred}} uses uncalibrated self prediction S^{pred} ; **SCALE**_{S^{few}} uses few shot score S^{few} ; **SCALE**_{S^{conf}} uses self-confidence S^{conf} , defined as $S_{t+1}^{\text{conf}} = \phi(W_i, \hat{P}_{t+1}^{\text{optimizer}})$, where $\hat{P}_{t+1}^{\text{optimizer}}$ appends “Output your confidence on the answer” to the optimizer prompt $P_{t+1}^{\text{optimizer}}$. Unlike S^{pred} , S^{conf} isolates the effect of our dedicated evaluation prompt P^{eval} by contrasting it with a minimal modification of the generation prompt.

Metrics The main metrics reported are: test performance and overall LLM token number incurred excluding test-time execution. Each benchmark’s performance metric is the same as in (Zhang et al., 2024b). Importantly, the token number is computed differently across methods to reflect their distinct optimization paradigms: for *task-level* approaches, it includes all tokens consumed in evaluating candidate workflows on the validation set to select the single best task-level policy; for *query-level* method, it includes tokens used in generating data for training the optimizer ϕ , evaluating query-level workflows during validation, and generating the query-level workflow at test time.

5.2 Main Results and Ablation Study

Across all six benchmarks, the results in Table 2 show that our method substantially reduces token number while maintaining test performance. Compared with Aflow, our method **SCALE** yields an average test performance drop of only 0.61%, while reducing total token number by 54% to 83% across all benchmarks. This demonstrates that full validation execution is not necessary for discovering high-quality workflows. Relative to

522 AgentPrune (Zhang et al., 2024a), which lowers
 523 token cost through structural pruning but still re-
 524 lies on repeated execution-based evaluation, our
 525 method achieves further token number reduction
 526 while delivering comparable performance, indicat-
 527 ing that replacing the evaluation paradigm itself
 528 yields greater savings than modifying the search
 529 structure alone.

530 The ablation variants further clarify where these
 531 gains come from. **SCALE**– S^{pred} drastically re-
 532 duces cost but exhibits noticeable performance
 533 degradation, suggesting that self prediction alone
 534 suffers from model bias. **SCALE**– S^{few} improves
 535 robustness but still requires larger execution bud-
 536 getts. **SCALE**– S^{conf} performs very bad across
 537 tasks, confirming that naive confidence signals are
 538 unreliable surrogates for workflow quality. In con-
 539 trast, using calibrated prediction as surrogate eval-
 540 uation **SCALE** strikes a stable balance between
 541 cost and performance by combining model-based
 542 prediction with few shot execution signals.

543 Overall, these results answer **Q1** affirmatively:
 544 **SCALE** maintains the test performance while re-
 545 ducing searching token number by up to 83%.

546 5.3 Comparing Different Surrogate 547 Evaluation Methods

548 To answer **Q2**, For every workflow generated in
 549 a full Aflow’s run, we log and compare S^{exec} to-
 550 gether four surrogate scores. Surrogate scores are
 551 intended to take place of $S_{i>M}^{\text{exec}}$ to guide selection
 552 and expansion, so a good surrogate should satisfy
 553 two properties: First, the value should be close to
 554 S^{exec} . If the scales differ, the search will unfairly
 555 favor one side of the two-stages. Second, it should
 556 induce a ranking consistent with S^{exec} to reliably
 557 distinguish workflows.

558 Let $\{x_t\}_{t=1}^T$ denote the sequence of S^{exec} along
 559 the search progress, and $\{y_t\}_{t=1}^T$ the corresponding
 560 surrogate scores. We quantify agreement between
 561 x and y with:

562 **(1) Pearson correlation:** $\text{Pearson}(x, y) =$
 563 $\frac{\sum_t (x_t - \bar{x})(y_t - \bar{y})}{\sqrt{\sum_t (x_t - \bar{x})^2} \sqrt{\sum_t (y_t - \bar{y})^2}}$, which measures linear
 564 ranking consistency. The larger pearson correla-
 565 tion means the better linear ranking consistency
 566 between two metrics.

567 **(2) First-order difference cosine similarity:**
 568 $\text{DiffCos}(x, y) = \frac{\sum_{t=2}^T \Delta x_t \Delta y_t}{\sqrt{\sum_{t=2}^T \Delta x_t^2} \sqrt{\sum_{t=2}^T \Delta y_t^2}}$, where
 569 $\Delta x_t = x_t - x_{t-1}$ and $\Delta y_t = y_t - y_{t-1}$. It captures
 570 whether the direction of round-to-round changes is
 571 aligned between two workflow’s measure metrics.

Method	Pearson \uparrow	DiffCos \uparrow	MAE \downarrow
S^{conf}	-0.0576	0.0377	0.0802
S^{pred}	0.0517	0.1275	0.0511
S^{few}	0.6827	0.6192	0.2160
\hat{S}	0.5217	0.5545	0.1634

572 Table 3: Agreement between surrogate evalua-
 573 tion metrics and full execution. \hat{S} strikes a good balance between
 574 value accuracy and ranking consistency.

