

# FUSIONFLOW: Enabling Deep Structural Exploration for Automated Agentic Workflow Generation

Anonymous ACL submission

## Abstract

Agentic workflows are commonly used to guide large language models in solving complex reasoning tasks. However, existing automated workflow generation methods primarily rely on stepwise local refinement or tree-based search over a single evolving workflow. Under limited optimization budgets, this paradigm constrains **structural depth**, hindering the discovery of workflows that require deep compositional structure. To address this limitation, we propose **FUSIONFLOW**, a framework centered on workflow fusion. Unlike incremental refinement, fusion enables structural leaps by synthesizing multiple independently evolved workflows, allowing exploration of deeper regions of the workflow space within a finite budget. To make fusion effective, **FUSIONFLOW** integrates local optimization, task-specific differentiation, and a dynamic scheduling mechanism. Experiments on six reasoning benchmarks demonstrate that **FUSIONFLOW** consistently outperforms existing automated workflow generation methods. Further ablation and analysis confirm that fusion is the key driver of deep structural exploration, highlighting fusion-driven exploration as an effective approach for overcoming depth limitations in automated workflow generation.

## 1 Introduction

Large language models (LLMs) (Touvron et al., 2023; Achiam et al., 2023; Anthropic, 2025b) have demonstrated strong capabilities in reasoning, planning, and decision-making (Wang et al., 2024a; Zaharia et al., 2024; Xi et al., 2025), giving rise to LLM-based agent systems that autonomously tackle complex tasks across diverse domains with minimal human intervention (Hong et al., 2023; Qian et al., 2024a; Ridnik et al., 2024). Early agent systems emphasize free-form reasoning, enabling autonomous planning and dynamic interaction with environments (Chen et al., 2023; Qiao et al., 2024;

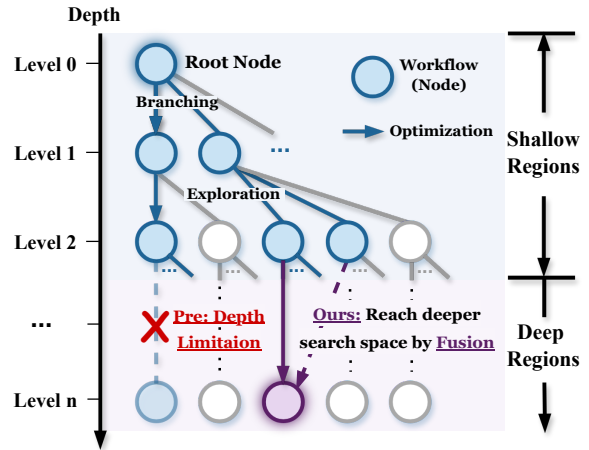


Figure 1: Illustration of **structural depth** limitation under a finite optimization budget. Incremental optimization methods are restricted to a shallow **structural depth** proportional to the iteration count. In contrast, fusion-based exploration enables structural leaps that allow the search to reach deeper regions of the workflow space within the same budget.

Zhu et al., 2024). Despite their adaptability and creativity, such systems often suffer from instability and limited controllability. To address these limitations, recent work has introduced *agentic workflows*, which structure reasoning into repeatable, explicitly defined pipelines of LLM calls (Khatab et al., 2024; Zhou et al., 2024; Yuksekgonul et al., 2024). By constraining decision-making within well-defined workflows, these approaches provide more stable and controllable foundations for complex problem solving.

Existing automated workflow generation frameworks explore large combinatorial spaces of code- or graph-based workflows and iteratively refine candidate workflows based on execution feedback (Zhuge et al., 2024; Hu et al., 2024; Liu et al., 2025), shifting the role of LLMs from solving individual tasks to constructing structured problem-solving procedures. Representative methods such as ADAS (Hu et al., 2024) and AFlow (Zhang et al., 2024) operationalize this paradigm through heuris-

064 tic optimization and tree-based exploration. Under  
065 practical optimization budgets, however, such in-  
066 cremental local refinement often constrains explo-  
067 ration to limited **structural depth**, i.e., the number  
068 of effective structural modifications accumulated  
069 along an optimization path, making it difficult to  
070 discover workflows with sufficiently complex struc-  
071 tures, as illustrated in Figure 1 and further analyzed  
072 in Appendix D.

073 To overcome this structural depth limitation, we  
074 introduce the **Fusion Operator**. Unlike incremen-  
075 tal local refinement methods that optimize a single  
076 workflow through step-wise modifications, fusion  
077 synthesizes effective components from multiple  
078 candidate workflows, enabling *structural leaps* dur-  
079 ing search. However, fusion alone is insufficient for  
080 effective structural exploration. Its success depends  
081 on the availability of diverse and complementary  
082 workflow structures to integrate. To this end, we  
083 leverage empirical observations about the organi-  
084 zation of complex reasoning tasks. An analysis  
085 of the MATH benchmark (Figure 2A) reveals that  
086 such tasks often exhibit clear sub-domain struc-  
087 ture, with distinct reasoning patterns across prob-  
088 lem categories. This motivates a search strategy in  
089 which workflows are first specialized for different  
090 sub-domains and subsequently fused to construct a  
091 more robust global solution. Guided by this obser-  
092 vation, we introduce the **Differentiation Operator**,  
093 which explicitly encourages the emergence of di-  
094 verse, specialized workflow variants that serve as  
095 complementary building blocks for effective fu-  
096 sion.

097 Based on these insights, we propose **FUSION-**  
098 **FLOW**, a workflow generation framework cen-  
099 tered on fusion-driven structural exploration. The  
100 core of this framework is the **Fusion Operator**,  
101 which serves as the primary mechanism for break-  
102 ing **structural depth** limitations, while other com-  
103 ponents are designed to support its effectiveness.  
104 Specifically, we incorporate an **Optimization Op-**  
105 **erator** to perform stable local refinement and a **Dif-**  
106 **ferentiation Operator** to generate diverse, com-  
107plementary workflows for fusion. Coordinating  
108 these operators introduces a challenge, as they in-  
109volve different trade-offs between exploitation and  
110 exploration. To address this issue, we introduce a  
111 **Dynamic Scheduling Algorithm** as the core con-  
112 trol mechanism of **FUSIONFLOW**. Modeled as  
113 a non-stationary Markov process, this algorithm  
114 adaptively selects operators, balancing stable local  
115 improvement with high-variance structural evolu-

tion.

116 Experiments on six reasoning benchmarks  
117 demonstrate that **FUSIONFLOW** consistently out-  
118 performs existing automated workflow generation  
119 methods. Further ablation and analysis confirm  
120 that fusion is the key driver of deep structural ex-  
121 ploration, highlighting fusion-driven exploration as  
122 an effective approach for overcoming depth limita-  
123 tions in automated workflow generation.

