

METAFLOW: A META APPROACH OF TRAINING LLMs INTO GENERALIZABLE WORKFLOW GENERATORS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) excel across a wide range of tasks, yet their instance-specific solutions often lack the structural consistency needed for reliable deployment. Workflows that encode recurring algorithmic patterns at the task level provide a principled framework, offering robustness across instance variations, interpretable traces for debugging, and reusability across problem instances. However, manually designing such workflows requires significant expertise and effort, limiting their broader application. While automatic workflow generation could address this bottleneck, existing methods either produce instance-specific solutions without learning task-level patterns, or cannot generalize beyond their training configurations. We present **MetaFlow**, which casts workflow generation as a meta-learning problem: given a task and an operator set, the model learns to compose solution strategies. MetaFlow trains in two stages—supervised fine-tuning on synthetic workflow data, followed by reinforcement learning with verifiable rewards (RLVR) that uses execution feedback across problem instances in the task to improve end-to-end success. The resulting model produces effective workflows for trained tasks and exhibits strong generalization to untrained tasks and novel operator sets. Across benchmarks in question answering, code generation, and mathematical reasoning, MetaFlow achieves performance comparable to state-of-the-art baselines on in-domain tasks with single inference, while demonstrating remarkable zero-shot generalization capabilities on out-of-domain tasks and operator sets.

1 INTRODUCTION

Large Language Models (LLMs) have demonstrated significant performance across a wide range of tasks, including code generation, question answering, and mathematical reasoning (Austin et al., 2021; Chen et al., 2021; Yang et al., 2018; Dua et al., 2019; Ding et al., 2024; Jiang et al., 2025; Cobbe et al., 2021; OpenAI, 2023; Zhu et al., 2024). However, because these models generate instance-specific solutions, they lack the structural consistency and transparency needed for reliable deployment, while also being difficult to adapt to similar tasks. *Workflows* that encode recurring algorithmic patterns provide a principled alternative, decomposing complex challenges into structured, manageable steps. However, manually designing such workflows requires significant expertise and effort, limiting their broader application.

To address this challenge, recent effort have focused on the automatic workflow generation (Khatab et al., 2023; Li et al., 2024; Song et al., 2024; Zhang et al., 2024a). Endowing LLMs with this strategy planning capability means lowering the barrier for complex task automation from requiring manual programming by experts to merely providing high-level task descriptions, thereby greatly liberating productivity. Nevertheless, representing the workflow as static graph (Zhuge et al., 2024) or neural network (Liu et al., 2024) in many of these methods limits the flexibility of generatable workflows.

A promising direction emerges from works like ADAS (Hu et al., 2024), AFlow (Zhang et al., 2024b), ScoreFlow (Wang et al., 2025) and FlowReasoner (Gao et al., 2025), which represent workflows as code (a structured combination of predefined *operators*), making the automatic generation of workflows more flexible and expressive, where *operator* is an encapsulation of common agentic

operations introduced by Zhang et al. (2024b). Within this code-based framework, current approaches adopt two different paradigms for workflow generation.

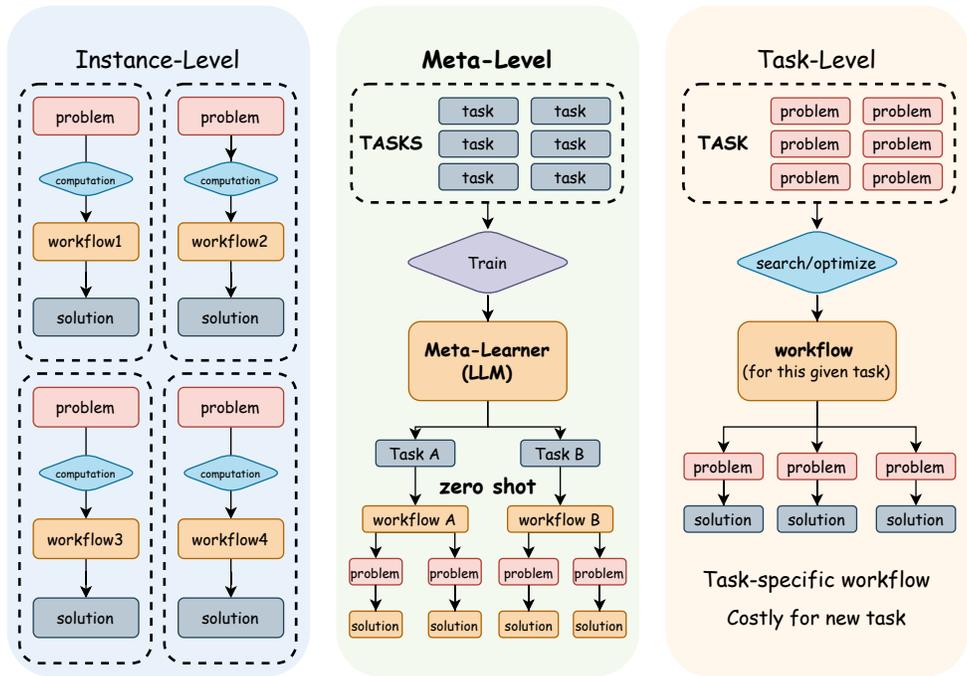


Figure 1: Illustration of **instance-level**, **task-level**, and **meta-level (ours)** workflow generation approaches. Unlike **instance-level** methods that generate workflows for individual problems or **task-level** methods that require costly search for each new task, our meta-learning approach enables zero-shot workflow generation across tasks.

The first paradigm comprises **task-level** approaches, exemplified by ADAS (Hu et al., 2024) and AFlow (Zhang et al., 2024b), which formulate workflow generation as a search problem within pre-defined task-operator set combination. Both methods employ iterative search strategies, with ADAS using evolutionary algorithms and AFlow using Monte Carlo Tree Search (MCTS), to discover high-performing workflows through repeated refinement. However, this search-based paradigm inherently constrains them to specific task and predetermined operator set. When encountering new tasks or operators, these methods require complete re-optimization from scratch, incurring substantial computational costs (Wang et al., 2025).

Conversely, the second paradigm consists of **instance-level** approaches, exemplified by ScoreFlow (Wang et al., 2025) and FlowReasoner (Gao et al., 2025), which generate workflows tailored to individual problem instances in the task. Both methods dynamically construct workflows at inference time, with ScoreFlow leveraging gradient-based optimization to refine agentic workflows, and FlowReasoner employing reasoning chains distilled from advanced models to design query-specific multi-agent systems. While these instance-level methods excel at tailoring workflows to specific problem instances, this granularity comes at the cost of reusability and deployment efficiency. They cannot capture task-level patterns that recur across similar problem instances, leading to redundant workflow generation for each query. In deployment scenarios, this approach foregoes the benefits of having optimized, reusable workflow templates that could be consistently applied to entire task.

The limitations of existing paradigms underscore two fundamental challenges that must be addressed to achieve truly general-purpose automatic workflow generation. **(1) How can we overcome the re-optimization requirement of task-level approaches when facing new domains (task-operator set combinations)? (2) How can we learn generalizable patterns that avoid redundant instance-level generation while adapting effectively to untrained domains?**

To systematically overcome the challenges, we propose **MetaFlow**, which formulates workflow generation as a meta-learning problem. As illustrated in Figure 1, unlike task-level search-based meth-

ods that require expensive re-optimization for new task-operator set combinations, and instance-level methods that generate workflows for individual queries, **MetaFlow** learns to directly synthesize workflows from task descriptions and operator set specifications, enabling zero-shot generation through a single model inference.

To achieve robust zero-shot generalization, **MetaFlow** employs a two-stage training paradigm with diverse task-operator pairs. Adopting the code-based workflow representation from prior works (Hu et al., 2024; Zhang et al., 2024b; Wang et al., 2025; Gao et al., 2025), we first synthesize thousands of workflows using Qwen-Max (Team, 2024) across four tasks and a single operator set to finetune Qwen3-8B (Yang et al., 2025), establishing the foundation for understanding how tasks and operators relate to workflow structures. Subsequently, we apply online reinforcement learning with GRPO (Shao et al., 2024) across expanded domains, where execution feedback on problem instances directly optimizes the generation policy. This training ensures the model learns generalizable workflow construction principles rather than memorizing patterns. At inference, **MetaFlow** zero-shot generates effective workflows for novel configurations with only a single forward pass.

Our main contributions are:

- **Meta-learning Framework:** We introduce **MetaFlow**, a novel approach that reformulates workflow generation from discrete search within fixed configurations to continuous learning across diverse task-operator set combinations. By conditioning workflow generation on task descriptions and operator specifications, our framework achieves strong zero-shot generalization to unseen domains without any re-optimization, reducing computational cost from thousands of API calls to a single model inference.
- **Scalable Training Pipeline:** We design a two-stage training framework combining supervised learning with online reinforcement learning, utilizing diverse domains to ensure robust generalization to untrained domains.
- **Comprehensive Evaluation:** Extensive experiments demonstrate that **MetaFlow** achieves competitive performance on in-domain benchmarks while exhibiting remarkable zero-shot generalization to out-of-domain task classes and operator sets, including solving programming problems with novel operator combinations never seen during training and solving question answering problem using the vector database search operator.

2 RELATED WORKS

2.1 AGENTIC WORKFLOW

Agentic workflows decompose complex tasks into structured steps through predefined operators and dependencies (Zhang et al., 2024b; Wang et al., 2025; Gao et al., 2025). Unlike autonomous agents that learn through environment interaction (Zhuge et al., 2024; Hong et al., 2024), workflows provide interpretable and consistent execution patterns. Recent works adopt code-based representations for superior expressiveness (Hu et al., 2024; Zhang et al., 2024b; Wang et al., 2025; Gao et al., 2025), supporting applications in code generation, question answering, and mathematical reasoning (Austin et al., 2021; Chen et al., 2021; Yang et al., 2018; Dua et al., 2019; Ding et al., 2024; Jiang et al., 2025; Cobbe et al., 2021; OpenAI, 2023; Zhu et al., 2024). However, manual workflow design remains a significant bottleneck requiring deep expertise.

2.2 AUTOMATIC WORKFLOW GENERATION

Recent advances have explored automating workflow generation for improving LLM performance (Chen et al., 2023; Zhang et al., 2024b; Wang et al., 2025; Li et al., 2024; Song et al., 2024). While some methods optimize prompts within fixed workflows (Guo et al., 2023; Khattab et al., 2023), we focus on optimizing workflow structures directly.

Current structural optimization follows two paradigms. Task-level approaches like ADAS (Hu et al., 2024) and AFlow (Zhang et al., 2024b) search for optimal workflows through evolutionary algorithms or MCTS, but require complete re-optimization for new domains. Instance-level methods including ScoreFlow (Wang et al., 2025) and FlowReasoner (Gao et al., 2025) generate query-specific workflows but fail to extract reusable patterns.

Our Methods, **MetaFlow** reformulates workflow generation as meta-learning over diverse task-operator combinations during training. Through two-stage optimization combining supervised learning with reinforcement learning, it achieves true zero-shot generation—producing effective workflows for novel domains via single model inference, eliminating both re-optimization and adaptation overhead.

3 PROBLEM DEFINITION

Existing works often formulate automatic workflow generation as optimization problems (Xu et al., 2025; Li et al., 2025), requiring separate optimizations for each task. We elevate this perspective by reformulating it as a meta-learning problem (Finn et al., 2017; Franceschi et al., 2018). To ground this formulation, we first define our core concepts:

- **PROBLEM INSTANCE** p : A single, concrete problem to be solved.
- **TASK** C : A family of problem instances sharing a common structure and solution strategy (e.g., GSM8K mathematical reasoning, DROP reading comprehension).
- **OPERATOR SET** Ops : A collection of fundamental, reusable operations (Zhang et al., 2024b) (e.g., Generate, Revise, Ensemble).
- **DOMAIN** (C, Ops) : The combination of a task and an operator set, defining a complete problem-solving context.