572 Similar to pearson correlation, the larger DiffCos
 573 means the two metrics are better aligned.

574 **(3) Mean absolute error (MAE):** $\text{MAE}(x, y) =$
 575 $\frac{1}{T} \sum_{t=1}^T |x_t - y_t|$ which directly measures value-
 576 level approximation quality of the surrogate evalua-
 577 tion methods.

578 Table 3 reports these metrics for different surro-
 579 gates. S^{conf} performs poorly on all metrics. S^{pred}
 580 achieves lowest MAE but almost zero correlation,
 581 indicating that it roughly matches the average scale
 582 of S^{exec} yet fails to order different workflows. S^{few}
 583 shows strong Pearson correlation and DiffCos but
 584 suffers from large MAE due to high variance from
 585 small sample size. \hat{S} combines the strengths of
 586 both: it improves correlation over self prediction
 587 while reducing MAE compared with few shot exe-
 588 cution.

589 Overall, the answer to **Q2** is that the calibrated
 590 prediction \hat{S} serves as the most effective surrogate
 591 for the full-execution score S^{exec} . As a *surrogate*
 592 *method*, it aligns best with S^{exec} in both value
 593 agreement and ranking consistency; as a *evalua-*
 594 *tion score* for task-level workflow generation, it
 595 maintains a competitive test performance while
 596 substantially reducing token cost.

597 6 Conclusion

598 We revisited agentic workflow generation and ar-
 599 rived at two main insights: (1) query-level work-
 600 flow generation is not always necessary, as a
 601 small pool of task-level workflows already achieves
 602 strong coverage, and (2) execution-based task-level
 603 evaluation is extremely costly while providing lim-
 604 ited benefit. Motivated by these findings, we de-
 605 veloped a task-level workflow generation frame-
 606 work **SCALE** to replace costly evaluation with
 607 calibrated self prediction. Across multiple bench-
 608 marks, it maintains competitive test performance
 609 with an average degradation of just 0.61% com-
 610 pared to existing approach while reducing overall
 611 token number by up to 83%.

612 7 Limitations

613 This work still has some limitations. Although our
614 method avoids the main drawbacks of both query-
615 level and task-level workflow generation methods,
616 it remains primarily based on task-level workflow
617 generation and has not yet fully leveraged the
618 generalization capability of task-level approaches
619 together with the fine-grained adaptability of query-
620 level methods. We also do not assess cross-domain
621 generalization in this work.

622 References

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1 others. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Mark Chen. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. *arXiv preprint arXiv:1903.00161*.
- Hongcheng Gao, Yue Liu, Yufei He, Longxu Dou, Chao Du, Zhijie Deng, Bryan Hooi, Min Lin, and Tianyu Pang. 2025. Flowreasoner: Reinforcing query-level meta-agents. *arXiv preprint arXiv:2504.15257*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Shengran Hu, Cong Lu, and Jeff Clune. 2024. Automated design of agentic systems. *arXiv preprint arXiv:2408.08435*.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, and 1 others. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. 2024. Encouraging divergent thinking in large language models through multi-agent debate. In *Proceedings of the 2024 conference on empirical methods in natural language processing*, pages 17889–17904.
- Boye Niu, Yiliao Song, Kai Lian, Yifan Shen, Yu Yao, Kun Zhang, and Tongliang Liu. 2025. Flow: Modularized agentic workflow automation. *arXiv preprint arXiv:2501.07834*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652.
- Yinjie Wang, Ling Yang, Guohao Li, Mengdi Wang, and Bryon Aragam. 2025a. Scoreflow: Mastering llm agent workflows via score-based preference optimization. *arXiv preprint arXiv:2502.04306*.
- Zhexuan Wang, Yutong Wang, Xuebo Liu, Liang Ding, Miao Zhang, Jie Liu, and Min Zhang. 2025b. Agentdropout: Dynamic agent elimination for token-efficient and high-performance llm-based multi-agent collaboration. *arXiv preprint arXiv:2503.18891*.
- Shengxiang Xu, Jiayi Zhang, Shimin Di, Yuyu Luo, Liang Yao, Hanmo Liu, Jia Zhu, Fan Liu, and Min-Ling Zhang. 2025. Robustflow: Towards robust agentic workflow generation. *arXiv preprint arXiv:2509.21834*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengan Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 2369–2380.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*.
- Rui Ye, Shuo Tang, Rui Ge, Yaxin Du, Zhenfei Yin, Siheng Chen, and Jing Shao. 2025. Mas-gpt: Training llms to build llm-based multi-agent systems. *arXiv preprint arXiv:2503.03686*.
- Guibin Zhang, Luyang Niu, Junfeng Fang, Kun Wang, Lei Bai, and Xiang Wang. 2025. Multi-agent architecture search via agentic supernet. *arXiv preprint arXiv:2502.04180*.

718 Guibin Zhang, Yanwei Yue, Zhixun Li, Sukwon Yun,
 719 Guancheng Wan, Kun Wang, Dawei Cheng, Jef-
 720 frey Xu Yu, and Tianlong Chen. 2024a. Cut the
 721 crap: An economical communication pipeline for
 722 llm-based multi-agent systems. *arXiv preprint*
 723 *arXiv:2410.02506*.