124 Our contributions are summarized as follows: 125

- We characterize the **structural depth** limi- 126  
127 tation of incremental workflow optimization  
128 methods and introduce the **Fusion Operator**  
129 to enable structural exploration beyond single-  
130 step refinement.
- We present **FUSIONFLOW**, a fusion-centered 131  
132 workflow generation framework in which  
133 optimization, differentiation, and dynamic  
134 scheduling are designed to support effective  
135 fusion.
- Experiments on multiple benchmarks validate 136  
137 that fusion is the key driver for overcoming  
138 depth limitations and achieving superior per-  
139 formance.

## 2 Related Work 140

### 2.1 LLM-based Autonomous Agent Planning 141

142 Early agent systems were mainly built for specific  
143 domains such as code generation (Hong et al., 2023;  
144 Qian et al., 2024a; Ridnik et al., 2024), question  
145 answering (Zhou et al., 2022, 2023; Xu et al., 2024),  
146 and robotics (Liang et al., 2022; Singh et al., 2022;  
147 Song et al., 2023). These systems followed fixed  
148 pipelines or rule-based designs for task planning  
149 and reasoning (Wang et al., 2024a; Xi et al., 2025).  
150 As large language models (LLMs) became more ca-  
151 pable (Achiam et al., 2023; Anthropic, 2025b), re-  
152 search shifted toward agents that can plan and make  
153 decisions more flexibly (Chen et al., 2023; Zhuge  
154 et al., 2023; Huang et al., 2024; Qiao et al., 2024;  
155 Zhu et al., 2024; Liu et al., 2025). Such LLM-based  
156 agents can decompose and solve problems through  
157 natural language reasoning, showing strong flexibil-  
158 ity and creativity. However, their decision process  
159 is not explicitly structured, so their behavior and  
160 results are often unpredictable before execution.  
161 This makes it difficult to optimize them for specific  
162 tasks or to ensure consistent performance. Current  
163 improvements mainly focus on higher-level strate-  
164 gies, such as using multiple agents to cooperate

(Guo et al., 2024; Wang et al., 2024a; Qian et al., 2024b) or adding external feedback loops (Shinn et al., 2023; Wang et al., 2024b,c). While these methods improve stability to some extent, they still lack precise control over the reasoning process.

## 2.2 Automated Agentic Workflow Optimization

Agentic workflows provide a structured alternative to fully autonomous agents by organizing multiple LLM calls into explicit, executable reasoning pipelines (Zhuge et al., 2024; Hu et al., 2024; Liu et al., 2025). By constraining decision making within predefined workflow structures, these approaches improve controllability and reproducibility compared to free-form agent behaviors.

Early work primarily focused on improving performance within fixed workflows through prompt or parameter optimization, without modifying the underlying structure (Chen et al., 2023; Xu et al., 2023; Yang et al., 2023; Fernando et al., 2023; Khattab et al., 2024; Zhou et al., 2024; Yuksekogonul et al., 2024). More recent methods extend this line of research by directly searching over workflow structures, representing workflows as graphs or executable programs and refining them through iterative optimization (Li et al., 2024; Cheng et al., 2024; Zhuge et al., 2024; Hu et al., 2024; Zhang et al., 2024). Despite differences in representation and search strategy, most existing approaches incrementally refine a single candidate workflow through local modifications, limiting achievable search depth under practical budgets. Our work departs from this paradigm by enabling structural exploration beyond single-path refinement, addressing the depth limitation of prior methods.

## 3 Method

### 3.1 Preliminary

**Agentic Workflow.** Formally, an agentic workflow  $\mathcal{W}$  is defined as an executable program modeled by a directed graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ , where  $\mathcal{N}$  represents the set of functional nodes and  $\mathcal{E}$  represents the control flow logic connecting them.

Each node  $n_i \in \mathcal{N}$  encapsulates a specific atomic operation performed by a LLM. A node is parameterized as  $n_i = (M, P, \tau, \mathcal{F})$ , where  $M$ ,  $P$ ,  $\tau$  and  $\mathcal{F}$  denotes the backbone model, the natural language prompt instruction, the sampling temperature, and the output format constraints respectively.

The edge set  $\mathcal{E}$  is implemented via programming language constructs, supporting sequential execution, conditional branching, loops, and recursive calls. We provide representative code examples of the workflow graph structure, operator interfaces, and prompt organization used in **FUSIONFLOW** in Appendix G.

Given an input task  $x$ , the workflow executes by traversing the directed graph according to the control-flow edges, invokes a sequence of LLM-based nodes, and finally returns a single output as the final answer. Mathematically, a workflow  $\mathcal{W}$  acts as a mapping function  $f_{\mathcal{W}} : \mathcal{X} \rightarrow \mathcal{Y}$ , which takes an input task  $x \in \mathcal{X}$  and produces a final answer  $y \in \mathcal{Y}$ , where  $\mathcal{X}$  denotes the input space of task instances and  $\mathcal{Y}$  denotes the output space of answers.

**Problem Definition.** We formulate the automated generation of agentic workflows as a search problem within the infinite workflow space  $\mathbb{S}$ . Given a task dataset  $\mathcal{D} = \{(x_j, y_j)\}_{j=1}^{|\mathcal{D}|}$  and an evaluation metric  $G(\hat{y}, y)$ , for each task instance  $(x, y)$  sampled from the dataset, the workflow produces a prediction  $\hat{y} = f_{\mathcal{W}}(x)$ , which is then compared against the ground-truth label  $y$  via the metric  $G(\hat{y}, y)$  to obtain a scalar score. The overarching objective is to identify an optimal workflow  $\mathcal{W}^*$  that maximizes the expected performance over the data distribution:

$$\mathcal{W}^* = \arg \max_{\mathcal{W} \in \mathbb{S}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [G(\hat{y}, y)] \quad (1)$$

Here,  $\mathbb{E}_{(x,y) \sim \mathcal{D}} [\cdot]$  denotes the empirical expectation over the dataset.