Within this meta-learning framework, an LLM serves as the meta-learner (the planner π_θ). Its core responsibility is to learn a meta-strategy that enables fast adaptation: given any domain (C, Ops) , the planner rapidly generates an efficient and reusable **WORKFLOW** W —a structured sequence of operators from Ops . When executed on any **PROBLEM INSTANCE** p from task C , this workflow produces a high-quality **SOLUTION** s .

Unlike traditional meta-learning approaches such as Model-Agnostic Meta-Learning (MAML) (Finn et al., 2017), which rely on gradient updates for adaptation, our method achieves fast adaptation through the synthesis of workflows without requiring any gradient-based fine-tuning.

As a meta-learning task, our goal is to optimize the meta-parameters θ of the planner π_θ such that the generated workflows maximize expected rewards across the domain distribution $\mathcal{D}_{(C, Ops)}$. We adopt the classic **bi-level optimization** framework Franceschi et al. (2018):

$$\theta^* = \arg \max_{\theta} \underbrace{\mathbb{E}_{(C, Ops) \sim \mathcal{D}_{(C, Ops)}}}_{\text{Outer Loop: Meta-Optimization}} \left[\underbrace{\mathbb{E}_{W \sim \pi_\theta(\cdot | C, Ops)} [\mathbb{E}_{p \sim C} [R(\text{Exec}(W, p))]]}_{\text{Inner Loop: Fast Adaptation \& Evaluation}} \right]$$

where:

- $(C, Ops) \sim \mathcal{D}_{(C, Ops)}$: Sample a domain (task-operator pair) from the distribution.
- $W \sim \pi_\theta(\cdot | C, Ops)$: The planner performs *fast adaptation*—given task description C and operator specifications Ops , it generates a customized workflow W without gradient updates.
- $p \sim C$: Sample problem instances from task C to evaluate the workflow.
- $\text{Exec}(W, p)$: Execute workflow W on problem p to produce a solution.
- $R(\cdot)$: Reward function evaluating solution quality (e.g., test pass rate, answer correctness).

The **outer loop** optimizes performance across the entire domain distribution $\mathcal{D}_{(C, Ops)}$ rather than on individual tasks or instances, driving the meta-learner π_θ to acquire cross-domain generalization capabilities. This contrasts sharply with prior works:

- **Instance-level methods** (e.g., ScoreFlow (Wang et al., 2025), ComfyUI-R1 (Xu et al., 2025)) optimize workflows for individual problem instances: $\arg \max_{W_p} \mathbb{E}_{p \sim C} [R(\text{Exec}(W_p, p))]$.
- **Task-level methods** (e.g., AFlow (Zhang et al., 2024b)) optimize a single workflow per task but require separate optimization for each new task-operator combination.

Our framework elevates to the **meta-level**, learning a planner that generalizes across tasks and operators: $W \sim \pi_\theta(\cdot | C, Ops)$.

4 METHODOLOGY

Our **MetaFlow** framework aims to train a large language model through a two-stage learning process to automatically generate efficient and reusable workflows for given domains (C, Ops), where C represents a task and Ops denotes the available operator set. The framework addresses a key challenge in automated workflow generation: how to enable a model to quickly adapt to new domains (task-operator combinations) without requiring extensive training data for each. This section presents our system architecture (Section 4.1) and the two-stage training algorithm combining supervised fine-tuning and reinforcement learning with verifiable rewards (Section 4.2).

4.1 METAFLOW ARCHITECTURE

The core architecture of **MetaFlow** consists of two components: (1) a **planner LLM** π_θ that generates workflows conditioned on domain specifications (C, Ops), and (2) an **execution-evaluation environment** that orchestrates workflow execution using the MetaGPT framework (Hong et al., 2023), invokes operators on sampled problem instances, and computes verifiable reward scores based on correctness evaluation.

Workflow Representation: A workflow is represented as a structured script based on the MetaGPT framework (Hong et al., 2023), composed of a series of predefined operator calls. As illustrated in the left part of Figure 2, we adopt the foundational text-processing operators from AFlow (Zhang et al., 2024b), including Generate, Summarize, Revise, Ensemble, and Programmer for basic tool invocation. Building upon the same MetaGPT interface, we further introduce additional operators for advanced text processing (e.g., Decompose 1, SelfConsistency) and domain-specific tools (e.g., VectorSearch 3). This standardized operator interface serves two purposes: (1) supervised fine-tuning provides a cold start that ensures the model generates syntactically valid operator calls with higher probability, and (2) the introduction of new operators demonstrates the necessity of our meta-learning approach for efficient and transferable workflow generation across diverse operator sets. Implementation details of the operators are provided in Appendix A.1 for reference.

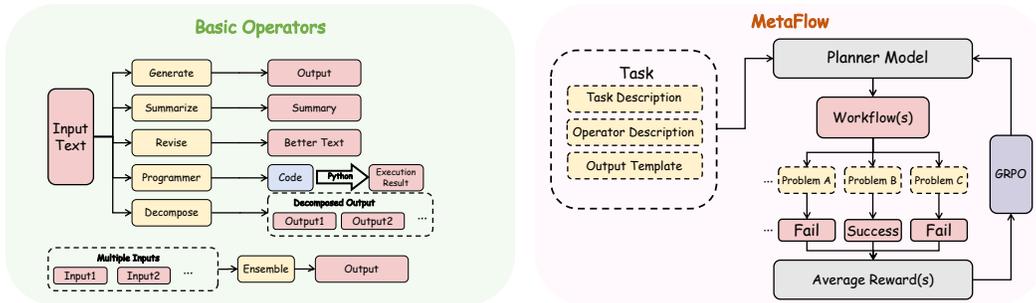


Figure 2: Overview of MetaFlow architecture. **Left:** Basic operators and their functionalities, including text-processing operators (Generate, Summarize, Revise, Ensemble) and basic tool-invocation operators (Programmer). **Right:** The MetaFlow training framework showing how the planner model generates workflows conditioned on domain (C, Ops), which are then evaluated on problem instances to compute rewards for GRPO optimization.

Input-Output: As illustrated in the right part of Figure 2, the MODEL INPUT consists of two parts: TASK DESCRIPTION 4, which elucidates the characteristics of the target task C and the input-output formats of problem instances; OPERATOR DESCRIPTIONS 5, which define the functions, parameters, and input-output formats of each available operator in Ops (whether pre-defined or user-defined). The OUTPUT 6 is a structured WORKFLOW aimed at efficiently solving all problem instances from task C using operator set Ops. This design paradigm allows users to introduce new operators and tasks through natural language descriptions at inference time, enabling the model to rapidly adapt without additional training. Implementation details are provided in Appendix A.2.

Execution & Evaluation Environment: To achieve end-to-end optimization, we build an automated environment. This environment receives a candidate workflow W_i and a set of N problem instances p_1, \dots, p_N sampled from task C as input. In the execution phase, the environment uses the MetaGPT framework to orchestrate the workflow (Hong et al., 2023). Each operator in the workflow (such as Generate) completes its specific subtask by calling an external lightweight language model API (e.g., Qwen-Turbo or GPT-4o-mini-0718) and processes each problem instance. Upon completion, an automated evaluator module verifies the correctness of each execution result, for example, by running unit tests or comparing outputs with ground truth answers. Based on the evaluation results, the environment computes a verifiable reward score $R(W_i)$ for workflow W_i , calculated as the average success rate over N instances:

$$R(W_i) = \frac{1}{N} \sum_{j=1}^N \mathbb{I}(\text{is correct}(\text{Exec}(W_i, p_j))) \in [0, 1] \quad (1)$$

where $\mathbb{I}(\cdot)$ is the indicator function. This reward score $R(W_i)$ is then passed to the training algorithm as the basis for policy updates.

4.2 TRAINING ALGORITHM

The training process consists of two stages: **supervised fine-tuning (SFT)** and **Reinforcement Learning with Verifiable Reward (RLVR)**. This design leverages supervised data to establish syntactic correctness, then employs reinforcement learning to optimize workflow effectiveness, separating structural learning from performance optimization.

4.2.1 PHASE ONE: SFT INITIALIZATION

To address the cold start problem in reinforcement learning, we first perform **supervised fine-tuning (SFT)** on the base LLM using a dataset of $([C, \text{Ops}], W)$ pairs, which is the pairs of domain $[C, \text{Ops}]$ and the corresponding workflow W . This phase teaches the model to generate syntactically correct workflows following the required template structure (Listing 6) and proper operator invocation patterns (Listing 5). By ensuring the model π_θ can produce well-formed code, the subsequent RL phase can focus solely on optimizing workflow effectiveness rather than learning basic syntax, significantly narrowing the search space and accelerating convergence. Details of the SFT dataset construction, including the four-stage pipeline and LoRA configuration, are provided in Appendix B.

4.2.2 PHASE TWO: RLVR OPTIMIZATION

After SFT, we employ policy gradient algorithms to perform end-to-end self-improvement on the planner π_θ . The core of this phase is the **Reinforcement Learning with Verifiable Reward (RLVR)** loop, with the complete algorithm presented in Algorithm 1. The specific process consists of the following steps:

1. **Policy Sampling:** For a task C sampled from the training set, the planner π_θ generates a batch of k candidate workflows W_1, W_2, \dots, W_k .
2. **Execution & Reward Calculation:** Each candidate workflow W_i is tested in the execution and evaluation environment described in Section 4.1, obtaining its corresponding reward score $R(W_i)$ (Equation 1) that reflects generalization capability:

$$C \in \mathcal{D}_{\text{train}} \xrightarrow{\pi_\theta} \{W_i | 1 \leq i \leq k\} \xrightarrow{\text{each } W_i} \{[W_i, R(W_i)], \text{ executed on } p_1^{(i)}, \dots, p_N^{(i)} | 1 \leq i \leq k\}$$

3. **Policy Update:** We use this batch of $[W_i, R(W_i)]$ pairs to update the parameters of the planner π_θ using the **Group Relative Policy Optimization (GRPO)** algorithm (Shao et al., 2024). The core idea of GRPO is to use the average performance within the group as a baseline to estimate advantages, thereby avoiding training an independent value network.

$$\text{advantage} = \hat{A}(W_i) = R(W_i) - \mu_R, \quad \mu_R = \sum_{j=1}^k R(W_j) / k$$

This group-relative advantages makes the optimization signal derive from whether the workflow performance is better or worse than the current batch’s average level, rather than an absolute, potentially noisy value estimate.

4. Variance Reduction: Evaluation Based on Common Random Numbers: The effectiveness of policy gradients largely depends on the accuracy of the advantage function $\hat{A}(W_i)$ estimation. In our GRPO method, the advantage is computed relative to the batch average reward μ_R . If each workflow W_i in the batch is evaluated on a set of independently and randomly sampled problem instances, the variance of the reward $R(W_i)$ will include not only policy randomness but also environmental randomness. This additional variance propagates to μ_R and $\hat{A}(W_i)$, producing noisier gradients and reducing learning efficiency.

To address this issue, we adopt a classic variance reduction technique: **Common Random Numbers (CRN)** (Kleijnen, 1975). In specific implementation, we ensure that all k candidate workflows W_1, \dots, W_k in the same training batch are evaluated on the exact same set of problem instances $\{p_1, \dots, p_N\}$. By fixing the random variable of problem instances when comparing workflows, we eliminate noise arising from differences in problem sampling. This keeps the expected value of $R(W_i) - R(W_j)$ unchanged but significantly reduces its variance. Ultimately, this ensures that our computed relative advantages more accurately reflect the intrinsic performance differences between workflows, leading to a more stable policy update direction and accelerating model convergence.