724 Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng,
 725 Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin
 726 Cheng, Sirui Hong, Jinlin Wang, and 1 others. 2024b.
 727 Aflow: Automating agentic workflow generation.
 728 *arXiv preprint arXiv:2410.10762*.

729 Mingchen Zhuge, Wenyi Wang, Louis Kirsch,
 730 Francesco Faccio, Dmitrii Khizbulin, and Jürgen
 731 Schmidhuber. 2024. Gptswarm: Language agents
 732 as optimizable graphs. In *Forty-first International*
 733 *Conference on Machine Learning*.

734 A Appendix

735 A.1 Self Prediction Prompt Template

736 The following prompt template is used to guide an
 737 LLM-based evaluator in performing *self-prediction*,
 738 i.e., estimating the expected accuracy of a candi-
 739 date workflow over the entire evaluation dataset
 740 before actual execution. The prompt is carefully
 741 structured to enforce rigorous static analysis while
 742 leveraging historical execution feedback for cali-
 743 bration. Key components include: (1) a clear task
 744 definition requiring both justification and a cali-
 745 brated probability score; (2) contextual informa-
 746 tion about the dataset, few-shot examples, and the
 747 workflow’s code structure; (3) explicit descriptions
 748 of allowed LLM-based operators to validate cor-
 749 rect usage; (4) reference experiences from prior
 750 rounds to enable prediction refinement through er-
 751 ror reflection; and (5) strict validation rules (e.g.,
 752 package imports, prompt definitions, operator inter-
 753 faces) that trigger immediate failure (score = 0.0)
 754 if violated. The output format enforces structured
 755 reasoning via `<reason>` and `<box>` tags to ensure
 756 parseable and consistent responses.

```

1 SELF_PREDICTION_PROMPT = """
2 You are an expert evaluator of workflows
3 .
4 Your task is to predict the probability
5     that a given workflow will correctly
6     execute on the WHOLE DATASET,
7     which represents your estimation of its
8     overall accuracy.
9 Respond with a brief explanation first,
10    followed by a single floating-point
11    number between 0.0 and 1.0.

12
13 Dataset Description:
14 <dataset>
15 {dataset_description}
16 </dataset>

```

13	Few-shot samples of the Dataset (jsonl format):	775
14	{dataset_few_shots}	776
15		777
16	Workflow to evaluate (python code):	778
17	<workflow>	779
18	{workflow}	780
19	</workflow>	781
20		782
21	Prompt used in the workflow (python code):	783
22	<prompt>	784
23	{prompt}	785
24	</prompt>	786
25		787
26	The workflow is Python code; the key function is <code>__call__(question)</code> , which produces the workflow’s response.	788
27	The workflow may call LLM-based operators described below:	789
28	<operator_description>	790
29	{operator_description}	791
30	</operator_description>	792
31		793
32	Reference experiences:	794
33	During the warm-up rounds, several workflows have been executed and evaluated.	795
34	Each record includes the round (the iteration number of the workflow), the score (the actual reward obtained after execution), the prediction (the reward you predicted in the previous round), and the python code of the workflow and prompt. (these workflows use the same operators as shown above in the <operator_description></ operator_description>)	796
35	These experiences are provided to help you calibrate your future predictions by comparing your past predicted rewards with the actual scores, you can adjust your estimation strategy to make your predicted rewards as close as possible to the real execution results.	797
36	<experience>	798
37	{experience}	799
38	</experience>	800
39		801
40	**General Instructions for evaluation:**	802
41	1. Step by step, carefully check for critical errors that could prevent execution:	803
42	- Package check (VERY IMPORTANT): The workflow code imports the required packages (for example: import numpy, asyncio and other commonly used Python packages). If any package is used in the workflow but missing or commented out, output 0.0 directly and do not continue other checks.	804
43	- Prompt check (VERY IMPORTANT): The workflow code uses prompts written in Python format.	805
44		806

```

844      45 If the workflow uses no prompts then
845          just continue other checks.
846      46 If the workflow uses prompts,
847          For every prompt referenced in the
848              workflow, you must verify that this
849                  prompt is properly defined in the
850                      prompt.py file (commonly imported as
851                          prompt_custom).
852      48 A prompt is considered properly
853          defined only if: It appears in the
854              prompt.py file without being
855                  commented out AND the prompt name
856                      matches exactly (including
857                          capitalization, underscores, and
858                              punctuation).
859      49 If any prompt used in the workflow is
860          missing, misspelled, or commented
861              out in prompt.py, you must
862                  immediately output 0.0.
863      50 Check ALL prompts used in the
864          workflow following the same rule.

865      51 - Operator check (VERY IMPORTANT):
866          The operator is provided in text
867              description format. If the workflow
868                  uses an operator, it must be among
869                      the operators defined in <
870                          operator_description>. If the
871                              workflow uses an undefined operator
872                                  (including mismatched names,
873                                      incorrect parameters, or improper
874                                          usage) OR the parameters passed when
875                                              using an operator do not comply
876                                                  with the interface requirements
877                                                      defined in operator_description,
878                                                          output 0.0 directly and do not
879                                                              continue other checks.

880      53 - Workflow check (VERY IMPORTANT):
881          The __call__ funciton must return
882              the output string of the workflow
883                  and the token usage only, more or
884                      less is totally wrong. The input of
885                          the workflow are only the input
886                              string, more or less is totally
887                                  wrong.

888      54 2. Step by step, Analyze whether the
889          workflow can logically solve the
890              queries in the WHOLE DATASET. Please
891                  carefully analyze the function of
892                      each operator and whether the
893                          position of each operator can
894                              smoothly promote the resolution of
895                                  the problem.

896      55 3. Consider potential hallucinations
897          from the operators, for now, we use
898              {backbone_model} as backbone model
899                  in operators.

900      56 4. Evaluate carefully and
901          comprehensively across all query
902              types in the WHOLE DATASET.

903      57 5. Be fair and rational: do not easily
904          assign 0.0 unless there is a severe
905              problem, and avoid scoring 1.0 with
906                  overconfidence.