**Structural Depth.** We define the **structural depth**  $d(\mathcal{W})$  of a workflow as the number of effective modifications it accumulates during the search process. Intuitively, each optimization step adds one modification, differentiation preserves the parent’s depth, and fusion aggregates modifications from multiple parent workflows. Formally, let  $\mathcal{N}(\mathcal{W})$  denote the set of modification nodes in  $\mathcal{W}$ , then  $d(\mathcal{W}) = |\mathcal{N}(\mathcal{W})|$ . Under purely incremental optimization, the maximum achievable depth is bounded by the length of a single optimization chain. In contrast, fusion enables  $d(\mathcal{W}_{\text{fuse}}) \geq \max_i d(\mathcal{W}_i)$  by merging modification histories from multiple trajectories, thereby allowing workflows to reach deeper structural regions within a finite iteration budget. The formal recur-

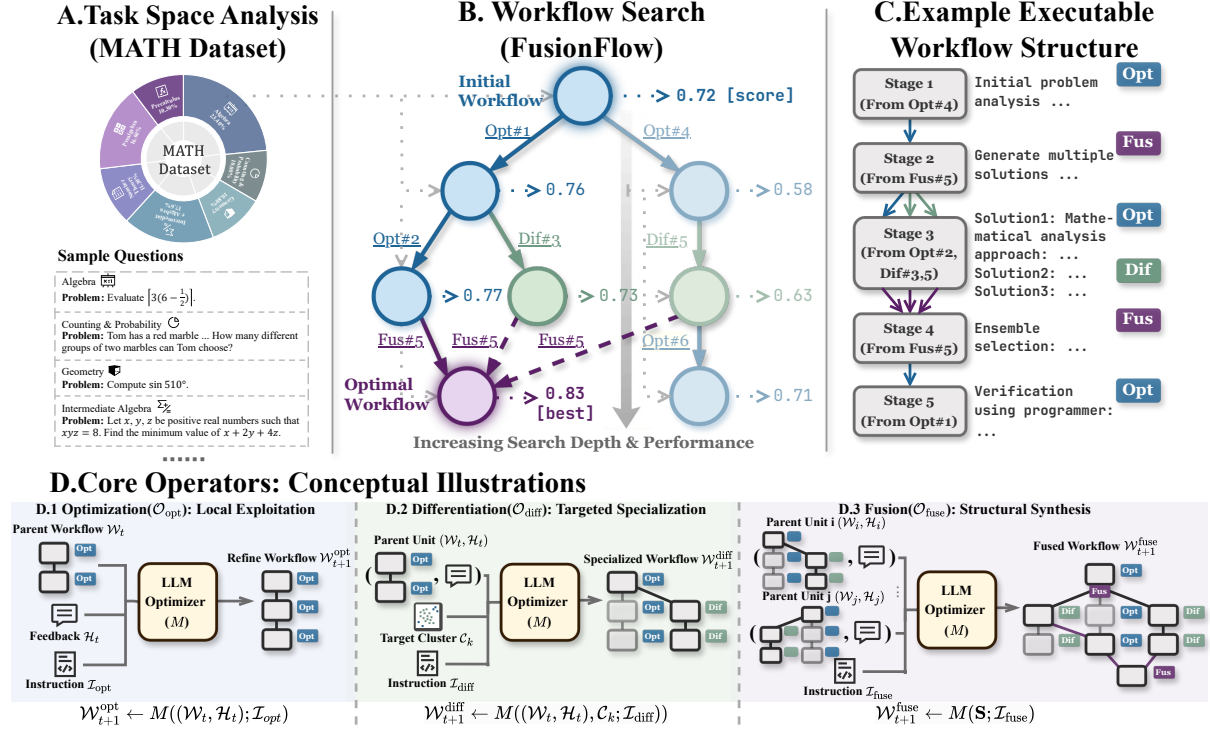


Figure 2: Overall framework of **FUSIONFLOW**. (A) The task space exhibits a clustered structure (e.g., Algebra, Goemetry in MATH Dataset), motivating targeted search strategies. (B) The evolutionary search tree where nodes represent workflows and edges denote operations. The illustration shows how the framework progressively improves by leveraging information accumulated from previously explored workflows through Optimization, Differentiation, and Fusion, evolving from a baseline (0.72) to the optimal workflow (0.83). (C) An example structure of an executable workflow. (D) Conceptual illustrations of the three core operators.

sive definition and empirical analysis are provided in Appendix D.1.

### 3.2 Overview of FUSIONFLOW

Figure 2 illustrates the overall design of **FUSIONFLOW**. The framework operates over a search tree  $\mathcal{T} = (\mathcal{W}, \mathcal{O})$  where each node  $\mathcal{W}$  represents a fully executable workflow and edges  $\mathcal{O}$  correspond to transformation operators. At each iteration  $t$ , the framework selects an operator  $\mathcal{O}_t \in \{\mathcal{O}_{\text{opt}}, \mathcal{O}_{\text{diff}}, \mathcal{O}_{\text{fuse}}\}$  and one or more parent workflows from the previously explored nodes, then applies the operator to generate a new workflow node, thereby expanding the search tree.

In our framework, agentic workflows are composed of atomic operators connected by code control logic, which allows a large language model to synthesize a composite workflow by identifying shared structures and reconciling complementary components across multiple workflows. Fusion itself is therefore not difficult to apply (as demonstrated in Appendix B.5). The main challenge, however, lies in determining when fusion should be applied and which workflows should be fused, as indiscriminate fusion can result in redun-

dancy, incompatible coordination, or performance degradation.

To address these challenges, **FUSIONFLOW** incorporates a set of auxiliary components that jointly ensure the effectiveness of fusion:

- **Optimization** refines individual workflows to produce candidates with reliable local behaviors.
- **Differentiation** intentionally introduces structural and functional diversity, creating complementary candidate workflows.
- **Dynamic Scheduling Algorithm** coordinates these components with the fusion operator, adaptively regulating fusion timing and candidate selection to maximize depth expansion under limited iteration budgets.

### 3.3 Fusion Operator

We formalize workflow fusion as a structure level composition operator  $\mathcal{O}_{\text{fuse}}$  that synthesizes multiple agentic workflows into a single complex workflow. Let  $\mathbf{S} = \{(\mathcal{W}_i, \mathcal{H}_i)\}_{i=1}^K$  denote a selected set of  $K$  parent units ( $K \geq 2$ ), where  $\mathcal{H}_i$  denotes