Algorithm 1 GRPO Optimization with CRN

Require: Training tasks $\mathcal{D}_{\text{train}}$, initial policy π_θ , group size k , batch size B , iterations T

```

1: for  $t = 1, \dots, T$  do
2:   Batch Sampling: Sample tasks  $\{C_1, \dots, C_B\} \sim \mathcal{D}_{\text{train}}$ 
3:   for each  $C_b$  in batch do
4:     Policy Sampling: Generate  $k$  workflows  $\mathcal{W}_b = \{W_{b,1}, \dots, W_{b,k}\}$ 
5:       where  $W_{b,i} \sim \pi_\theta(\cdot | C_b)$ 
6:     Common Random Numbers: Fix problem set  $\mathcal{P}_b = \{p_1, \dots, p_N\} \sim C_b$ 
7:     Reward Computation: For all  $i \in [1, k]$ :
8:        $R(W_{b,i}) = \frac{1}{N} \sum_{j=1}^N \mathbb{I}[W_{b,i} \text{ solves } p_j]$ 
9:     Group-Relative Advantage: Compute baseline  $\mu_{R,b} = \frac{1}{k} \sum_{i=1}^k R(W_{b,i})$ 
10:    For all  $i$ :  $\hat{A}(W_{b,i}) = R(W_{b,i}) - \mu_{R,b}$ 
11:  end for
12:  Batch GRPO Update:
13:  Collect all  $\{(W_{b,i}, \hat{A}(W_{b,i}))\}_{b=1, \dots, B; i=1, \dots, k}$ 
14:  Update:  $\theta \leftarrow \text{GRPO}(\theta, \{(W_{b,i}, \hat{A}(W_{b,i}))\})$ 
15: end for
16: return Optimized policy  $\pi_\theta$ 

```

5 EXPERIMENTS

This section evaluates **MetaFlow** through systematic experiments addressing three core questions: (1) Can **MetaFlow** achieve competitive performance with single-inference generation, eliminating the computational overhead of per-instance optimization? (2) Does the framework generalize to out-of-distribution tasks and novel operators unseen during training? (3) How do design choices—task-level versus instance-level generation, simple versus complex workflows—affect the performance-cost trade-off? We present training configuration, operator integration experiments, zero-shot generalization results, and main performance analysis.

5.1 TRAINING CONFIGURATION

Stage 1: Supervised Fine-Tuning (SFT) Initialization. We first conduct supervised fine-tuning (SFT) on the base model Qwen3-8B to teach the model π_θ^{Base} to generate syntactically correct workflows (Listing 6) and narrow the search space. We construct an expert dataset containing approximately 1,300 high-quality ($[C, \text{Ops}], W$) pairs, covering the tasks $\mathcal{D}_{\text{train}}^{\text{SFT}} = \{\text{GSM8K}, \text{DROP}, \text{MBPP}, \text{Humaneval}\}$ with the corresponding operator set for each task $\text{Ops}^{\text{SFT}} =$

(Generate, Summarize, Revise, Ensemble). The dataset is synthesized using Qwen-Max API through a four-stage pipeline (Appendix B). As shown in Figure 3 (left and middle), we train for one epoch with batch size 16 using LoRA (rank-16) (Hu et al., 2022) to prevent mode collapse on this limited dataset (Appendix B.2).

Stage 2: Reinforcement Learning with Verifiable Rewards (RLVR) Optimization. Following SFT initialization, we employ RLVR for end-to-end optimization of the planner $\pi_{\theta}^{\text{SFT}}$ with the same task set $\mathcal{D}_{\text{train}}^{\text{RLVR}} = \mathcal{D}_{\text{train}}^{\text{SFT}} = \{\text{GSM8K, DROP, MBPP, Humaneval}\}$. Critically, to enhance generalization to diverse operator configurations, we train with four different operator sets $\{\text{Ops}^{\text{RLVR}, i}\}_{i=1}^4$ that progressively introduce novel operators (Programmer, Decompose) beyond the base SFT set, randomly sampling domain combinations (C, Ops) at each iteration. The training employs Algorithm 1 for 137 steps. To control training costs, the execution of generated workflow W_i on problem instance p_j calls the Qwen-Turbo API. As shown in Figure 3 (right), the average reward rapidly increases from 0.67 to approximately 0.88 within the first 35 training steps, then stabilizes around 0.85–0.90, demonstrating the effectiveness of GRPO optimization with CRN. We select the checkpoint at step 100 for evaluation. The details are provided in Appendix C.

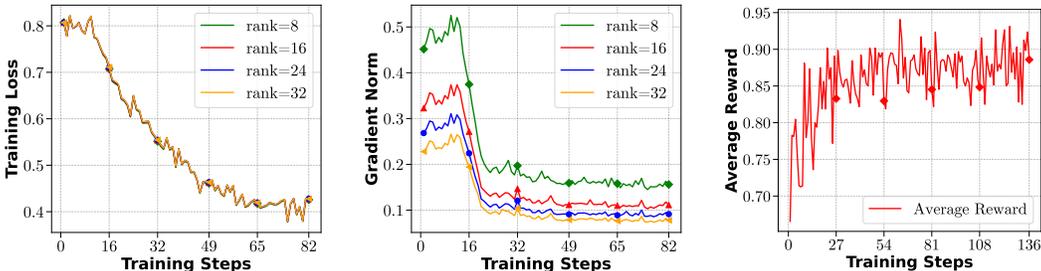


Figure 3: Training dynamics of MetaFlow. **Left:** SFT training loss across different LoRA ranks (with the same learning rate $lr = 0.00002$). **Middle:** SFT gradient norms showing convergent behavior. **Right:** RLVR average reward demonstrating rapid policy improvement and stabilization.

5.2 LOW-COST OPERATOR INTEGRATION: CONTINUOUS LIBRARY EXPANSION

A key advantage of **MetaFlow** over task-level search methods (Zhang et al., 2024b; Hu et al., 2024) is its ability to continuously expand the operator library without costly re-optimization. To validate this capability, we demonstrate a practical workflow: abstracting a manually-tested reasoning pattern into a reusable operator. Specifically, we introduce the **SelfConsistency** operator (see Appendix 2 for complete implementation), which encapsulates the Self-Consistency pattern (Wang et al., 2022) (parallel generation with majority voting), directly derived from one of our experimental baselines (CoT SC in Table 1). Critically, incorporating this operator requires no model retraining, only defining its natural language interface following the format described in Listing 5.

To demonstrate immediate operator utilization, we generate workflows for GSM8K mathematical reasoning after integrating **SelfConsistency**. Listing 13 shows an example generated workflow that successfully invokes the newly integrated operator with parallel sampling ($n=5$) and similarity-based voting, demonstrating that **MetaFlow** can immediately compose effective workflows with novel operators through single-inference generation without any retraining.

Implications. This validates **MetaFlow**’s *low-cost continuous extensibility*: (1) **Minimal integration cost**: introducing a new operator requires only defining its interface, without model retraining or iterative search (hours for MCTS/evolutionary methods vs. minutes for **MetaFlow**); (2) **Rapid abstraction of proven patterns**: practitioners can quickly encapsulate manually-tested workflow logic discovered through experimentation into reusable operators; (3) **Immediate high-quality generation**: generating effective workflows with the new operator costs merely a single inference; (4) **Scalable library growth**: each new operator becomes immediately available across all tasks without proportional computational overhead, enabling continuous knowledge accumulation. This paradigm transforms workflow optimization from isolated task-specific searches into a scalable knowledge base where proven patterns become reusable primitives.

5.3 ZERO-SHOT GENERALIZATION TO NOVEL OPERATORS

To validate true out-of-distribution (OOD) generalization, we evaluate **MetaFlow** on HotpotQA (Yang et al., 2018) multi-hop question answering with the VectorSearch operator (Listing 3)—entirely unseen during training. This tests two critical OOD dimensions: (1) **domain shift** from math/code reasoning to retrieval-based QA, and (2) **novel operator** requiring hybrid document and sentence-level vector search. We generate 100 candidate workflows with Ops = {Generate, Summarize, Revise, Ensemble, VectorSearch} evaluated on 100 HotpotQA instances (Validation), compared against a CoT+RAG baseline using the same API (Qwen-Turbo).

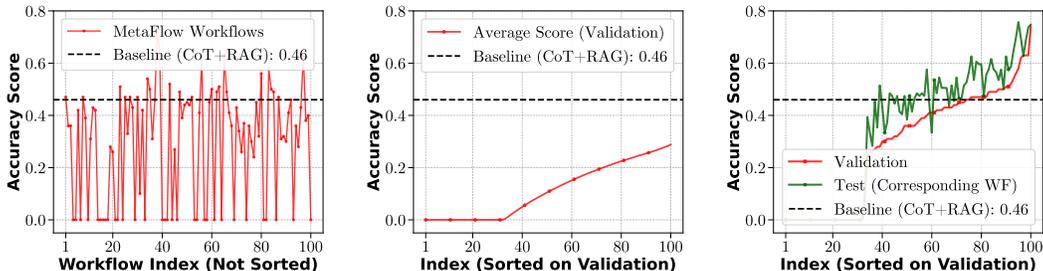


Figure 4: Zero-shot generalization with novel VectorSearch operator on HotpotQA. **Left:** Individual workflow performance across 100 generations shows high variance, 31 workflows fail completely while the best achieves 0.74 accuracy. **Middle:** Average score converges to 0.29, reflecting the exploration cost of zero-shot generation. **Right:** Best-of- k score reaches 0.74 (60.9% above baseline 0.46), surpassing baseline at $k = 76$.

Figure 4 reveals the characteristics of zero-shot workflow generation. The **left panel** shows substantial performance variance: while 31% of workflows fail (score=0.0, likely due to syntax errors or incorrect operator usage), successful workflows demonstrate strong performance, with the best achieving 0.74 accuracy. The **middle panel** tracks the average score across generated workflows, which converges to 0.29, below the 0.46 baseline due to the high failure rate inherent in zero-shot generation. However, the **right panel** demonstrates that the best-of- k selection strategy efficiently identifies high-quality solutions: the best score monotonically increases, surpassing the baseline at $k = 76$ workflows and reaching 0.74 (**60.9% improvement**) with the test set (another 100 instances) with the same trend of performance. This validates that moderate sampling suffices to discover effective workflows despite zero-shot exploration risks. The highest-scoring workflow (Listing 14 in Appendix F) implements sophisticated multi-hop reasoning through iterative retrieval-generation cycles—extracting entities, performing two-stage document search with connection analysis, and synthesizing information across retrieved contexts—demonstrating **MetaFlow**’s ability to compose complex operator sequences for unseen tools. The results confirm our meta-learning approach: train once on diverse tasks, then rapidly adapt to new operators and domains without re-optimization.

5.4 MAIN RESULTS AND ANALYSIS

Experimental Setup. We evaluate **MetaFlow** on benchmarks spanning question answering (HotpotQA, DROP: 1,000 instances each), code generation (MBPP, HumanEval: full sets), and mathematical reasoning (GSM8K: 1,000 instances, MATH: Level-5 problems), following ScoreFlow’s configuration with 1:4 train-test splits. We compare against two baseline categories: (1) **manually designed workflows**—IO, CoT (Wei et al., 2022), CoT SC (Wang et al., 2022), MedPrompt (Nori et al., 2023), MultiPersona (Wang et al., 2024), Self-Refine (Madaan et al., 2024); (2) **automatically optimized workflows**—ADAS (Hu et al., 2024), AFlow (Zhang et al., 2024b), ScoreFlow (Wang et al., 2025). The Planner Qwen3-8B uses base operators {Generate, Summarize, Revise, Ensemble} with dynamic prompt rewriting, plus novel operators {Decompose, Programmer} for OOD testing (natural language descriptions provided at inference). For each task, we generate 20 candidate workflows, validate on 50 instances, and select the best for testing. All methods use GPT-4o-mini-0718 as executor and judge. Workflow execution is orchestrated by MetaGPT.