```

```

64 6. There's no need to focus heavily on
914      output details, such as formatting
915          inconsistencies, since extraction is
916              handled simply or by specialized
917                  subsequent steps.

65 66 Output format:
918      - Provide a brief explanation of your
919          reasoning in a <reason> tag.
920      - Wrap your final probability in a <box>
921          tag **after** the <reason>.

66 67 For example:
922      <reason>The workflow correctly calls all
923          operators and uses only defined
924              prompts.</reason>
925      <box>0.85</box>
926      """
927
928
929
930
931

```

Listing 1: The prompt template for self prediciton.

A.2 Workflows Generated By SCALE

Here we show the best workflows generated by our method SCALE across six datasets.

```

1 1 from typing import Literal
2 2 import workspace_SCALE.DROP.workflows_43
3 3     .template.operator as operator
4 4 import workspace_SCALE.DROP.workflows_43
5 5     .round_6.prompt as prompt_custom
6 6 from scripts.async_llm import
7 7     create_llm_instance
8
9
10 from scripts.evaluator import
11     DatasetType
12
13 class Workflow:
14     def __init__(
15         self,
16         name: str,
17         llm_config,
18         dataset: DatasetType,
19     ) -> None:
20         self.name = name
21         self.dataset = dataset
22         self.llm = create_llm_instance(
23             llm_config)
24         self.answer_generate = operator.
25             AnswerGenerate(self.llm)
26         self.custom = operator.Custom(
27             self.llm)
28         self.sc_ensemble = operator.
29             ScEnsemble(self.llm)
30         self.verify_output = operator.
31             Custom(self.llm) # Reusing Custom
32                 as VerifyOutput
33
34     async def __call__(self, problem:
35         str):
36         """
37             Implementation of the workflow
38         """
39
40             # Step-by-step generation using
41                 AnswerGenerate
42                     initial_solution = await self.
43                         answer_generate(input=problem)

```

```

977      initial_thought =
978      initial_solution['thought']
979      initial_answer =
980      initial_solution['answer']
981
982      # Refine using custom method
983      # with clearer instruction and context
984      refined_solutions = []
985      for _ in range(5): # Increase
986          number of refinements for better
987          consensus
988          refined_sol = await self.
989          custom(
990              input=f"{problem}\n"
991              nInitial Thought: {initial_thought}\n"
992              nInitial Answer: {initial_answer}\",
993              instruction=
994              prompt_custom.
995              REFINE_WITH_CONTEXT_AND_ACCURATE_PERCENTAGE_CALCULATION_PROMPT
996
997          )
998          refined_solutions.append(
999          refined_sol['response'])
1000
1001      # Use self-consistency ensemble
1002      # to select the best solution
1003      ensemble_result = await self.
1004      sc_ensemble(solutions=
1005      refined_solutions)
1006      raw_final_response =
1007      ensemble_result['response']
1008
1009      # Verification step to ensure
1010      # final output format compliance
1011      verified_solution = await self.
1012      verify_output(
1013          input=raw_final_response,
1014          instruction=prompt_custom.
1015          VERIFY_OUTPUT_FORMAT_PROMPT
1016          )
1017
1018      return verified_solution[
1019          'response'], self.llm.
1020          get_usage_summary()['total_tokens']

```

Listing 2: The best workflow generated by SCALE for DROP

```

1021 VERIFY_OUTPUT_FORMAT_PROMPT = """Ensure
1022     the response is properly formatted
1023     with the answer() wrapper. If the
1024     answer is not wrapped in answer(),
1025     add it around the final result. For
1026     example:
1027     - If the answer is a number: answer(42)
1028     - If the answer is text: answer(Wilson)
1029     - If the answer is a vector: answer(\\
1030         begin{pmatrix} 1 \\\ 2 \\\end{\
1031         pmatrix})
1032     - If the answer is a range: answer(1-10)
1033     - If the answer is a percentage: answer
1034         (95%)
1035
1036     The response should contain only the
1037     final answer wrapped in answer()
1038     with no additional text or
1039     explanation."""
1040
1041 VERIFY_MATH_REASONING_PROMPT = """Check
1042     the mathematical reasoning in the

```

```

11     solution. Verify that:
12     1. All calculations are correct
13     2. The logic follows mathematical
14     principles
15     3. The steps lead to the correct
16     conclusion
17     4. The final answer is consistent with
18     the reasoning
19
20     If any errors are found, correct them
21     and provide the corrected solution.
22     Ensure the final answer is wrapped
23     in answer() with no additional text.
24 """