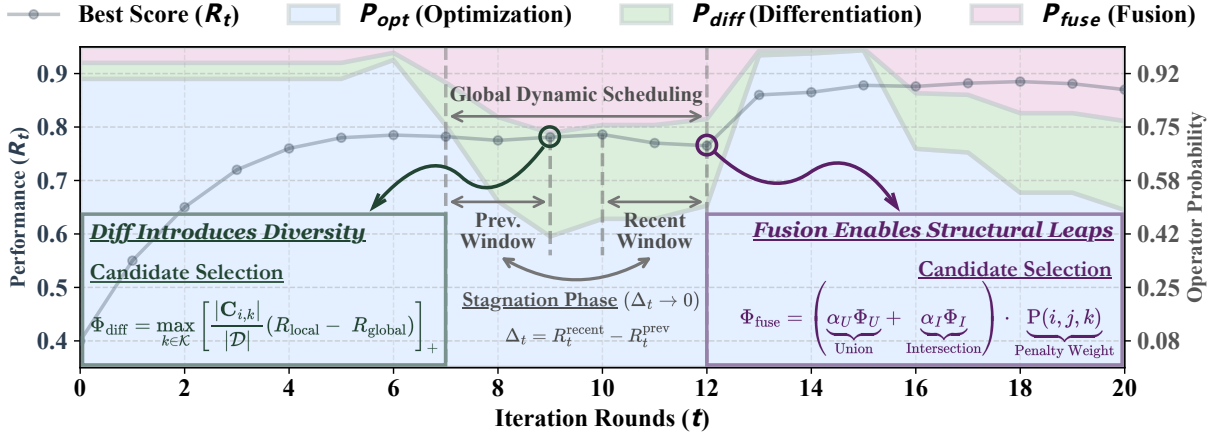


Figure 3: Conceptual illustration of the dynamic scheduling process in **FUSIONFLOW** (real quantitative results are reported in Appendix F). The curve shows the evolution of the best workflow performance  $R_t$  over iteration rounds, while shaded regions indicate the probabilistic allocation among optimization ( $p_{\text{opt}}$ ), differentiation ( $p_{\text{diff}}$ ), and fusion ( $p_{\text{fuse}}$ ). The middle panel highlights the two sliding windows used to detect stagnation. The formula in the lower-left corner defines the differentiation potential. The formula in the lower-right corner defines the fusion potential. Circled iterations indicate rounds where differentiation and fusion are applied.

execution feedback of  $\mathcal{W}_i$  on the task dataset. Formally, the fusion process is defined as:

$$\mathcal{W}_{t+1}^{\text{fuse}} \leftarrow M(\mathcal{S}; \mathcal{I}_{\text{fuse}}) \quad (2)$$

Unlike optimization, fusion integrates the internal structures of multiple workflows. Concretely, each parent workflow  $\mathcal{W}_i$  is represented as a directed graph  $\mathcal{G}_i = (\mathcal{N}_i, \mathcal{E}_i)$ . The fused workflow  $\mathcal{W}_{t+1}^{\text{fuse}}$  integrates selected subgraphs from  $\mathcal{G}_i$  through newly constructed control-flow nodes. As a result,  $\mathcal{W}_{t+1}^{\text{fuse}}$  preserves the functional behaviors of its parents while introducing additional coordination logic. The core logic of the fusion operator is summarized in Appendix H.2.

### 3.4 Support Components

#### 3.4.1 Optimization Operator

The optimization operator  $\mathcal{O}_{\text{opt}}$  performs local refinement on a single workflow to improve reliability and correctness. Formally, it is defined as:

$$\mathcal{W}_{t+1}^{\text{opt}} \leftarrow M((\mathcal{W}_t, \mathcal{H}_t); \mathcal{I}_{\text{opt}}) \quad (3)$$

where  $\mathcal{I}_{\text{opt}}$  specifies optimization instructions that constrain the modification to prompt content and node parameters. Within **FUSIONFLOW**, optimization primarily serves as a preparatory operator that stabilizes local behaviors and ensures that workflows are sufficiently mature before being considered as reliable building blocks for fusion. Appendix H.1 provides the main logic used to produce and validate a single local refinement.

#### 3.4.2 Differentiation Operator

The differentiation operator  $\mathcal{O}_{\text{diff}}$  generates specialized variants of a workflow with respect to different task subspaces. Given a workflow  $\mathcal{W}_t$ , differentiation produces a specialized workflow:

$$\mathcal{W}_{t+1}^{\text{diff}} \leftarrow M((\mathcal{W}_t, \mathcal{H}_t), \mathcal{C}_k; \mathcal{I}_{\text{diff}}) \quad (4)$$

where  $\mathcal{C}_k \subset \mathcal{D}$  represents a task sub-cluster and  $\mathcal{I}_{\text{diff}}$  guides structural or reasoning-level modifications conditioned on the characteristics of  $\mathcal{C}_k$ . Differentiation increases functional and structural diversity in the pool, aiming to produce complementary workflows that make subsequent fusion more likely to yield non-redundant gains. We summarize the differentiation main execution logic in Appendix H.3.

#### 3.4.3 Dynamic Scheduling Algorithm

We model the search process as a partially observable Markov decision process which is formally illustrated in Appendix B.7. To coordinate the overall search process, we introduce a dynamic scheduling mechanism (as shown in Figure 3) that regulates the application of all operators in the framework. Since workflow fusion constitutes the core contributor to depth expansion in **FUSIONFLOW**, we focus our presentation here on fusion scheduling, including when fusion is applied and how fusion candidates are selected. The complete scheduling procedure over all operators is provided in the Appendix C and Appendix E.

**Probabilistic Fusion Scheduling** Within the Markov decision framework, the scheduler must

367 decide when to transition from exploitation to ex- 412  
 368 ploration. We use recent performance dynamics 413  
 369 as an observable proxy for the latent advantage of 414  
 370 each operator. Let  $R_t$  denote the best performance 415  
 371 achieved at iteration  $t$ . Over a sliding window of 416  
 372 size  $k$ , we define 417

$$373 \begin{aligned} R_t^{\text{recent}} &= \max_{u \in [t-k+1, t]} R_u, \\ R_t^{\text{prev}} &= \max_{u \in [t-2k+1, t-k]} R_u. \end{aligned} \quad (5)$$

374 Using the window maximum provides a robust 422  
 375 progress indicator that tracks upper-bound improve- 423  
 376 ment while filtering evaluation noise. The improve- 424  
 377 ment difference  $\Delta_t = R_t^{\text{recent}} - R_t^{\text{prev}}$  serves as an 425  
 378 empirical estimator for the marginal return of con- 426  
 379 tinued local refinement. Based on this quantity, we 427  
 380 define a fusion activation signal 428

$$381 \pi_t = \frac{1}{1 + \exp[\kappa \cdot \Delta_t]} \quad (6)$$

382 where  $\kappa$  controls sensitivity to performance stagna- 432  
 383 tion. A higher  $\pi_t$  indicates reduced marginal gains 433  
 384 from local refinement, under which the scheduler 434  
 385 assigns higher probability to selecting fusion (as 435  
 386 illustrated in Appendix B.7 and Appendix C).

387 **Fusion Candidate Selection** Once fusion is acti- 436  
 388 vated, the scheduler must select which workflows 437  
 389 to fuse. This corresponds to choosing the action 438  
 390 operands within the Markov decision process. Can- 439  
 391 didate selection is based on structural complemen- 440  
 392 tarity and behavioral consensus, which approxi- 441  
 393 mate the latent compatibility between workflows 442  
 394 that determines fusion success.