Performance Analysis. Table 1 shows **MetaFlow** achieves 78.8 average accuracy, competitive with ScoreFlow (82.5), AFlow (78.3), and ADAS (73.1) while requiring only single-inference generation versus their resource-intensive per-instance optimization. Notably, **MetaFlow** matches or

exceeds manually designed workflows and surpasses many automated methods on individual tasks (e.g., GSM8K: 93.8). The framework demonstrates strong performance across diverse domains—mathematical reasoning (GSM8K, MATH), reading comprehension (DROP), and code generation (MBPP)—validating our meta-learning approach’s cross-task generalization. All baseline scores are from ScoreFlow (Wang et al., 2025). And FlowReasoner (Gao et al., 2025) achieves 82.19 on MBPP with GPT-4o-mini-0718 as executor.

Method	DROP	MBPP	GSM8K	MATH	Avg
IO	81.6	69.5	89.1	52.2	73.1
CoT (Wei et al., 2022)	83.2	70.4	88.3	53.4	73.8
CoT SC (Wang et al., 2022)	83.2	71.3	88.6	53.8	74.2
MedPrompt (Nori et al., 2023)	83.0	69.2	88.1	53.7	73.5
MultiPersona (Wang et al., 2024)	81.3	70.4	89.8	51.9	73.4
Self Refine (Madaan et al., 2024)	82.5	70.0	87.5	50.0	72.5
ADAS (Hu et al., 2024)	81.3	68.7	90.5	51.7	73.1
AFlow (Zhang et al., 2024b)	83.5	82.9	90.8	55.8	78.3
ScoreFlow (Wang et al., 2025)	86.2	84.7	94.6	64.4	82.5
Ours	82.8	77.5	93.8	61.0	78.8

Table 1: Performance comparison across multiple benchmarks. **MetaFlow** achieves 78.8 average score with single-inference generation, demonstrating competitive performance and strong cross-domain generalization.

Understanding the Performance Gap: Paradigm Differences and Trade-offs. Our average performance (78.8) trails ScoreFlow (82.5) by 3.7 points, which reflects fundamental differences in problem formulation rather than algorithmic limitations. This gap stems from three interacting factors. *First, risk amplification from task-level generation:* As a cross-task generator, **MetaFlow** produces one workflow W_C per task—any syntax error or malformed operator call yields zero accuracy across all test instances. In contrast, ScoreFlow’s instance-level approach generates W_p per problem—a failed generation affects only one data point. This risk materializes in our VectorSearch experiments (Figure 4 left): 31% of generated workflows fail to execute, directly penalizing task-level metrics. *Second, the complexity-performance trade-off:* While simple workflows dominate on tasks solvable by direct reasoning, this advantage inverts for complex tool integration scenarios. On HotpotQA with the novel VectorSearch operator (Section 5.3), the simple CoT+RAG baseline achieves 0.46 accuracy. MetaFlow’s best workflow—requiring sophisticated orchestration of retrieval, decomposition, and reasoning operators—reaches 0.74 (+60.9% improvement), demonstrating that complex operator composition becomes essential when tool complexity increases. Critically, even with best-of-20 validation selection, our total cost remains orders of magnitude lower than iterative search methods: 20 generations yield a reusable task-level workflow versus thousands of MCTS evaluations (AFlow) or N per-instance generations (ScoreFlow). Figure 4 (right) shows best-of- k scaling stabilizes around $k \approx 80$, where marginal sampling cost ($< \$1$ with Qwen3-8B) far outweighs other automatically designed workflows baselines.

6 CONCLUSION

We introduce MetaFlow, a meta-learning framework that trains language models to generate task-level workflows through a two-stage paradigm combining supervised fine-tuning with reinforcement learning across diverse task-operator combinations. Across benchmarks in question answering, code generation, and mathematical reasoning, MetaFlow achieves competitive performance (78.8 average accuracy) with single-inference generation, while demonstrating strong zero-shot generalization to novel operators and domains. The observed 31% syntax error rate in zero-shot generation highlights a key limitation of single-inference approaches—future work could explore multi-turn reinforcement learning where the model iteratively refines workflows through interaction with execution feedback, potentially combining the efficiency of learned meta-strategies with the robustness of adaptive generation.

REFERENCES

- 540
541
542 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
543 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language
544 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 545
546 Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje F Karlsson, Jie Fu, and Yemin
547 Shi. Autoagents: A framework for automatic agent generation. *arXiv preprint arXiv:2309.17288*,
2023.
- 548
549 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared
550 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
551 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,
552 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,
553 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-
554 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex
555 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,
556 Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa,
557 Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob
558 McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating
large language models trained on code. *CoRR*, abs/2107.03374, 2021.
- 559
560 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
561 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
562 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
2021.
- 563
564 Hongxin Ding, Yue Fang, Runchuan Zhu, Xinke Jiang, Jinyang Zhang, Yongxin Xu, Xu Chu, Jun-
565 feng Zhao, and Yasha Wang. 3ds: Decomposed difficulty data selection’s case study on llm
566 medical domain adaptation. *arXiv preprint arXiv:2410.10901*, 2024.
- 567
568 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner.
569 DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In
NAACL-HLT (1), pp. 2368–2378. Association for Computational Linguistics, 2019.
- 570
571 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation
572 of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, pp.
1126–1135. PMLR, 2017.
- 573
574 Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. Bilevel
575 programming for hyperparameter optimization and meta-learning. In *Proceedings of the 35th*
576 *International Conference on Machine Learning*, pp. 1568–1577. PMLR, 2018.
- 577
578 Hongcheng Gao, Yue Liu, Yufei He, Longxu Dou, Chao Du, Zhijie Deng, Bryan Hooi, Min
579 Lin, and Tianyu Pang. Flowreasoner: Reinforcing query-level meta-agents. *arXiv preprint*
arXiv:2504.15257, 2025.
- 580
581 Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian,
582 and Yujiu Yang. Connecting large language models with evolutionary algorithms yields powerful
583 prompt optimizers. *arXiv preprint arXiv:2309.08532*, 2023.
- 584
585 Ilgee Hong, Zichong Li, Alexander Bukharin, Yixiao Li, Haoming Jiang, Tianbao Yang, and Tuo
586 Zhao. Adaptive preference scaling for reinforcement learning with human feedback. *Advances in*
Neural Information Processing Systems, 37:107249–107269, 2024.
- 587
588 Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang,
589 Steven Ka Shing Zhang, Zijuan Gui, et al. Metagpt: Meta programming for a multi-agent collab-
orative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- 590
591 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
592 Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- 593
Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. *arXiv preprint*
arXiv:2408.08435, 2024.

- 594 Bowen Jiang, Runchuan Zhu, Jiang Wu, Zinco Jiang, Yifan He, Junyuan Gao, Jia Yu, Rui Min,
595 Yinfan Wang, Haote Yang, et al. Evaluating large language model with knowledge oriented
596 language specific simple question answering. *arXiv preprint arXiv:2505.16591*, 2025.
597
- 598 Omar Khattab, Bhanukiran Vinzamuri Akula, et al. Dspy: Expressive, modular prompting for
599 language models. *arXiv preprint arXiv:2310.01348*, 2023.
- 600 Jack P. C. Kleijnen. Antithetic variates, common random numbers and optimal computer time
601 allocation in simulation. *Management Science*, 21(10):1176–1185, 1975.
602
- 603 Chenyang Li, Ziqiang Wang, Dong Zhang, Xue Zhao, Cheng Wang, Xingwu Wang, Yuan Wang,
604 Haifeng Zhang, and Wenwu Zhu. Scoreflow: Mastering llm agent workflows via score-based
605 preference optimization. *arXiv preprint arXiv:2502.04306*, 2025.
- 606 Zelong Li, Shuyuan Xu, Kai Mei, Wenyue Hua, Balaji Rama, Om Raheja, Hao Wang, He Zhu, and
607 Yongfeng Zhang. Autoflow: Automated workflow generation for large language model agents.
608 *arXiv preprint arXiv:2407.12821*, 2024.
609
- 610 Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. A dynamic llm-powered agent network
611 for task-oriented agent collaboration. In *First Conference on Language Modeling*, 2024.
- 612 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
613 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
614 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
615
- 616 Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King,
617 Jonathan Larson, Yuanzhi Li, Weishung Liu, et al. Can generalist foundation models outcompete
618 special-purpose tuning? case study in medicine. *arXiv preprint arXiv:2311.16452*, 2023.
- 619 OpenAI. Aime benchmark for mathematical reasoning. <https://openai.com/research>, 2023. Ac-
620 cessed 2024.
621
- 622 Zhihong Shao, Peiyi Yuan, Hongsheng Li, Yizhe Wang, Yubo Xu, Xiaoke Sun, Ke Liu, Yuanhan
623 Lin, Chunyan Yue, Kun Chen, et al. Deepseekmath: Pushing the limits of mathematical reasoning
624 in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 625 Linxin Song, Jiale Liu, Jieyu Zhang, Shaokun Zhang, Ao Luo, Shijian Wang, Qingyun Wu, and
626 Chi Wang. Adaptive in-conversation team building for language model agents. *arXiv preprint
627 arXiv:2405.19425*, 2024.
628
- 629 Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2, 2024.
- 630 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
631 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
632 *The Eleventh International Conference on Learning Representations*, 2022.
633
- 634 Yinjie Wang, Ling Yang, Guohao Li, Mengdi Wang, and Bryon Aragam. Scoreflow: Mastering
635 llm agent workflows via score-based preference optimization. *arXiv preprint arXiv:2502.04306*,
636 2025.
- 637 Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. Unleashing the
638 emergent cognitive synergy in large language models: A task-solving agent through multi-persona
639 self-collaboration. In *Proceedings of the 2024 Conference of the North American Chapter of the
640 Association for Computational Linguistics: Human Language Technologies (Volume 1: Long
641 Papers)*, pp. 257–279, 2024.
- 642 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
643 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
644 neural information processing systems*, 35:24824–24837, 2022.
645
- 646 Zhenran Xu, Yiyu Wang, Xue Yang, Longyue Wang, Weihua Luo, Kaifu Zhang, Baotian Hu, and
647 Min Zhang. Comfyui-r1: Exploring reasoning models for workflow generation. *arXiv preprint
arXiv:2506.09790*, 2025.

648 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
649 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
650 *arXiv:2505.09388*, 2025.

651
652 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov,
653 and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question
654 answering. *arXiv preprint arXiv:1809.09600*, 2018.

655 Guibin Zhang, Yanwei Yue, Xiangguo Sun, Guancheng Wan, Miao Yu, Junfeng Fang, Kun Wang,
656 Tianlong Chen, and Dawei Cheng. G-designer: Architecting multi-agent communication topolo-
657 gies via graph neural networks. *arXiv preprint arXiv:2410.11782*, 2024a.

658 Jiayi Zhang, Yihang Xiang, Chao Wang, Amina Zhou, Jiaqi Lu, Di Chen, Chaowei He, Yanshuai
659 Wang, Bin Ding, Dacheng Gao, et al. Aflow: Automating agentic workflow generation. *arXiv*
660 *preprint arXiv:2410.10762*, 2024b.

661
662 Zeyu Zhu et al. Olympiadbench: A benchmark for mathematical reasoning at the olympiad level.
663 *arXiv preprint arXiv:2402.00000*, 2024.

664
665 Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen
666 Schmidhuber. Gptswarm: Language agents as optimizable graphs. In *Forty-first International*
667 *Conference on Machine Learning*, 2024.