```

Listing 3: The prompt used in the executor agents of best workflow generated by SCALE for DROP

```

2500
2501 from typing import Literal
2502 import workspace_SCALE.HotpotQA.
2503     workflows.template.operator as
2504     operator
2505 import workspace_SCALE.HotpotQA.
2506     workflows.round_15.prompt as
2507     prompt_custom
2508 from scripts.async_llm import
2509     create_llm_instance
2510
2511
2512 from scripts.evaluator import
2513     DatasetType
2514
2515 class Workflow:
2516     def __init__(
2517         self,
2518         name: str,
2519         llm_config,
2520         dataset: DatasetType,
2521     ) -> None:
2522         self.name = name
2523         self.dataset = dataset
2524         self.llm = create_llm_instance(
2525             llm_config)
2526         self.custom = operator.Custom(
2527             self.llm)
2528         self.answer_generate = operator.
2529             AnswerGenerate(self.llm)
2530         self.sc_ensemble = operator.
2531             ScEnsemble(self.llm)
2532         self.refine = operator.Custom(
2533             self.llm) # New operator for post-
2534             ensemble refinement
2535
2536     async def __call__(self, problem:
2537         str):
2538         """
2539             Implementation of the workflow
2540         """
2541         # Generate multiple reasoning
2542         paths
2543         solutions = []
2544         for _ in range(3):
2545             solution = await self.
2546             answer_generate(input=problem)
2547             solutions.append(solution['
2548                 answer'])"""

```

```

1107      # Self-consistency ensemble to
1108      choose most frequent answer
1109      ensemble_response = await self.
1110      sc_ensemble(solutions=solutions)
1111      selected_answer =
1112      ensemble_response['response']

1113      # Verify ensemble result for
1114      validity and confidence before
1115      refining
1116      verify_response = await self.
1117      custom(
1118          input=f"Question: {problem}\nCandidate Answer: {selected_answer}"
1119          ,
1120          instruction=prompt_custom.
1121          VERIFY_ENSEMBLE_CONFIDENCE_PROMPT
1122          )
1123          is_valid = "answer(valid)" in
1124          verify_response['response']

1125          if not is_valid:
1126              fallback_response = await
1127              self.custom(
1128                  input=problem,
1129                  instruction=
1130                  prompt_custom.
1131                  FALLBACK_ANSWER_GENERATION_PROMPT
1132                  )
1133                  selected_answer =
1134                  fallback_response['response']

1135          # Refine the selected answer
1136          # with a stricter categorical/entity-
1137          # focused prompt
1138          refined_response = await self.
1139          refine(
1140              input=f"Question: {problem}\nSelected Answer: {selected_answer}"
1141              ,
1142              instruction=prompt_custom.
1143              REFINE_TO_ENTITY_OR_CATEGORY_PROMPT
1144              )
1145              refined_answer =
1146              refined_response['response']

1147          # Check if refined answer is
1148          # valid; if not, trigger fallback
1149          if "answer(None)" in
1150          refined_answer or not refined_answer
1151          .strip():
1152              fallback_response = await
1153              self.custom(
1154                  input=problem,
1155                  instruction=
1156                  prompt_custom.
1157                  FALLBACK_ANSWER_GENERATION_PROMPT
1158                  )
1159                  refined_answer =
1160                  fallback_response['response']

1161          # Final formatting verification
1162          verified_result = await self.
1163          custom(
1164              input=f"Question: {problem}\nCandidate Answer: {refined_answer}"
1165              ,
1166              instruction=prompt_custom.
1167              FINAL_FORMAT_VERIFICATION_PROMPT
1168              )

```

```

72
73      final_output = verified_result['
74      response']
75      return final_output, self.llm.
76      get_usage_summary()["total_tokens"]

```

Listing 4: The best workflow generated by SCALE for HotpotQA

```

1  VERIFY_ENSEMBLE_CONFIDENCE_PROMPT = """
2      You are given a question and a
3      candidate answer derived via
4      ensemble. Determine whether the
5      candidate answer is logically
6      consistent with the question and
7      shows sufficient confidence. If the
8      answer is relevant and confident,
9      respond with answer(valid).
10     Otherwise, respond with answer(
11         invalid). Do not explain or rephrase
12         just evaluate confidence and
13         relevance."""
14
15  REFINE_TO_ENTITY_OR_CATEGORY_PROMPT = """
16      You are given a question and a
17      candidate answer. Your task is to
18      reformulate the candidate answer
19      into a precise categorical label or
20      named entity that best fits the
21      question. Avoid explanatory or vague
22      language. Return only the most
23      accurate concise form, such as a
24      person's name, a location, a
25      historical event, or a categorical
26      label. Do not add punctuation or
27      quotation marks. Always wrap your
28      final output in answer(...).
29     Examples: answer(Polish independence
30     ), answer(William Shakespeare),
31     answer(no), answer(chronological
32     collection of critical quotations).
33 """
34
35  FINAL_FORMAT_VERIFICATION_PROMPT = """
36      You
37      are tasked with extracting and
38      formatting the final answer from a
39      candidate answer such that it
40      precisely matches the expected
41      format. Avoid including any
42      descriptive or explanatory text.
43      Focus on named entities, binary
44      responses, or categorical labels as
45      appropriate. For named entities (
46      people, places, works), return only
47      the name. For yes/no questions,
48      return exactly "yes" or "no". For
49      categorical responses, return the
50      exact category. Always wrap your
51      final response in answer(...).
52     Examples: answer(Limbo), answer(no),
53     answer(Southern Isles). If the
54     candidate answer contains multiple
55     possibilities, choose the most
56     likely one based on the question. If
57     the candidate is unclear, make a
58     best-guess effort to extract the
59     intended answer. Do not add quotes
60     or extra punctuation. Do not explain
61     your choice."""
62
63
64
65
66
67
68
69
70
71