395 Concretely, the scheduler first filters a small 443  
 396 subset  $\mathcal{P}' \subset \mathcal{P}_t$  of good performance workflows. 444  
 397 We then evaluate candidate workflow groups  $\mathcal{S} =$  445  
 398  $\{\mathcal{W}_1, \dots, \mathcal{W}_m\} \subset \mathcal{P}'$  using a fusion potential that 446  
 399 jointly captures complementary coverage and be- 447  
 400 havioral consensus. For each workflow  $\mathcal{W}_i$ , let 448  
 401  $\mathcal{C}_i \subseteq \mathcal{D}$  denote the subset of task instances in the 449  
 402 dataset that  $\mathcal{W}_i$  successfully solves. Complemen- 450  
 403 tarity and consensus are measured by both pairwise 451  
 404 set operations within the group and group-level ag- 452  
 405 gregation, and the group with the highest combined 453  
 406 score  $\Phi_{\text{fuse}}$  is selected for fusion (as illustrated in 454  
 407 Appendix C.4). 455

## 408 4 Experiments

### 409 4.1 Experimental Setup

410 **Datasets.** Following the convention of previous 459  
 411 work on automated workflow generation (Hu et al., 460

2024; Zhang et al., 2024), we evaluate our method 412  
 on six public benchmarks, including HotpotQA 413  
 (Yang et al., 2018), DROP (Dua et al., 2019), 414  
 HumanEval (Chen, 2021), MBPP (Austin et al., 415  
 2021), GSM8K (Cobbe et al., 2021), and MATH 416  
 (Hendrycks et al., 2021). As is common practice, 417  
 these methods require a small annotated set for 418  
 workflow optimization; therefore, we split each 419  
 dataset into a validation set and a test set at a 1:4 420  
 ratio, using the former for optimization. Further 421  
 details on the datasets and their splits are provided 422  
 in Appendix A.1. 423

**Baselines.** We compare our method against two 424  
 types of baselines. First, we consider a range 425  
 of manually designed prompting methods, includ- 426  
 ing standard IO (Direct LLM invocation), Chain- 427  
 of-Thought (Wei et al., 2022), Self-Consistency 428  
 with CoT (5 answers) (Wang et al., 2022), Multi- 429  
 Persona Debate (Wang et al., 2024d), Self-Refine 430  
 (max 3 iterations) (Madaan et al., 2023), and Med- 431  
 Prompt (Nori et al., 2023). Second, we compare 432  
 against automated workflow optimization methods, 433  
 specifically the recent approach such as ADAS (Hu 434  
 et al., 2024) and AFlow (Zhang et al., 2024). 435

**Inference Details.** We use different models for 436  
 the optimization and execution phases. Specifi- 437  
 cally, we employ Claude-3-7-Sonnet (Anthropic, 438  
 2025a) as the optimizer and GPT-4o-mini (OpenAI, 439  
 2024) as the executor. All models are accessed via 440  
 APIs with a temperature of 0. Further details are 441  
 provided in Appendix A.2. 442

### 443 4.2 Main Results

Table 1 summarizes the performance of all methods 444  
 across six benchmarks. API costs are reported 445  
 separately in Appendix A.3. 446

**Comparison with Manual Prompting.** Over- 447  
 all, FUSIONFLOW outperforms most manually 448  
 designed prompting strategies and achieves the 449  
 highest average accuracy across benchmarks. No- 450  
 tably, it yields substantial improvements on MATH, 451  
 DROP, and MBPP. For example, on MATH, FU- 452  
 SIONFLOW surpasses MedPrompt by +5.1 points 453  
 (67.5 vs. 62.4), indicating that automatically dis- 454  
 covered workflows can outperform static prompt 455  
 designs under identical model settings. 456

**Comparison with Automated Optimization.** 457  
 We further compare FUSIONFLOW with auto- 458  
 mated workflow optimization baselines. Across all 459  
 benchmarks, FUSIONFLOW consistently achieves 460

Method	Benchmarks						Avg.
	HotpotQA	DROP	HumanEval	MBPP	GSM8K	MATH	
<b>Manual Prompting Methods</b>							
IO (Direct)	69.1	78.3	87.0	72.1	91.5	56.4	75.7
CoT (Wei et al., 2022)	66.6	79.4	91.6	71.6	90.3	59.3	76.5
CoT SC (5-shot) (Wang et al., 2022)	63.8	77.0	87.0	73.9	<b>92.5</b>	62.1	76.1
MultiPersona (Wang et al., 2024d)	68.8	78.5	90.8	69.5	91.4	62.6	76.9
Self-Refine (Madaan et al., 2023)	64.3	75.0	<b>92.4</b>	73.9	85.5	56.4	74.6
MedPrompt (Nori et al., 2023)	70.6	78.9	90.1	74.2	91.2	62.4	77.9
<b>Automated Workflow Optimization</b>							
ADAS (Hu et al., 2024)	74.3	75.6	85.5	70.1	88.4	52.9	74.5
AFlow (Zhang et al., 2024)	68.8	83.3	89.3	81.3	91.2	62.6	79.4
<b>FUSIONFLOW(ours)</b>	<b>77.0</b>	<b>84.9</b>	90.8	<b>84.5</b>	90.8	<b>67.5</b>	<b>82.6</b>

Table 1: Performance comparison on six public benchmarks. All methods utilize GPT-4o-mini as the executor. We report accuracy (%) for all tasks. Best results are highlighted in bold.

Variant	Accuracy(%)
IO	56.4
CoT	59.3
<b>w/ FUSIONFLOW</b>	
Optimization Only	62.6 (↑ 6.2)
+ Differentiation	61.7 (↑ 5.3)
+ Fusion	64.8 (↑ 8.4)
+ Diff. & Fusion	67.5 (↑ 11.1)

Table 2: Ablation study of **FUSIONFLOW** components on the MATH dataset, using GPT-4o-mini as the executor model.

stronger performance, outperforming AFlow by an average margin of +3.2 points. While AFlow remains competitive on several benchmarks, its overall performance is limited under the same optimization budget. In contrast, **FUSIONFLOW** demonstrates more robust gains, suggesting the effectiveness of its workflow generation strategy beyond incremental refinement.

**Summary of Observations.** Taken together, the results in Table 1 show that **FUSIONFLOW** delivers consistent improvements over both manually designed prompting methods and existing automated optimization approaches. These gains motivate a deeper analysis of how different optimization strategies explore the workflow space, which we examine in detail through optimization dynamics (Section 4.5) and structural depth analyses (Appendix D) in subsequent sections.