668 A METAFLow ARCHITECTURE

669 A.1 OPERATORS DESIGN

670
671 For the basic operators defined in the left part of Figure 2, their implementations are identical to
672 those in the AFlow (Zhang et al., 2024b) codebase, and thus we omit their detailed descriptions
673 here. In the following, we introduce the newly defined operators that we have designed for this
674 work.
675
676

677 **Decompose Operator:** Beyond the basic text-processing operators, we introduce the Decompose
678 operator to handle complex problem-solving scenarios that require hierarchical decomposition. This
679 operator breaks down intricate problems into structured subproblems with explicit dependency re-
680 lationships, enabling workflows to tackle multi-step reasoning tasks systematically. The operator is
681 invoked as: `await self.decompose(instruction: str, context: str)` and returns a list of
682 subproblem dictionaries, each containing an ID, description, and dependencies. The following code
683 block shows the implementation details of this operator.

```
684 class Decompose(Operator):
685     """
686     Core Operator: Decompose
687     Breaks down complex problems into manageable subproblems.
688     """
689     async def __call__(self, instruction: str = "", context: str = "") -> List[Dict[str, str]]:
690         prompt = f"""You are an expert at problem decomposition. Break down the complex problem into
691         manageable subproblems.
692
693         **Instruction on decomposition strategy:**
694         {instruction}
695
696         **Problem/Context to Decompose:**
697         {context if context else "No context provided."}
698
699         **Original Problem:**
700         {self.problem_text}
701
702         Your response MUST be valid XML with 'think' and 'subproblems' fields.
703         - In "think", explain your decomposition strategy
704         - In "subproblems", provide a list where each item has:
705           - id: unique identifier (e.g., "sub1", "sub2")
706           - description: clear description of the subproblem
707           - dependencies: comma-separated IDs of prerequisite subproblems (empty if none)
708
709         **EXAMPLE:**
710         <think>This problem requires finding area then volume.</think>
711         <subproblems>
712         [
713           [{"id": "sub1", "description": "Calculate the radius", "dependencies": ""}],
714           [{"id": "sub2", "description": "Calculate the area", "dependencies": "sub1"}],
715           [{"id": "sub3", "description": "Calculate the volume", "dependencies": "sub2"}]
716         ]
717         """
```

```

702 ]
703 </subproblems>"""
704
705     response = await self._fill_node(DecomposeOp, prompt, mode="xml_fill")
706     response = DecomposeOp(**response)
707     return [sub.dict() for sub in response.subproblems]

```

Listing 1: Implementation of the Decompose operator for hierarchical problem decomposition.

SelfConsistency Operator: To demonstrate **MetaFlow**'s capability for low-cost operator integration and continuous library expansion (Section 5.2), we introduce the SelfConsistency operator—a practical abstraction of the Self-Consistency reasoning pattern (Wang et al., 2022). This operator encapsulates the parallel generation with majority voting strategy, which was originally one of our experimental baselines (CoT SC in Table 1). Critically, incorporating this operator into **MetaFlow** requires *no model retraining*; the planner learns to effectively utilize it through its natural language interface description alone. The following code block presents the complete implementation:

```

715 class SelfConsistency(Operator):
716     """Self-Consistency: multi-path sampling + majority voting for better reasoning."""
717     def __init__(self, llm, problem_text: str = "", n_samples: int = 5,
718                 similarity_threshold: float = 0.85, return_full_info: bool = False):
719         super().__init__(llm, problem_text)
720         self.n_samples = n_samples
721         self.similarity_threshold = similarity_threshold
722         self.return_full_info = return_full_info
723
724     async def __call__(self, instruction: str = "", context: str = "", answer_type: str = "auto") -> str:
725         silent = os.environ.get('SCOREFLOW_SILENT', 'false').lower() == 'true'
726         if not silent:
727             print(f"[SelfConsistency] samples={self.n_samples}, threshold={self.similarity_threshold}")
728
729         paths = await self._parallel_sampling(instruction, context)
730         if not silent:
731             print(f"Generated {len(paths)} paths")
732
733         raw_answers = await self._extract_answers_batch(paths, answer_type)
734         if not silent:
735             print(f"Extracted {len([a for a in raw_answers if a])} answers")
736
737         normalized = [self._normalize_answer(a, answer_type) for a in raw_answers]
738         clustered = self._cluster_similar_answers(normalized)
739         votes = Counter(clustered)
740
741         if not votes:
742             if not silent:
743                 print("No valid answers, fallback")
744             return await self._fallback_generate(instruction, context)
745
746         answer, count = votes.most_common(1)[0]
747         conf = count / len(clustered) if clustered else 0
748         if not silent:
749             print(f"Votes: {dict(votes)} | Answer: {answer} | Conf: {conf:.1%}")
750
751         if self.return_full_info:
752             return self._format_full_response(paths, raw_answers, votes, answer, conf)
753         return answer
754
755     async def _parallel_sampling(self, instruction: str, context: str) -> List[str]:
756         prompt = f"""Solve step by step.
757         **Problem:** {self.problem_text}
758         **Instruction:** {instruction or "Solve carefully."}
759         **Context:** {context or "None"}
760         Show reasoning. End with "Final Answer:" or "The answer is:"
761         **Solution:**"""
762         tasks = [self._sample_single_path(prompt) for _ in range(self.n_samples)]
763         results = await asyncio.gather(*tasks, return_exceptions=True)
764         return [r for r in results if isinstance(r, str) and len(r.strip()) > 10]
765
766     async def _sample_single_path(self, prompt: str) -> str:
767         try:
768             return (await self._fill_node(GenerateOp, prompt, mode="single_fill")).get("response", "")
769         except:
770             return ""
771
772     async def _extract_answers_batch(self, paths: List[str], answer_type: str) -> List[str]:
773         answers = []
774         for p in paths:
775             ans = self._rule_based_extraction(p, answer_type)
776             answers.append(ans if ans else await self._llm_extract_answer(p, answer_type))
777         return answers
778
779     def _rule_based_extraction(self, text: str, answer_type: str) -> str:
780         patterns = [
781             r"(?:final answer|the answer is|answer:)\s*[:\s]*(.+?)\s*(?:\n|\s|\s|$)\s*(?:\s|$)",
782             r"(?:therefore|thus|so|hence),?\s*(?:the answer is)?\s*[:\s]*(.+?)\s*(?:\n|\s|\s|$)\s*(?:\s|$)",
783             r"\boxed\{(.+)\}", r"\*(.+?)\*\s*\$";

```

```

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

for p in patterns:
    m = re.search(p, text.lower(), re.IGNORECASE | re.MULTILINE)
    if m:
        ans = re.sub(r'^[:\s]+[[:\s\.\.]]+$', '', m.group(1).strip())
        if ans and len(ans) < 200:
            return ans
    if answer_type == "choice":
        m = re.search(r"(?:answer|option|choice)[[:\s]]*([A-Za-d])\b", text, re.IGNORECASE)
        if m: return m.group(1).upper()
    if answer_type == "numeric":
        nums = re.findall(r'(?!\d+)(?:\.\d+)?(?:/\d+)?', text)
        if nums: return nums[-1]
    return ""

async def _llm_extract_answer(self, path: str, answer_type: str) -> str:
    hints = {"numeric": "number only", "choice": "letter only", "short": "brief", "boolean": "yes/no only"}
    prompt = f"Extract final answer. {hints.get(answer_type, '')} Return ONLY the answer.\n**Solution **: ** {path[-200:]} \n**Answer:**"
    try:
        ans = (await self._fill_node(GenerateOp, prompt, mode="single_fill")).get("response", "").strip()
        return re.sub(r"^(the answer is|answer:|final answer:)\s*", '', ans, flags=re.IGNORECASE).strip()
    except:
        return ""

def _normalize_answer(self, answer: str, answer_type: str) -> str:
    if not answer: return ""
    n = re.sub(r"^(the |a |an )'", '', answer.strip().lower())
    n = re.sub(r'[,!?:;]+$', '', n)
    if answer_type == "numeric" or re.match(r"^\d+$", n):
        n = re.sub(r'[$%\s]', '', n)
        if '/' in n:
            try: p = n.split('/'); n = str(float(p[0]) / float(p[1])) if len(p) == 2 else n
            except: pass
    elif answer_type == "choice":
        m = re.search(r'[A-Za-d]', n)
        if m: n = m.group(0).upper()
    elif answer_type == "boolean":
        n = 'yes' if n in ['yes', 'true', '1', 'correct', 'right'] else ('no' if n in ['no', 'false', '0', 'incorrect', 'wrong'] else n)
    return n

def _cluster_similar_answers(self, answers: List[str]) -> List[str]:
    valid = [a for a in answers if a]
    if not valid or self.similarity_threshold >= 1.0: return valid
    clusters = []
    for ans in valid:
        matched = False
        for rep, members in clusters:
            if SequenceMatcher(None, ans, rep).ratio() >= self.similarity_threshold:
                members.append(ans); matched = True; break
        if not matched: clusters.append((ans, [ans]))
    return [Counter(m).most_common(1)[0][0] for a in valid for r, m in clusters if a in m]

async def _fallback_generate(self, instruction: str, context: str) -> str:
    prompt = f"Solve: {self.problem_text}\nInstruction: {instruction}\nContext: {context or 'None'}"
    return (await self._fill_node(GenerateOp, prompt, mode="single_fill")).get("response", "No answer")

def _format_full_response(self, paths: List[str], raw: List[str], votes: Counter, answer: str, conf: float) -> str:
    lines = [f"## Self-Consistency Result", f"**Final Answer:** {answer}", f"**Confidence:** {conf:.1%}"]
    lines.append(f"### Votes")
    lines.append([f"- {a}: {c}" for a, c in votes.most_common()])
    return "\n".join(lines)

```

Listing 2: Implementation of the SelfConsistency operator, which abstracts the Self-Consistency pattern (Wang et al., 2022) into a reusable component. This operator performs parallel sampling of multiple reasoning paths (default $n = 5$) and applies majority voting with similarity-based answer clustering to select the most consistent solution. The implementation demonstrates **MetaFlow**'s low-cost extensibility: practitioners can encapsulate manually-tested reasoning patterns into operators without model retraining.

VectorSearch Operator: To evaluate the generalization capability of **MetaFlow** to completely out-of-distribution (OOD) operators, we introduce the VectorSearch operator 3, which is a complex tool-calling operator that was entirely unseen during both the SFT and RLVR training phases. Unlike the basic text-processing operators used in training, VectorSearch requires sophisticated external API interactions with a vector database (ChromaDB) and involves multiple parameters for controlling retrieval behavior. The operator is invoked as: `await self.vector_search(instruction: str, context: str, top_k: int)`, where it performs hybrid retrieval combining document-level and sentence-level semantic search over the HotpotQA knowledge base. Despite its complexity and complete absence from training data, **MetaFlow** successfully learns to incorporate this operator into