```

```

1243
1244     7 FALBACK_ANSWER_GENERATION_PROMPT = """
1245         Given the original question, please
1246             generate a concise and direct answer
1247                 focusing strictly on the key entity
1248                     or fact requested. Avoid
1249                         explanations or additional
1250                             commentary. Always wrap your final
1251                                 output in answer(...). Example:
1252                                     Question: Who wrote Pride and
1253                                         Prejudice? Output: answer(Jane
1254                                         Austen)"""

```

Listing 5: The prompt used in the executor agents of best workflow generated by SCALE for HotpotQA

```

1255
1256     1 from typing import Literal
1257     2 import workspace_calibrated_prediction.
1258         GSM8K.workflows.template.operator as
1259             operator
1260     3 import workspace_calibrated_prediction.
1261         GSM8K.workflows.round_7.prompt as
1262             prompt_custom
1263     4 from scripts.async_llm import
1264         create_llm_instance
1265
1266     5
1267     6
1268     7 from scripts.evaluator import
1269         DatasetType
1270
1271     9 class Workflow:
1272     10     def __init__(self,
1273         name: str,
1274         llm_config,
1275         dataset: DatasetType,
1276     ) -> None:
1277         self.name = name
1278         self.dataset = dataset
1279         self.llm = create_llm_instance(
1280             llm_config)
1281         self.custom = operator.Custom(
1282             self.llm)
1283         self.programmer = operator.
1284             Programmer(self.llm)
1285         self.sc_ensemble = operator.
1286             ScEnsemble(self.llm)
1287
1288     23     async def __call__(self, problem:
1289         str):
1290         """
1291             Implementation of the workflow
1292             """
1293
1294     27         # Step 1: Generate multiple
1295             solutions using Programmer for
1296                 diverse computation paths
1297
1298         solutions = []
1299         for _ in range(3):
1300             solution = await self.
1301             programmer(problem=problem, analysis
1302                 ="Solve the following math problem
1303                     precisely. Return only the final
1304                         numeric result.")
1305             solutions.append(solution['
1306                 output'])
1307
1308         # Step 2: Use ScEnsemble to
1309             select the most consistent solution
1310                 among the generated ones

```

```

34             ensemble_result = await self.
1309             sc_ensemble(solutions=solutions,
1310                 problem=problem)
1311
1312             # Step 3: Format the selected
1313             result properly using Custom to
1314                 ensure it meets expected structure
1315             formatted_solution = await self.
1316             custom(
1317                 input=f"Problem: {problem}\\
1318                     nComputed Result: {ensemble_result['
1319                         response']]}",
1320                         instruction=prompt_custom.
1321             FORMAT_ANSWER_PROMPT
1322             )
1323
1324             # Step 4: Validate that the
1325             result makes sense in context (e.g.,
1326                 not negative where inappropriate,
1327                     correct order of magnitude)
1328             validated_solution = await self.
1329             custom(
1330                 input=f"Problem: {problem}\\
1331                     nFormatted Result: {{
1332                         formatted_solution['response']]},
1333                         instruction=prompt_custom.
1334             VALIDATE_NUMERIC_RESULT_PROMPT
1335             )
1336
1337             # Step 5: Final formatting
1338             verification to ensure answer is
1339                 boxed correctly
1340             final_result = await self.custom(
1341                 input=f"Problem: {problem}\\
1342                     nValidated Result: {{
1343                         validated_solution['response']]},
1344                         instruction=prompt_custom.
1345             FINAL_BOILING_CHECK_PROMPT
1346             )
1347
1348             return final_result['response'],
1349                     self.llm.get_usage_summary()["
1350                         total_tokens"]

```

Listing 6: The best workflow generated by SCALE for GSM8K

```

1 FORMAT_ANSWER_PROMPT = """You are given a
1352         math problem and its computed
1353             numeric result. Your task is to
1354                 format the result in a standardized
1355                     way by placing it inside \boxed{}.
1356             Only return the final formatted
1357                 answer without any additional text
1358                     or explanation. For example, if the
1359                         result is 123, return \boxed{123}.
1360             """
1361
1362
1363
1364 VALIDATE_NUMERIC_RESULT_PROMPT = """You
1365         are given a math word problem and a
1366             formatted numeric result. Check
1367                 whether the result is logically
1368                     reasonable in the context of the
1369                         problem (e.g., not negative when
1370                             expecting a count, correct magnitude
1371 ). If it seems incorrect, estimate a
1372                 plausible value and return it in
1373                     the same \boxed{} format. Otherwise

```

```

1374     , return the original result in \\
1375     boxed{} format. Only return the
1376     final result in \\boxed{}."""
1377
1378
1379 7 FINAL_BOILING_CHECK_PROMPT = """You are
1380     given a math problem and a validated
1381     numeric result. Ensure that the
1382     final result is enclosed in \\boxed
1383     {} and represents a clean numeric
1384     answer without any extra commentary
1385     or formatting issues. Return only
1386     the properly boxed result."""