### 4.3 Ablation Study

To investigate the contribution of each component in **FUSIONFLOW**, we conduct an ablation study on the MATH dataset. We define four variants: (1) **Op-**

Backbone	Method	Accuracy (%)
GPT-4o-mini	IO	56.4
	AFlow	62.6
	<b>Ours</b>	<b>67.5</b>
GPT-4o	IO	63.6
	AFlow	68.1
	<b>Ours</b>	<b>74.7</b>
Qwen-2.5-14b	IO	60.5
	AFlow	64.4
	<b>Ours</b>	<b>68.1</b>
Llama-3-8b	IO	10.3
	AFlow	28.8
	<b>Ours</b>	<b>29.8</b>

Table 3: **Performance Comparison across Different Models.** We compare our method with IO and AFlow using four different foundation models.

**timization Only:** which degenerates to a standard iterative search similar to AFlow; (2) **+ Differentiation:** which enables Differentiation Operator; (3) **+ Fusion:** which enables Fusion Operator and (4) **+ Diff. & Fusion:** which activates both operators and represents the full version of **FUSIONFLOW**.

As presented in Table 2, optimization alone quickly saturates, confirming the limited depth reachable by purely local search. Differentiation yields modest gains by increasing exploration diversity, but improvements remain constrained.

In contrast, enabling fusion leads to a substantial performance gain (+8.4 points), demonstrating that fusion is the primary driver for escaping local optima. Combining differentiation and fusion achieves the best performance, as diverse workflow variants provide complementary structures for effective fusion. These results directly support our

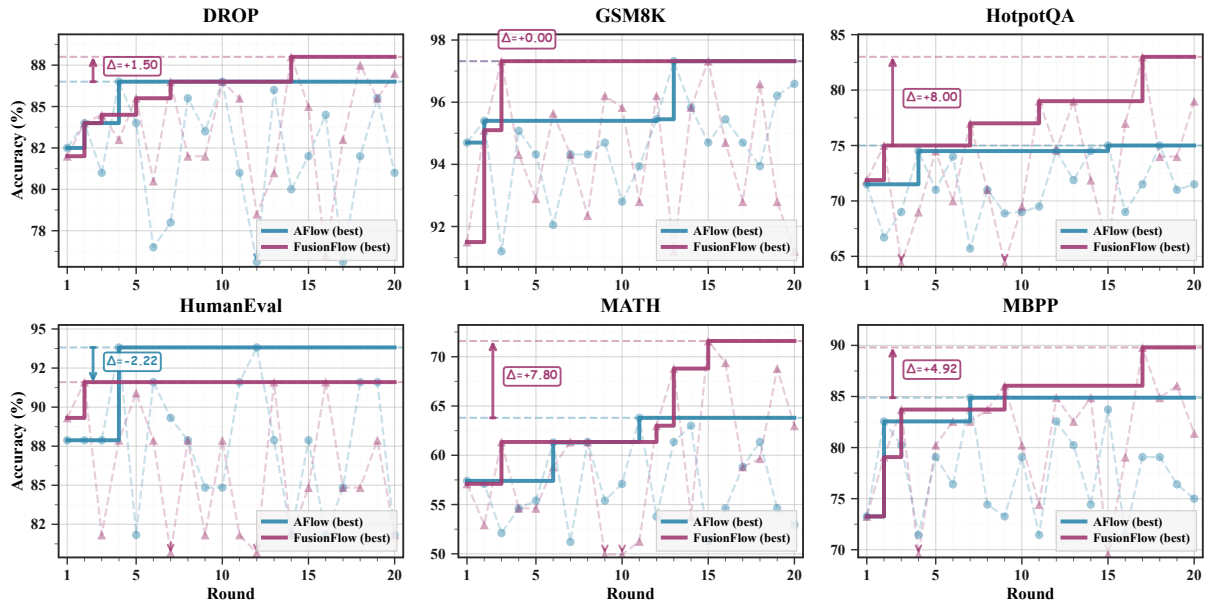


Figure 4: Validation accuracy during optimization across 20 rounds. Scatter points show the validation accuracy of individual candidate workflows at each round, while solid step curves report the best-so-far validation accuracy, which serves as the workflow selection criterion. Horizontal dashed lines indicate the final best validation accuracy achieved by each method, with  $\Delta$  denoting the absolute improvement (percentage points) of **FUSIONFLOW** over AFlow.

claim that structural exploration through fusion is an effective way to overcome depth limitations.

#### 4.4 Robustness Across Execution Models

To evaluate robustness, we test **FUSIONFLOW** with four different executor models on the MATH dataset: GPT-4o-mini, GPT-4o, Qwen-2.5-14b, and Llama-3-8b. As shown in Table 3, **FUSIONFLOW** consistently outperforms Direct IO and AFlow across all executors. While absolute accuracy varies with model capability, the relative improvement over AFlow remains stable, indicating that the benefits of fusion-based structural exploration are largely model-agnostic.

#### 4.5 Optimization Dynamics

We analyze optimization trajectories under identical iteration budgets of 20 rounds across six benchmarks. The primary metric is the **best-so-far validation accuracy** achieved across all rounds (step curves in Figure 4), since the final workflow is selected based on validation performance and evaluated only once on the test set.

As shown in Figure 4, two distinct regimes emerge. On relatively simple benchmarks (GSM8K and HumanEval), both methods converge rapidly, with best-so-far curves plateauing within a few iterations, indicating that shallow workflows are sufficient and leaving limited room for structural

exploration. In contrast, on more complex benchmarks (HotpotQA, MBPP, and MATH), AFlow stagnates early, whereas **FUSIONFLOW** continues to discover improved workflows, achieving absolute accuracy gains of +8.0%, +4.9%, and +7.8%, respectively.

Notably, **FUSIONFLOW** exhibits higher per-round variance, which is a direct consequence of the exploratory operators (differentiation and fusion) that go beyond single-step refinement. This exploration mechanism enables **FUSIONFLOW** to break through the structural depth limitations, escaping local optima and ultimately attaining higher best-so-far performance.

## 5 Conclusion

Existing automated workflow generation methods based on incremental refinement are constrained by structural depth limitations, often leading to early stagnation under finite optimization budgets. To address this issue, we propose **FUSIONFLOW**, a framework that enables structural exploration through workflow fusion, supported by optimization, differentiation, and dynamic scheduling. Experiments on six benchmarks demonstrate consistent improvements over existing methods, showing that fusion-driven exploration is an effective approach for overcoming depth limitations in automated workflow generation.