generated workflows based solely on its natural language description provided at inference time, as demonstrated in Appendix F.

```

810
811
812
813 class VectorSearch(Operator):
814     """
815     Core Operator: Vector Search
816     Retrieves relevant documents from HotpotQA vector database using RAG system.
817     """
818     def __init__(self, llm, problem_text: str = "", db_config: Dict = None):
819         super().__init__(llm, problem_text)
820         self.config = self._load_config()
821         if db_config:
822             self.config.update(db_config)
823             self._init_chromadb()
824
825     def _load_config(self) -> Dict:
826         """Load configuration from db.config file"""
827         from pathlib import Path
828         config = {}
829         config_file = Path(__file__).parent / "db.config"
830
831         if config_file.exists():
832             with open(config_file, 'r') as f:
833                 for line in f:
834                     line = line.strip()
835                     if line and not line.startswith('#') and '=' in line:
836                         key, value = line.split('=', 1)
837                         key, value = key.strip(), value.strip()
838                         if key in ('DOC_TOP_K', 'SENT_TOP_K'):
839                             config[key.lower()] = int(value)
840                         elif key == 'HYBRID_MODE':
841                             config[key.lower()] = value.lower() == 'true'
842                         else:
843                             config[key.lower()] = value
844         else:
845             config = {
846                 'db_path': '...chroma-db', 'model_path': '...all-MiniLM-L6-v2',
847                 'doc_top_k': 3, 'sent_top_k': 5, 'hybrid_mode': True
848             }
849         return config
850
851     def _init_chromadb(self):
852         """Initialize ChromaDB connection and collections"""
853         try:
854             import chromadb
855             from chromadb.utils import embedding_functions
856             self.client = chromadb.PersistentClient(path=self.config['db_path'])
857             model = self.config['model_path'] if os.path.exists(self.config['model_path']) else "all-MiniLM-L6-v2"
858             self.embedding_function = embedding_functions.SentenceTransformerEmbeddingFunction(model_name=model)
859             self.doc_collection = self.client.get_collection("documents")
860             self.sent_collection = self.client.get_collection("sentences")
861         except Exception as e:
862             self.doc_collection = self.sent_collection = None
863
864     async def __call__(self, instruction: str = "", context: str = "", top_k: int = None) -> str:
865         """Execute vector search and return formatted retrieved documents."""
866         if not self.doc_collection or not self.sent_collection:
867             return "Error: Vector database not initialized."
868         doc_k = top_k or self.config['doc_top_k']
869         query = await self._process_query(instruction, context)
870         retrieved_data = self._hybrid_retrieval(query, doc_k, self.config['sent_top_k'])
871         return self._format_context(retrieved_data)
872
873     async def _process_query(self, instruction: str, context: str) -> str:
874         """Process and combine query from instruction and context."""
875         if instruction and context:
876             return f"{instruction} {context}"
877         return instruction or (context[:200] if context else self.problem_text[:200] or "general information")
878
879     def _hybrid_retrieval(self, query: str, doc_k: int, sent_k: int) -> Dict:
880         """Perform hybrid document and sentence level retrieval."""
881         retrieved_data = {'documents': [], 'scores': []}
882         try:
883             if self.config.get('hybrid_mode', True):
884                 doc_results = self.doc_collection.query(query_texts=[query], n_results=doc_k)
885                 for i, (doc_id, doc_text, meta, dist) in enumerate(zip(
886                     doc_results['ids'][0], doc_results['documents'][0],
887                     doc_results['metadatas'][0], doc_results['distances'][0])):
888                     retrieved_data['documents'].append({
889                         'doc_id': doc_id, 'title': meta.get('title', 'Unknown'),
890                         'text': doc_text[:500], 'type': 'document', 'rank': i + 1})
891                     retrieved_data['scores'].append(float(dist))
892                 sent_results = self.sent_collection.query(query_texts=[query], n_results=sent_k)
893                 for i, (sid, stxt, meta, dist) in enumerate(zip(
894                     sent_results['ids'][0], sent_results['documents'][0],
895                     sent_results['metadatas'][0], sent_results['distances'][0])):
896                     if i < 3:
897                         retrieved_data['documents'].append({

```

```

864         'doc_id': sid, 'title': meta.get('title', 'Unknown'),
865         'text': stxt, 'type': 'sentence', 'rank': i + 1})
866         retrieved_data['scores'].append(float(dist))
867     except Exception:
868         pass
869     return retrieved_data
870
871     def _format_context(self, retrieved_data: Dict) -> str:
872         """Format retrieved documents into readable context string."""
873         if not retrieved_data['documents']:
874             return "No relevant documents found."
875         parts = ["**Retrieved Information:**\n"]
876         for d in retrieved_data['documents']:
877             parts.append(f"[{d['type']}.title()]: {d['title']} {d['text']}\n")
878         return "\n".join(parts)

```

Listing 3: Implementation of the VectorSearch operator for retrieval-augmented generation.

875 A.2 INPUT AND OUTPUT OF THE MODEL

876 The MODEL INPUT consists of two parts: (1) TASK TYPE DESCRIPTION, which elucidates the domain characteristics of the target problem family and the input-output formats of each belonging problem instance. The block below shows the complete TASK TYPE DESCRIPTION for GSM8K.

```

881 ### Problem Domain Overview
882 This domain tests multi-step mathematical reasoning using basic arithmetic operations.
883
884 #### Key Characteristics & Requirements
885 - **Answer Format**: Single numerical value (integer or decimal)
886 - **Solution Steps**: 2-8 step reasoning chains using +, -, *, /
887 - **Critical**: Track intermediate results and units throughout
888 - **Validation**: Final answer must be numerically exact
889 - **No Complex Math**: Only elementary arithmetic, no algebra or calculus
890
891 #### Common Problem Types & Solution Strategies
892 - **Sequential Operations**: Step-by-step calculations building on previous results
893 - **Rate Problems**: Distance/speed/time, work rates, unit prices
894 - **Distribution**: Dividing quantities, equal sharing, remainders
895 - **Proportions**: Percentages, fractions, ratios, scaling
896 - **Multi-entity**: Track different quantities for multiple people/objects
897
898 #### Workflow Focus Points
899 1. Extract all numerical values and their context
900 2. Identify what the question asks for
901 3. Build step-by-step calculation chain
902 4. Show intermediate results explicitly
903 5. Return final numerical answer only
904
905 #### Input Format
906 ...
907 ---
908 **QUESTION:**
909 [Complete word problem text]
910 ---
911 ...
912 Multiple problems follow the same structure if provided.

```

Listing 4: Complete TASK TYPE DESCRIPTION of GSM8K including domain overview, key requirements and the format of problem instance.

902 (2) OPERATOR DESCRIPTIONS, which define the functions, parameters, and input-output formats of each available operator (whether pre-set or user-defined). The block below shows the complete OPERATOR DESCRIPTIONS of the operator set Ops = (Generate, Summarize, Revise, Ensemble).

```

913 ### Available Operators & Building Blocks
914 All operators follow a consistent interface pattern and are initialized with the problem text.
915
916 #### **Core Operators**
917
918 **Generate: CREATE new information**
919 - **Signature**: 'await self.generate(instruction: str, context: str = "") -> str'
920 - **Purpose**: Produces new text, analysis, or reasoning based on strategic instructions
921
922 **Revise: IMPROVE existing information**
923 - **Signature**: 'await self.revise(instruction: str, context: str) -> str'
924 - **Purpose**: Critiques and refines existing text based on specific improvement criteria
925
926 **Summarize: COMPRESS information**
927 - **Signature**: 'await self.summarize(instruction: str, context: str) -> str'
928 - **Purpose**: Condenses text while preserving key information relevant to the problem
929
930 **Ensemble: DECIDE between or synthesize options**
931 - **Signature**: 'await self.ensemble(instruction: str, contexts: List[str]) -> str'
932 - **Purpose**: Evaluates, compares, or merges multiple candidate solutions

```

Listing 5: Complete OPERATOR DESCRIPTIONS of an operator set.

The OUTPUT requires a structured WORKFLOW aimed at efficiently solving all problem instances under the task type using the given operator set based on the template below.

```

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
'''python
# --- DO NOT IMPORT HERE ---
class Workflow:
    def __init__(self, config, problem) -> None:
        # --- DO NOT MODIFY THIS SECTION ---
        self.config = config
        self.problem_text = problem
        self.llm = create(config)

        self.generate = operator.Generate(self.llm, self.problem_text)
        self.revise = operator.Revise(self.llm, self.problem_text)
        self.summarize = operator.Summarize(self.llm, self.problem_text)
        self.ensemble = operator.Ensemble(self.llm, self.problem_text)

    async def run_workflow(self):
        """
        Implement the core problem-solving logic here.
        """
        import asyncio
        # --- YOUR WORKFLOW LOGIC HERE ---
...

```

Listing 6: OUTPUT template of given Ops = (Generate, Summarize, Revise, Ensemble).

B DETAILS OF SUPERVISED FINE-TUNING

The supervised fine-tuning (SFT) stage initializes the planner model π_θ to generate syntactically correct workflows following the required template structure (see Appendix A.2 for the detailed input-output format). This stage addresses the cold start problem by teaching the model the basic grammar of workflow construction before reinforcement learning optimization.

Dataset Construction Pipeline. We construct approximately 1,300 high-quality $([C, Ops], W)$ pairs using Qwen-Max as the expert model. The construction follows a four-stage pipeline:

1. **Prompt Crafting:** Construct input prompts combining task descriptions C (domain characteristics, input-output formats) and operator specifications Ops (function signatures, purposes) following the format defined in Appendix A.2.
2. **Expert Generation:** Feed prompts to Qwen-Max API to obtain comprehensive responses including workflow explanations and complete executable code.
3. **Code Extraction:** Parse API responses to extract clean workflow code W , discarding natural language explanations.
4. **Quality Verification:** Execute each extracted workflow on 1-2 randomly sampled problem instances $p \sim C$ to ensure (a) no syntax errors and (b) correct solutions. Only validated workflows are retained.

This pipeline is applied across four tasks (GSM8K, DROP, MBPP, Humaneval) with the basic operator set $Ops^{SFT} = (\text{Generate, Summarize, Revise, Ensemble})$, producing diverse workflow patterns that establish the foundation for RLVr optimization. To ensure stable training on this limited dataset, we employ parameter-efficient fine-tuning with LoRA configuration, as detailed in Section B.2.

B.1 DATASET CONSTRUCTION EXAMPLES

To illustrate the construction pipeline, we present a concrete example for the GSM8K mathematical reasoning task. The input prompt (Listing 7) combines task description and operator specifications:

```

966
967
968
969
970
971
### 1. Problem Domain Overview
This domain tests multi-step mathematical reasoning using basic arithmetic operations.

#### Key Characteristics & Requirements
- **Answer Format**: Single numerical value (integer or decimal)
- **Solution Steps**: 2-8 step reasoning chains using +, -, *, /
- **Critical**: Track intermediate results and units throughout
- **Validation**: Final answer must be numerically exact
- **No Complex Math**: Only elementary arithmetic, no algebra or calculus

#### Common Problem Types & Solution Strategies
- **Sequential Operations**: Step-by-step calculations building on previous results

```

```

972 - **Rate Problems**: Distance/speed/time, work rates, unit prices
973 - **Distribution**: Dividing quantities, equal sharing, remainders
974 - **Proportions**: Percentages, fractions, ratios, scaling
975 - **Multi-entity**: Track different quantities for multiple people/objects
976
977 ##### Workflow Focus Points
978 1. Extract all numerical values and their context
979 2. Identify what the question asks for
980 3. Build step-by-step calculation chain
981 4. Show intermediate results explicitly
982 5. Return final numerical answer only
983
984 ##### 2. Available Operators & Building Blocks
985 All operators follow a consistent interface pattern and are initialized with the problem text.
986
987 ##### **Core Operators**
988 **Generate: CREATE new information**
989 - **Signature**: 'await self.generate(instruction: str, context: str = "") -> str'
990 - **Purpose**: Produces new text, analysis, or reasoning based on strategic instructions
991
992 **Revise: IMPROVE existing information**
993 - **Signature**: 'await self.revise(instruction: str, context: str) -> str'
994 - **Purpose**: Critiques and refines existing text based on specific improvement criteria
995
996 **Summarize: COMPRESS information**
997 - **Signature**: 'await self.summarize(instruction: str, context: str) -> str'
998 - **Purpose**: Condenses text while preserving key information relevant to the problem
999
1000 **Ensemble: DECIDE between or synthesize options**
1001 - **Signature**: 'await self.ensemble(instruction: str, contexts: List[str]) -> str'
1002 - **Purpose**: Evaluates, compares, or merges multiple candidate solutions
1003
1004 ##### Input Format
1005 '''
1006 ---
1007 **QUESTION**:
1008 [Complete word problem text]
1009 ---
1010 '''
1011 Multiple problems follow the same structure if provided.
1012
1013 ##### 3. Your Task: Complete the 'run_workflow' Method
1014 Your task is to write the Python code for the 'run_workflow' method within the provided template.
1015 Focus on creating a robust, reusable workflow that leverages detailed instructions.
1016
1017 **Base Template**:
1018 '''python
1019 # --- DO NOT IMPORT HERE ---
1020 class Workflow:
1021     def __init__(self, config, problem) -> None:
1022         # --- DO NOT MODIFY THIS SECTION ---
1023         self.config = config
1024         self.problem_text = problem
1025         self.llm = create(config)
1026
1027         self.generate = operator.Generate(self.llm, self.problem_text)
1028         self.revise = operator.Revise(self.llm, self.problem_text)
1029         self.summarize = operator.Summarize(self.llm, self.problem_text)
1030         self.ensemble = operator.Ensemble(self.llm, self.problem_text)
1031
1032     async def run_workflow(self):
1033         """
1034         Implement the core problem-solving logic here.
1035         """
1036         import asyncio
1037         # --- YOUR WORKFLOW LOGIC HERE ---
1038     ...