```

Listing 7: The prompt used in the executor agents of best workflow generated by SCALE for GSM8K

```

1387
1388 1 from typing import Literal
1389 2 import workspace_SCALE.MATH.workflows.
1390     template.operator as operator
1391 3 import workspace_SCALE.MATH.workflows.
1392     round_20.prompt as prompt_custom
1393 4 from scripts.async_llm import
1394     create_llm_instance
1395
1396
1397 7 from scripts.evaluator import
1398     DatasetType
1399
1400 9 class Workflow:
1401 10     def __init__(_
1402 11         self,
1403 12         name: str,
1404 13         llm_config,
1405 14         dataset: DatasetType,
1406 15     ) -> None:
1407 16         self.name = name
1408 17         self.dataset = dataset
1409 18         self.llm = create_llm_instance(
1410 19             llm_config)
1411
1412 20         self.custom = operator.Custom(
1413 21             self.llm)
1414         self.verify_format = operator.
1415             Custom(self.llm)
1416         self.verify_math = operator.
1417             Custom(self.llm)
1418
1419 24     async def __call__(self, problem:
1420         str):
1421 25         """
1422             Implementation of the workflow
1423         """
1424 28         solution = await self.custom(
1425             input=problem, instruction="")
1426
1427 30         # Verify mathematical reasoning
1428 31         verified_solution = await self.
1429             verify_math(input=solution['response'],
1430             instruction=prompt_custom.
1431             VERIFY_MATH_REASONING_PROMPT)
1432
1433 33         # Verify and format the output
1434             to ensure it has answer() wrapper
1435 34         formatted_solution = await self.
1436             verify_format(input=
1437             verified_solution['response']),

```

```

35             instruction=prompt_custom.
1438             VERIFY_OUTPUT_FORMAT_PROMPT)
1439
1440 36             return formatted_solution['
1441             response'], self.llm.
1442             get_usage_summary()['total_tokens']
1443

```

Listing 8: The best workflow generated by SCALE for MATH

```

1 1 VERIFY_OUTPUT_FORMAT_PROMPT = """Ensure
1444     the response is properly formatted
1445     with the answer() wrapper. If the
1446     answer is not wrapped in answer(),
1447     add it around the final result. For
1448     example:
1449
2 - If the answer is a number: answer(42)
1450
3 - If the answer is text: answer(Wilson)
1451
4 - If the answer is a vector: answer(\\
1452     begin{pmatrix} 1 \\\ 2 \\\end{\
1453         pmatrix})
1454
5 - If the answer is a range: answer(1-10)
1455
6 - If the answer is a percentage: answer
1456     (95%)
1457
7
8 The response should contain only the
1458     final answer wrapped in answer()
1459     with no additional text or
1460     explanation."""
1461
9
10 10 VERIFY_MATH_REASONING_PROMPT = """Check
1462     the mathematical reasoning in the
1463     solution. Verify that:
1464
11 1. All calculations are correct
1465
12 2. The logic follows mathematical
1466     principles
1467
13 3. The steps lead to the correct
1468     conclusion
1469
14 4. The final answer is consistent with
1470     the reasoning
1471
15
16 If any errors are found, correct them
1472     and provide the corrected solution.
1473     Ensure the final answer is wrapped
1474     in answer() with no additional text.
1475
1476 """
1477
1478
1479

```

Listing 9: The prompt used in the executor agents of best workflow generated by SCALE for MATH

```

1480
1481 1 from typing import Literal
1482 2 import workspace_SCALE.HumanEval.
1483     workflows.template.operator as
1484         operator
1485 3 import workspace_SCALE.HumanEval.
1486     workflows.round_11.prompt as
1487         prompt_custom
1488 4 from scripts.async_llm import
1489     create_llm_instance
1490
1491
1492 7 from scripts.evaluator import
1493     DatasetType
1494
1495 9 class Workflow:
1496 10     def __init__(_
1497 11         self,
1498 12         name: str,
1499