## 556 Limitations

557 This work shares several limitations common to au-  
558 tomated agentic workflow generation frameworks.  
559 First, introducing workflow optimization inevitably  
560 incurs additional computational overhead com-  
561 pared to directly applying a fixed model at infer-  
562 ence time. Although our experiments adopt rela-  
563 tively low-cost execution models, which may lead  
564 to comparable or even lower deployment costs than  
565 using stronger models directly, the optimization  
566 stage still introduces extra expense that may be  
567 undesirable in latency- or budget-critical settings.

568 Second, workflow optimization can reduce flex-  
569 ibility and generalization across diverse task dis-  
570 tributions. Effective workflows often require opti-  
571 mization on a representative subset of tasks before  
572 deployment, which limits their ability to adapt in-  
573 stantly to unseen or highly heterogeneous problem  
574 domains. As a result, automated workflow genera-  
575 tion is better suited to scenarios where task distribu-  
576 tions are relatively stable, domain specialization is  
577 desired, and low-cost models are preferred, while  
578 some degree of offline optimization or preprocess-  
579 ing is acceptable.

580 Addressing these limitations, particularly im-  
581 proving generalization efficiency and reducing opti-  
582 mization overhead, remains an important direction  
583 for future research.

## 584 References

585 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
586 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
587 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
588 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-  
589 cal report. *arXiv preprint arXiv:2303.08774*.

590 Anthropic. 2025a. [Claude 3.7 sonnet and claude code](#).

591 Anthropic. 2025b. [Introducing claude 4](#).

592 Jacob Austin, Augustus Odena, Maxwell I Nye,  
593 Maarten Bosma, Henryk Michalewski, David Do-  
594 han, Ellen Jiang, Carrie J Cai, Michael Terry, Quoc V  
595 Le, and 1 others. 2021. Program synthesis with large  
596 language models. corr abs/2108.07732 (2021). *arXiv*  
597 *preprint arXiv:2108.07732*.

598 Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang,  
599 Jaward Sesay, Börje F Karlsson, Jie Fu, and Yemin  
600 Shi. 2023. Autoagents: A framework for automatic  
601 agent generation. *arXiv preprint arXiv:2309.17288*.

602 Mark Chen. 2021. Evaluating large language models  
603 trained on code. *arXiv preprint arXiv:2107.03374*.

Ching-An Cheng, Allen Nie, and Adith Swaminathan. 604  
2024. Trace is the next autodiff: Generative opti- 605  
mization with rich feedback, execution traces, and 606  
llms. *Advances in Neural Information Processing* 607  
*Systems*, 37:71596–71642. 608

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 609  
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 610  
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 611  
Nakano, and 1 others. 2021. Training verifiers 612  
to solve math word problems. *arXiv preprint* 613  
*arXiv:2110.14168*. 614

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel 615  
Stanovsky, Sameer Singh, and Matt Gardner. 2019. 616  
Drop: A reading comprehension benchmark re- 617  
quiring discrete reasoning over paragraphs. *arXiv* 618  
*preprint arXiv:1903.00161*. 619

Chrisantha Fernando, Dylan Banarse, Henryk 620  
Michalewski, Simon Osindero, and Tim Rock- 621  
täschel. 2023. Promptbreeder: Self-referential 622  
self-improvement via prompt evolution. *arXiv* 623  
*preprint arXiv:2309.16797*. 624

Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, 625  
Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi- 626  
angliang Zhang. 2024. Large language model based 627  
multi-agents: A survey of progress and challenges. 628  
*arXiv preprint arXiv:2402.01680*. 629

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul 630  
Arora, Steven Basart, Eric Tang, Dawn Song, and Ja- 631  
cob Steinhardt. 2021. Measuring mathematical prob- 632  
lem solving with the math dataset. *arXiv preprint* 633  
*arXiv:2103.03874*. 634

Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu 635  
Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, 636  
Zili Wang, Steven Ka Shing Yau, Zijuan Lin, and 637  
1 others. 2023. Metagpt: Meta programming for a 638  
multi-agent collaborative framework. In *The Twelfth* 639  
*International Conference on Learning Representa-* 640  
*tions*. 641

Shengran Hu, Cong Lu, and Jeff Clune. 2024. Au- 642  
tomated design of agentic systems. *arXiv preprint* 643  
*arXiv:2408.08435*. 644

Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei 645  
Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruim- 646  
ing Tang, and Enhong Chen. 2024. Understanding 647  
the planning of llm agents: A survey. *arXiv preprint* 648  
*arXiv:2402.02716*. 649

Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, 650  
Zhiyuan Zhang, Keshav Santhanam, Saiful Haq, 651  
Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, 652  
Heather Miller, and 1 others. 2024. Dspy: Compiling 653  
declarative language model calls into state-of-the-art 654  
pipelines. In *The Twelfth International Conference* 655  
*on Learning Representations*. 656

Zelong Li, Shuyuan Xu, Kai Mei, Wenyue Hua, Bal- 657  
aji Rama, Om Raheja, Hao Wang, He Zhu, and 658

659	Yongfeng Zhang. 2024. Autoflow: Automated workflow generation for large language model agents. <i>arXiv preprint arXiv:2407.12821</i> .	Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2022. Progprompt: Generating situated robot task plans using large language models. <i>arXiv preprint arXiv:2209.11302</i> .	714
660			715
661			716
662	Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. 2022. Code as policies: Language model programs for embodied control. <i>arXiv preprint arXiv:2209.07753</i> .		717
663			718
664			719
665		Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. 2023. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 2998–3009.	720
666			721
667	Bang Liu, Xinfeng Li, Jiayi Zhang, Jinlin Wang, Tanjin He, Sirui Hong, Hongzhang Liu, Shaokun Zhang, Kaitao Song, Kunlun Zhu, and 1 others. 2025. Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems. <i>arXiv preprint arXiv:2504.01990</i> .		722
668			723
669			724
670			725
671		Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	726
672			727
673			728
674	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, and 1 others. 2023. Self-refine: Iterative refinement with self-feedback. <i>Advances in Neural Information Processing Systems</i> , 36:46534–46594.		729
675			730
676			731
677			732
678		Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, and 1 others. 2024a. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> , 18(6):186345.	733
679			734
680			735
681	Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King, Jonathan Larson, Yuanzhi Li, Weishung Liu, and 1 others. 2023. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. <i>arXiv preprint arXiv:2311.16452</i> .		736
682			737
683			738
684			739
685			740
686	OpenAI. 2024. <a href="#">Gpt-4o mini: advancing cost-efficient intelligence</a> .		741
687			742
688			743
689	Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, and 1 others. 2024a. Chatdev: Communicative agents for software development. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15174–15186.		744
690			745
691			746
692			747
693			748
694			749
695			750
696	Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Kunlun Zhu, Hanchen Xia, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, and 1 others. 2024b. Scaling large language model-based multi-agent collaboration. <i>arXiv preprint arXiv:2406.07155</i> .		751
697			752
698			753
699			754
700	Shuofei Qiao, Ningyu Zhang, Runnan Fang, Yujie Luo, Wangchunshu Zhou, Yuchen Eleanor Jiang, Chengfei Lv, and Huajun Chen. 2024. Autoact: Automatic agent learning from scratch for qa via self-planning. <i>arXiv preprint arXiv:2401.05268</i> .		755
701			756
702			757
703			758
704			759
705	Tal Ridnik, Dedy Kredo, and Itamar Friedman. 2024. Code generation with alphacodium: From prompt engineering to flow engineering. <i>arXiv preprint arXiv:2401.08500</i> .		760
706			761
707			762
708			763
709	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 36:8634–8652.		764
710			765
711			766
712			767
713			768
			769
			770
		Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. 2024d. Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 257–279.	751
			752
			753
			754
			755
			756
			757
			758
			759
			760
			761
			762
			763
			764
			765
			766
			767
			768
			769
			770