```

Listing 7: Input prompt for Qwen-Max to generate workflow examples (GSM8K task).

Given this input prompt, Qwen-Max generates a comprehensive response that includes both natural language explanation of the workflow design and the complete workflow code. Listing 8 shows the raw output from the API:

```

1017 Now let me understand the task better and implement the 'run_workflow' method.
1018 This is a workflow for solving multi-step mathematical reasoning problems using the provided operators.
1019 I've created a robust workflow for multi-step mathematical reasoning problems. Here's how it works:
1020
1021 '''python
1022 # --- DO NOT IMPORT HERE ---
1023 class Workflow:
1024     def __init__(self, config, problem) -> None:
1025         # --- DO NOT MODIFY THIS SECTION ---
1026         self.config = config
1027         self.problem_text = problem
1028         self.llm = create(config)
1029
1030         self.generate = operator.Generate(self.llm, self.problem_text)
1031         self.revise = operator.Revise(self.llm, self.problem_text)

```

```

1026     self.summarize = operator.Summarize(self.llm, self.problem_text)
1027     self.ensemble = operator.Ensemble(self.llm, self.problem_text)
1028
1029     async def run_workflow(self):
1030         """
1031         Implement the core problem-solving logic here.
1032         """
1033         # Step 1: Extract all numerical values and their context
1034         extracted_info = await self.generate(
1035             instruction="Extract all numerical values from the problem, including their units and
1036             contextual meaning (e.g., '5 apples', 'speed of 60 km/h'). List them clearly.",
1037             context=self.problem_text
1038         )
1039         # Step 2: Identify what the question asks for
1040         target_query = await self.generate(
1041             instruction="Based on the problem, state precisely what numerical value is being asked for,
1042             including its expected unit if any.",
1043             context=self.problem_text
1044         )
1045         # Step 3: Build a step-by-step calculation plan
1046         plan = await self.generate(
1047             instruction="Create a clear, step-by-step plan (2 to 8 steps) to solve the problem using only
1048             +, -, *, / operations. Reference the extracted numbers and specify what each step computes.",
1049             context=f"Extracted info: {extracted_info}\nTarget: {target_query}"
1050         )
1051         # Step 4: Execute the plan with explicit intermediate results
1052         calculation = await self.generate(
1053             instruction="Perform the arithmetic step-by-step as per the plan. Show each intermediate result
1054             with its unit or meaning. Ensure only elementary arithmetic is used.",
1055             context=plan
1056         )
1057         # Step 5: Extract and validate the final numerical answer
1058         final_answer_raw = await self.generate(
1059             instruction="From the calculation steps, extract only the final numerical answer (as an integer
1060             or decimal). Do not include units, explanations, or text -- only the number.",
1061             context=calculation
1062         )
1063         # Step 6: Revise to ensure numerical exactness and proper format
1064         final_answer_clean = await self.revise(
1065             instruction="Ensure the output is a single numerical value (e.g., 42 or 15.75), with no extra
1066             characters, units, or formatting. If it's a whole number, do not include decimal places unless
1067             required.",
1068             context=final_answer_raw
1069         )
1070         # Step 7: Return the final answer
1071         return final_answer_clean.strip()
1072     """
1073
1074     **Workflow Design:**
1075
1076     1. **Problem Extraction** - Uses 'generate' to parse the problem, extracting all numerical values, the goal
1077     , and relationships between quantities
1078
1079     2. **Parallel Solution Generation** - Generates three independent solutions using different approaches:
1080     - Standard step-by-step forward calculation
1081     - Alternative verification method
1082     - Working backwards from the goal
1083
1084     3. **Solution Revision** - Each solution is revised to catch arithmetic errors, verifying each calculation
1085     step
1086
1087     4. **Ensemble Decision** - Compares all three verified solutions to determine the correct answer through
1088     majority vote or error analysis
1089
1090     5. **Answer Extraction** - Uses 'summarize' to extract just the final numerical value
1091
1092     **Key Features:**
1093     - Uses 'asyncio.gather' for parallel execution (generating and revising solutions simultaneously)
1094     - Multiple independent solutions reduce the chance of systematic errors
1095     - Revision step catches arithmetic mistakes
1096     - Ensemble aggregates results for higher accuracy
1097     - Clean numerical output extraction at the end

```

Listing 8: Raw output from Qwen-Max API containing workflow explanation and code.

Finally, we parse the API response to extract only the executable workflow code, discarding the natural language explanations. This extracted code forms a single training example in our SFT dataset. Listing 9 shows the clean workflow code extracted from the above output:

```

1077 # --- DO NOT IMPORT HERE ---
1078 class Workflow:
1079     def __init__(self, config, problem) -> None:
1080         # --- DO NOT MODIFY THIS SECTION ---
1081         self.config = config

```

```

1080     self.problem_text = problem
1081     self.llm = create(config)
1082
1083     self.generate = operator.Generate(self.llm, self.problem_text)
1084     self.revise = operator.Revise(self.llm, self.problem_text)
1085     self.summarize = operator.Summarize(self.llm, self.problem_text)
1086     self.ensemble = operator.Ensemble(self.llm, self.problem_text)
1087
1088     async def run_workflow(self):
1089         """
1090         Implement the core problem-solving logic here.
1091         """
1092         # Step 1: Extract all numerical values and their context
1093         extracted_info = await self.generate(
1094             instruction="Extract all numerical values from the problem, including their units and
1095             contextual meaning (e.g., '5 apples', 'speed of 60 km/h'). List them clearly.",
1096             context=self.problem_text
1097         )
1098
1099         # Step 2: Identify what the question asks for
1100         target_query = await self.generate(
1101             instruction="Based on the problem, state precisely what numerical value is being asked for,
1102             including its expected unit if any.",
1103             context=self.problem_text
1104         )
1105
1106         # Step 3: Build a step-by-step calculation plan
1107         plan = await self.generate(
1108             instruction="Create a clear, step-by-step plan (2 to 8 steps) to solve the problem using only
1109             +, -, *, / operations. Reference the extracted numbers and specify what each step computes.",
1110             context=f"Extracted info: {extracted_info}\nTarget: {target_query}"
1111         )
1112
1113         # Step 4: Execute the plan with explicit intermediate results
1114         calculation = await self.generate(
1115             instruction="Perform the arithmetic step-by-step as per the plan. Show each intermediate result
1116             with its unit or meaning. Ensure only elementary arithmetic is used.",
1117             context=plan
1118         )
1119
1120         # Step 5: Extract and validate the final numerical answer
1121         final_answer_raw = await self.generate(
1122             instruction="From the calculation steps, extract only the final numerical answer (as an integer
1123             or decimal). Do not include units, explanations, or text -- only the number.",
1124             context=calculation
1125         )
1126
1127         # Step 6: Revise to ensure numerical exactness and proper format
1128         final_answer_clean = await self.revise(
1129             instruction="Ensure the output is a single numerical value (e.g., 42 or 15.75), with no extra
1130             characters, units, or formatting. If it's a whole number, do not include decimal places unless
1131             required.",
1132             context=final_answer_raw
1133         )
1134
1135         # Step 7: Return the final answer
1136         return final_answer_clean.strip()

```

Listing 9: Extracted workflow code for SFT training data.

The above example demonstrates the complete four-stage pipeline from input prompt to validated workflow code, illustrating how each training pair is constructed and verified.

B.2 LORA CONFIGURATION

To mitigate the risk of mode collapse when fine-tuning on this limited dataset of approximately 1,300 examples, we employ Low-Rank Adaptation (LoRA) with rank-16 (Hu et al., 2022) and train for one epoch with batch size of 16, as illustrated in Figure 3. This parameter-efficient approach serves as an effective regularization mechanism: preliminary experiments with full-parameter fine-tuning resulted in degenerate repetition patterns (Listing 10), where the model generates circular, non-terminating reasoning loops instead of producing valid workflows.

```

1127 To solve this problem, we need to first understand what the problem is asking. The problem is asking us to
1128 find the answer. To find the answer, we need to solve this problem. To solve this problem, we need to first
1129 understand what the problem is asking. The problem is asking us to find the answer. To find the answer, we
1130 need to solve this problem. To solve this problem, we need to first understand what the problem is asking.
1131 The problem is asking us to find the answer. To find the answer, we need to solve this problem...

```

Listing 10: Example of degenerate repetition pattern (mode collapse) observed during full-parameter fine-tuning on the limited SFT dataset.

C DETAILS OF REINFORCEMENT LEARNING WITH VERIFIABLE REWARDS

The task set remains the same as the SFT stage:

$$\mathcal{D}_{\text{train}}^{\text{RLVR}} = \mathcal{D}_{\text{train}}^{\text{SFT}} = \{\text{GSM8K, DROP, MBPP, Humaneval}\}$$

To enhance the model’s generalization to diverse operator configurations, we train with four different operator sets that progressively introduce novel operators beyond the base SFT set:

$$\begin{aligned} \text{Ops}^{\text{RLVR},1} &= (\text{Generate, Summarize, Revise, Ensemble}) \\ \text{Ops}^{\text{RLVR},2} &= (\text{Generate, Summarize, Revise, Ensemble, Programmer}) \\ \text{Ops}^{\text{RLVR},3} &= (\text{Generate, Summarize, Revise, Ensemble, Decompose}) \\ \text{Ops}^{\text{RLVR},4} &= (\text{Generate, Summarize, Revise, Ensemble, Programmer, Decompose}) \end{aligned}$$

During each training iteration, we randomly sample one operator set from $\{\text{Ops}^{\text{RLVR},i}\}_{i=1}^4$ along with a task from $\mathcal{D}_{\text{train}}^{\text{RLVR}}$, forming diverse domain combinations (C, Ops) for meta-learning. This diversity ensures the planner learns operator-agnostic workflow construction principles rather than memorizing fixed operator-task associations.

D CASE STUDY: WORKFLOW DESIGN TRANSFORMATION AFTER TRAINING

D.1 OVERVIEW

To better understand the effectiveness of our training paradigm, we conduct a critical ablation study analyzing changes in model behavior before and after training. We compare two randomly selected workflows generated by the untrained base model (Qwen3-8B) versus the fully trained **MetaFlow** model after our two-stage training.

D.2 IMPACT OF TRAINING ON WORKFLOW GENERATION BEHAVIOR

Before Training (Base Model Qwen3-8B): Without SFT and RLVR optimization, the base model generates workflows with multi-step logic following a linear process: extract function name → summarize → generate → revise → ensemble.

```

class Workflow:
    # ... (initialization code omitted) ...
    async def run_workflow(self):
        # Step 1: Extract function name from test cases
        func_name = await self.generate(
            instruction="Extract the function name from the test cases.
                Return ONLY the function name, nothing else.",
            context=self.problem_text
        )
        # Step 2: Summarize problem description and test cases
        problem_summary = await self.summarize(
            instruction="Summarize the problem description and test cases
                to create a clear problem statement.",
            context=self.problem_text
        )
        # Step 3: Generate code based on problem summary
        generated_code = await self.generate(
            instruction=f"Write a Python function named '{func_name}'...
                Return ONLY the function definition, nothing else.",
            context=problem_summary
        )
        # Steps 4 & 5: Revise and ensemble
        revised_code = await self.revise(...)
        final_code = await self.ensemble(contexts=[generated_code, revised_code])
        return final_code

```

Listing 11: Workflow generated by the untrained base model.

Analysis: While structured, this workflow has potential weaknesses:

- Information Bottleneck:** The summarize step, intended to condense information, risks losing critical details or constraints from the original problem description.

- 1188
1189 **2. Instruction Ambiguity:** Vague instructions like "Return ONLY the function definition"
1190 can be misinterpreted, leading to incomplete code that lacks necessary import statements
1191 or helper functions.