```

```

1499     llm_config,
1500     dataset: DatasetType,
1501   ) -> None:
1502     self.name = name
1503     self.dataset = dataset
1504     self.llm = create_llm_instance(
1505       llm_config)
1506     self.custom = operator.Custom(
1507       self.llm)
1508     self.custom_code_generate =
1509       operator.CustomCodeGenerate(self.llm)
1510   )
1511   self.sc_ensemble = operator.
1512   ScEnsemble(self.llm)
1513   self.test = operator.Test(self.
1514   llm)

1515   async def __call__(self, problem:
1516     str, entry_point: str):
1517     """
1518       Implementation of the workflow
1519       Custom operator to generate
1520       anything you want.
1521       But when you want to get
1522       standard code, you should use
1523       custom_code_generate operator.
1524     """
1525
1526   # Rephrase the problem for
1527   # clarity
1528   rephrased_problem = await self.
1529   custom(input="", instruction=
1530   prompt_custom.
1531   REPHRASE_PROBLEM_PROMPT + problem)
1532   clarified_problem =
1533   rephrased_problem['response']

1534   # Generate multiple solutions
1535   for ensemble
1536     solutions = []
1537     tested_solutions = []
1538     for _ in range(5):
1539       solution = await self.
1540       custom_code_generate(problem=
1541       clarified_problem, entry_point=
1542       entry_point, instruction="")
1543
1544   # Reflect on the generated
1545   # solution to improve it
1546   reflection_prompt = f"Review
1547   the following code solution and fix
1548   any logical or syntax errors:\\n\\n
1549   {solution['response']}"
1550   reflected_solution = await
1551   self.custom(input=clarified_problem,
1552   instruction=reflection_prompt)
1553
1554   solutions.append(
1555     reflected_solution['response'])

1556   # Pre-test each refined
1557   # solution to filter valid ones early
1558   test_result = await self.
1559   test(problem=problem, solution=
1560     reflected_solution['response'],
1561     entry_point=entry_point)
1562
1563   if test_result['result']:
1564     tested_solutions.append(
1565       test_result['solution'])

1566   # Prioritize validated solutions
1567
1568

```

```

52   ; fallback to all if none pass
53   if tested_solutions:
54     ensemble_input =
55       tested_solutions
56   else:
57     ensemble_input = solutions

58   # Use ScEnsemble to select the
59   # most consistent solution
60   ensemble_result = await self.
61   sc_ensemble(solutions=ensemble_input
62   , problem=problem)
63   final_solution = ensemble_result
64   ['response']

65   return final_solution, self.llm.
66   get_usage_summary()["total_tokens"]

```

Listing 10: The best workflow generated by SCALE for HumanEval

```

1 REPHRASE_PROBLEM_PROMPT = """Please
2   rephrase the following programming
3   problem in clearer terms, making
4   sure to highlight the key
5   requirements and expected output
6   format. Problem: """

```

Listing 11: The prompt used in the executor agents of best workflow generated by SCALE for HumanEval

```

1 from typing import Literal
2 import workspace SCALE.MBPP.workflows.
3   template.operator as operator
4 import workspace SCALE.MBPP.workflows.
5   round_7.prompt as prompt_custom
6 from scripts.async_llm import
7   create_llm_instance

8
9 class Workflow:
10   def __init__(
11     self,
12     name: str,
13     llm_config,
14     dataset: DatasetType,
15   ) -> None:
16     self.name = name
17     self.dataset = dataset
18     self.llm = create_llm_instance(
19       llm_config)

20     self.custom = operator.Custom(
21       self.llm)
22     self.custom_code_generate =
23       operator.CustomCodeGenerate(self.llm
24     )
25     self.test = operator.Test(self.
26       llm)

27   async def __call__(self, problem:
28     str, entry_point: str):
29     """
30       Implementation of the workflow

```

```

1629      27     Custom operator to generate
1630      28     anything you want.
1631      29     But when you want to get
1632      30     standard code, you should use
1633      31     custom_code_generate operator.
1634      32     """
1635      33     # Generate the initial solution
1636      34     solution = await self.
1637      35     custom_code_generate(problem=problem
1638      36     , entry_point=entry_point,
1639      37     instruction="Generate Python code
1640      38     that solves the given problem. Make
1641      39     sure to return the result of the
1642      40     function, not just print it.")
1643      41
1644      42     # Verify that the solution has
1645      43     proper format and return statements
1646      44     verified_solution = await self.
1647      45     custom(input=f"Problem: {problem}\n
1648      46     nEntry point: {entry_point}\n
1649      47     nSolution:\n{solution['response']}",
1650      48     instruction=prompt_custom.
1651      49     VERIFY_CODE_FORMAT_PROMPT)
1652      50
1653      51     # Test the verified solution to
1654      52     ensure it works correctly
1655      53     test_result = await self.test(
1656      54     problem=problem, solution=
1657      55     verified_solution['response'],
1658      56     entry_point=entry_point)
1659      57
1660      58     final_solution = test_result['
1661      59     solution'] if test_result['result']
1662      60     else verified_solution['response']
1663      61
1664      62     return final_solution, self.llm.
1665      63     get_usage_summary()["total_tokens"]

```

Listing 12: The best workflow generated by SCALE for MBPP

```

1666 1 VERIFY_CODE_FORMAT_PROMPT = """Verify
1667 2     that the provided code solution has
1668 3     the correct function signature and
1669 4     includes proper return statements.
1670 5     The solution must:
1671 6     1. Contain the function with the exact
1672 7         name specified in the entry_point
1673 8     2. Include a return statement that
1674 9         returns the result of the function (
1675 10        not just print it)
1676 11     3. Follow proper Python syntax
1677 12
1678 13     If the code is missing the function
1679 14         signature or return statement,
1680 15         please fix it. Return the corrected
1681 16         code."""

```

Listing 13: The prompt used in the executor agents of best workflow generated by SCALE for MBPP