1192 **After Training (MetaFlow Model):** In contrast, the **MetaFlow** model optimized through our two-
1193 stage training paradigm generates a workflow that, while maintaining complexity, is significantly
1194 more robust and methodical in its approach, and has learned the self-consistency operator unseen
1195 during training.

```

1196 % --- DO NOT IMPORT HERE ---
1197 class Workflow:
1198     def __init__(self, config, problem) -> None:
1199         # --- DO NOT MODIFY THIS SECTION ---
1200         self.config = config
1201         self.problem_text = problem
1202         self.llm = create(config)
1203
1204         self.generate = operator.Generate(self.llm, self.problem_text)
1205         self.revise = operator.Revise(self.llm, self.problem_text)
1206         self.summarize = operator.Summarize(self.llm, self.problem_text)
1207         self.ensemble = operator.Ensemble(self.llm, self.problem_text)
1208         self.selfconsistency = operator.SelfConsistency(self.llm, self.problem_text)
1209
1210     async def run_workflow(self):
1211         """
1212         Implement the core problem-solving logic here.
1213         Remember:
1214         - Use detailed, comprehensive instructions
1215         - Dynamic instruction construction is powerful
1216         - All operators expect (instruction: str, context: str) except Ensemble which takes contexts: List[
1217         str]
1218         """
1219         import asyncio
1220
1221         # Step 1: Extract function name from test cases
1222         function_name = await self.generate(
1223             instruction="""
1224             Analyze the test cases provided and extract the function name.
1225             The function name should match the naming convention used in the test cases.
1226             Return ONLY the function name, nothing else.
1227             """,
1228             context=self.problem_text
1229         )
1230
1231         # Step 2: Extract input types and expected outputs from test cases
1232         input_types_output = await self.generate(
1233             instruction="""
1234             Analyze the test cases provided and extract the input types and expected outputs.
1235             Format the response as a dictionary where keys are the parameter names and values are the types
1236             .
1237             Also, identify the expected output format.
1238             Return ONLY the structured information, nothing else.
1239             """,
1240             context=self.problem_text
1241         )
1242
1243         # Step 3: Generate multiple candidate implementations
1244         implementations = await asyncio.gather(
1245             self.generate(
1246                 instruction=f"""
1247                 Given the problem description: {self.problem_text}
1248                 Generate a Python function that satisfies the test cases.
1249                 Ensure that the function name is '{function_name}'.
1250                 Include all necessary imports at the beginning of the function.
1251                 Use the extracted input types and expected outputs: {input_types_output}.
1252                 """,
1253                 context=self.problem_text
1254             ),
1255             # ... (additional parallel generate calls omitted for brevity) ...
1256             self.generate(
1257                 instruction=f"""
1258                 Given the problem description: {self.problem_text}
1259                 Generate another Python function that satisfies the test cases.
1260                 Ensure that the function name is '{function_name}'.
1261                 Include all necessary imports at the beginning of the function.
1262                 Use the extracted input types and expected outputs: {input_types_output}.
1263                 """,
1264                 context=self.problem_text
1265             )
1266         )
1267
1268         # Step 4: Select the best implementation using SelfConsistency
1269         best_implementation = await self.selfconsistency(
1270             instruction=f"""
1271             Evaluate the following candidate implementations:
1272             {implementations[0]}
1273             {implementations[1]}
1274
1275             Select the best implementation based on the following criteria:

```

```

1242     - Consistency with the test cases
1243     - Logical correctness
1244     - Efficient and clean code
1245     - Proper use of imports and function signature
1246
1246     Return ONLY the best implementation, nothing else.
1247     """
1247     context=self.problem_text
1248 )
1248
1249     return best_implementation

```

Listing 12: Workflow generated by the trained **MetaFlow** model.

Analysis: The trained model’s workflow exhibits several key improvements contributing to its robustness:

1. **From Compression to Structured Extraction:** The model avoids the potentially lossy summarize step. Instead, it performs targeted extractions to parse unstructured text into structured data (e.g., function name, I/O types). This approach preserves information integrity and provides a more reliable foundation for subsequent steps.
2. **Parallel Exploration for Increased Robustness:** Rather than following a single generation-and-revision path, the model generates multiple, independent candidate solutions in parallel using `asyncio.gather`. This multi-path exploration inherently increases the robustness of the process by diversifying the solution space and reducing the risk of settling on a single, flawed implementation.
3. **More Explicit and Context-Aware Instructions:** The instructions are dynamically constructed using f-strings to include context from prior steps (e.g., ‘function_name’, ‘input_types_output’). This makes the prompts highly specific and unambiguous, directly mitigating the weaknesses of the base model and ensuring that each generative step is precisely guided.

E SELF CONSISTENCY OPERATOR: LOW-COST INTEGRATION EXAMPLE

This section demonstrates **MetaFlow**’s ability to immediately generate effective workflows with newly integrated operators without retraining. After introducing the `SelfConsistency` operator (which encapsulates parallel generation with majority voting), the model generates the following workflow for GSM8K mathematical reasoning through a single inference.

```

1276 # --- DO NOT IMPORT HERE ---
1276 class Workflow:
1277     def __init__(self, config, problem) -> None:
1277         # --- DO NOT MODIFY THIS SECTION ---
1278         self.config = config
1279         self.problem_text = problem
1279         self.llm = create(config)
1280
1280         self.generate = operator.Generate(self.llm, self.problem_text)
1281         self.revise = operator.Revise(self.llm, self.problem_text)
1282         self.summarize = operator.Summarize(self.llm, self.problem_text)
1283         self.ensemble = operator.Ensemble(self.llm, self.problem_text)
1284         self.self_consistency = operator.SelfConsistency(
1285             self.llm,
1286             self.problem_text,
1287             n_samples=5,
1288             similarity_threshold=0.9,
1289             return_full_info=False
1290         )
1291
1292     async def run_workflow(self):
1293         """
1294         Implement the core problem-solving logic here.
1295         Remember:
1296         - Use detailed, comprehensive instructions
1297         - Dynamic instruction construction is powerful
1298         - All operators expect (instruction: str, context: str) except Ensemble which takes contexts: List[str]
1299         """
1300         import asyncio
1301         sc_instruction = (
1302             "Solve this math problem step by step. "
1303             "Show all calculations clearly. "
1304             "At the end, state your final numerical answer."
1305         )

```

```

1296     sc_answer = await self.self_consistency(
1297         instruction=sc_instruction,
1298         context="",
1299         answer_type="numeric"
1300     )
1301     format_instruction = (
1302         f"The calculated answer is: {sc_answer}\n\n"
1303         "Extract ONLY the final numerical answer. "
1304         "If it's a decimal, round to a reasonable precision. "
1305         "If the answer represents a count of items/people, it should be a whole number. "
1306         "Return ONLY the number, nothing else."
1307     )
1308     final_answer = await self.generate(
1309         instruction=format_instruction,
1310         context=sc_answer
1311     )
1312     return final_answer.strip()

```

Listing 13: Example workflow generated by **MetaFlow** for GSM8K after integrating the SelfConsistency operator. The workflow demonstrates immediate operator utilization: (1) invoking SelfConsistency with $n=5$ parallel samples and similarity-based voting, and (2) formatting the numeric answer. This workflow was generated without any model retraining, showcasing low-cost operator extensibility.

F CASE STUDY: ZERO-SHOT WORKFLOW GENERATION WITH NOVEL VECTORSEARCH OPERATOR

This case study demonstrates **MetaFlow**'s zero-shot generalization to the novel VectorSearch operator on HotpotQA multi-hop question answering. The operator was never encountered during SFT or RLVR training, yet the model successfully generates workflows that orchestrate complex retrieval-reasoning pipelines. Listing 14 shows the highest-scoring workflow (0.74 accuracy) among 100 generated candidates, which achieves 60.9% improvement over the CoT+RAG baseline (0.46 accuracy).

```

1325 # --- DO NOT IMPORT HERE ---
1326 class Workflow:
1327     def __init__(self, config, problem) -> None:
1328         # --- DO NOT MODIFY THIS SECTION ---
1329         self.config = config
1330         self.problem_text = problem
1331         self.llm = create(config)
1332
1333         self.generate = operator.Generate(self.llm, self.problem_text)
1334         self.revise = operator.Revise(self.llm, self.problem_text)
1335         self.summarize = operator.Summarize(self.llm, self.problem_text)
1336         self.ensemble = operator.Ensemble(self.llm, self.problem_text)
1337         self.vector_search = operator.VectorSearch(self.llm, self.problem_text)
1338
1339     async def run_workflow(self):
1340         """
1341         Implement the core problem-solving logic here.
1342         Remember:
1343         - Use detailed, comprehensive instructions
1344         - Dynamic instruction construction is powerful
1345         - All operators expect (instruction: str, context: str) except Ensemble which takes contexts: List[str]
1346         """
1347         import asyncio
1348
1349         # Step 1: Extract key entities from the question
1350         extraction_instruction = (
1351             "Identify all named entities, key terms, and concepts in the question. "
1352             "Focus on entities that are likely to appear in the documents and are critical for answering the question."
1353         )
1354         question_entities = await self.generate(instruction=extraction_instruction, context=self.problem_text)
1355
1356         # Step 2: First Vector Search to retrieve initial relevant documents
1357         first_search_instruction = (
1358             f"Use the extracted entities: {question_entities}. "
1359             "Perform a vector search to find documents that mention these entities. "
1360             "Prioritize documents that provide foundational information about the question."
1361         )
1362         first_retrieved_docs = await self.vector_search(instruction=first_search_instruction, context=self.problem_text, top_k=5)
1363
1364         # Step 3: Analyze first retrieved documents to identify connections

```

```

1350     analysis_instruction = (
1351         f"Based on the retrieved documents: {first_retrieved_docs}. "
1352         "Identify any shared entities, relationships, or connections between the documents. "
1353         "Highlight how these connections can help answer the question."
1354     )
1355     connections_analysis = await self.generate(instruction=analysis_instruction, context=self.
1356         problem_text)
1357
1358     # Step 4: Second Vector Search to retrieve documents connecting the entities
1359     second_search_instruction = (
1360         f"Use the connections identified: {connections_analysis}. "
1361         "Perform a vector search to find documents that explicitly connect the entities or provide
1362         additional context."
1363     )
1364     second_retrieved_docs = await self.vector_search(instruction=second_search_instruction, context=
1365         self.problem_text, top_k=5)
1366
1367     # Step 5: Synthesize information from both document sets
1368     synthesis_instruction = (
1369         f"Combine the information from the first set: {first_retrieved_docs}, "
1370         f"and the second set: {second_retrieved_docs}. "
1371         "Extract the exact answer to the question based on the synthesized information. "
1372         "Ensure the answer is concise and factually accurate."
1373     )
1374     final_answer = await self.generate(instruction=synthesis_instruction, context=self.problem_text)
1375
1376     return final_answer

```

Listing 14: Highest-scoring workflow for HotpotQA using the novel VectorSearch operator. The workflow performs iterative multi-hop reasoning: (1) entity extraction, (2) first-hop document retrieval, (3) connection analysis, (4) second-hop retrieval, and (5) answer synthesis. This demonstrates **MetaFlow**'s ability to compose complex operator sequences for unseen tools through learned meta-strategies.

ANALYSIS AND PERFORMANCE

the successful integration and application of a new tool. In Step 2, the model dynamically constructs a search query from its initial analysis and invokes the VECTORSEARCH operator, effectively performing active information retrieval. When evaluated on the HotpotQA downstream task, this approach achieved a **60% search accuracy**. This result is significant as it confirms that our operator framework enables models to learn and effectively utilize unseen tools.

NOTE ON COMPARABILITY WITH BASELINES

It is crucial to highlight a fundamental difference between our evaluation and that of many previous works on HotpotQA. Our methodology requires the model to **actively perform a search** to find relevant information. In contrast, prior baselines are often provided with the ground-truth supporting documents as part of their input, thereby bypassing the challenging information retrieval step entirely. Because our system solves a more complete and realistic version of the task that includes an explicit search phase, a direct comparison of end-to-end accuracy with such baselines is not meaningful.