771	Benfeng Xu, An Yang, Junyang Lin, Quan Wang,	Yuqi Zhu, Shuofei Qiao, Yixin Ou, Shumin Deng,	826
772	Chang Zhou, Yongdong Zhang, and Zhendong Mao.	Shiwei Lyu, Yue Shen, Lei Liang, Jinjie Gu, Hua-	827
773	2023. Expertprompting: Instructing large language	jun Chen, and Ningyu Zhang. 2024. Knowa-	828
774	models to be distinguished experts. <i>arXiv preprint</i>	gent: Knowledge-augmented planning for llm-based	829
775	<i>arXiv:2305.14688</i> .	agents. <i>arXiv preprint arXiv:2403.03101</i> .	830
776	Shicheng Xu, Liang Pang, Huawei Shen, Xueqi Cheng,	Mingchen Zhuge, Haozhe Liu, Francesco Faccio, Dy-	831
777	and Tat-Seng Chua. 2024. Search-in-the-chain: Inter-	lan R Ashley, Róbert Csordás, Anand Gopalakrish-	832
778	actively enhancing large language models with search	nan, Abdullah Hamdi, Hasan Abed Al Kader Ham-	833
779	for knowledge-intensive tasks. In <i>Proceedings of the</i>	moud, Vincent Herrmann, Kazuki Irie, and 1 others.	834
780	<i>ACM Web Conference 2024</i> , pages 1362–1373.	2023. Mindstorms in natural language-based soci-	835
781	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao	eties of mind. <i>arXiv preprint arXiv:2305.17066</i> .	836
782	Liu, Quoc V Le, Denny Zhou, and Xinyun Chen.	Mingchen Zhuge, Wenyi Wang, Louis Kirsch,	837
783	2023. Large language models as optimizers. In	Francesco Faccio, Dmitrii Khizbullin, and Jürgen	838
784	<i>The Twelfth International Conference on Learning</i>	Schmidhuber. 2024. Gptswarm: Language agents	839
785	<i>Representations</i> .	as optimizable graphs. In <i>Forty-first International</i>	840
786	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,	<i>Conference on Machine Learning</i> .	841
787	William Cohen, Ruslan Salakhutdinov, and Christo-	<b>A Appendix: Datasets and Experimental</b>	842
788	pher D Manning. 2018. Hotpotqa: A dataset for	<b>Setup</b>	843
789	diverse, explainable multi-hop question answering.	In this appendix, we provide a detailed account of	844
790	In <i>Proceedings of the 2018 conference on empiri-</i>	the datasets, data partitioning, and inference details	845
791	<i>cal methods in natural language processing</i> , pages	in our experiments.	846
792	2369–2380.	<b>A.1 Dataset Descriptions</b>	847
793	Mert Yuksekgonul, Federico Bianchi, Joseph Boen,	We conducted our experiments on six widely-used	848
794	Sheng Liu, Zhi Huang, Carlos Guestrin, and James	public benchmarks spanning a diverse range of rea-	849
795	Zou. 2024. Textgrad: Automatic "differentiation" via	soning tasks, including question answering, code	850
796	text. <i>arXiv preprint arXiv:2406.07496</i> .	generation, and mathematical problem-solving. For	851
797	Matei Zaharia, Omar Khattab, Lingjiao Chen,	each dataset, we used a subset of instances, which	852
798	Jared Quincy Davis, Heather Miller, Chris Potts,	were partitioned into a validation set for workflow	853
799	James Zou, Michael Carbin, Jonathan Fran-	optimization and a test set for final performance	854
800	kle, Naveen Rao, and Ali Ghodsi. 2024. The	evaluation.	855
801	shift from models to compound ai systems.	<b>HotpotQA (Yang et al., 2018)</b> A multi-hop ques-	856
802	<a href="https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/">https://bair.berkeley.edu/blog/2024/02/</a>	tion answering dataset that requires finding and	857
803	<a href="https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/">18/compound-ai-systems/</a> .	reasoning over multiple supporting documents to	858
804	Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng,	answer a question.	859
805	Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin	• <b>Task Type:</b> Multi-hop QA	860
806	Cheng, Sirui Hong, Jinlin Wang, and 1 others. 2024.	• <b>Data Source:</b> 1,000 instances randomly sam-	861
807	Aflow: Automating agentic workflow generation.	pled from the original test set.	862
808	<i>arXiv preprint arXiv:2410.10762</i> .	• <b>Data Split:</b> 200 instances for the validation	863
809	Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman,	set and 800 for the test set.	864
810	Haohan Wang, and Yu-Xiong Wang. 2023. Lan-	<b>DROP (Dua et al., 2019)</b> A challenging reading	865
811	guage agent tree search unifies reasoning acting	comprehension benchmark requiring discrete rea-	866
812	and planning in language models. <i>arXiv preprint</i>	soning over paragraphs, such as addition, counting,	867
813	<i>arXiv:2310.04406</i> .	or sorting.	868
814	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei,	• <b>Task Type:</b> Reading Comprehension	869
815	Nathan Scales, Xuezhi Wang, Dale Schuurmans,	• <b>Data Source:</b> 1,000 instances randomly sam-	870
816	Claire Cui, Olivier Bousquet, Quoc Le, and 1 oth-	pled from the original development (dev) set.	871
817	ers. 2022. Least-to-most prompting enables complex	• <b>Data Split:</b> 200 instances for the validation	872
818	reasoning in large language models. <i>arXiv preprint</i>	set and 800 for the test set.	873
819	<i>arXiv:2205.10625</i> .		
820	Pei Zhou, Jay Pujara, Xiang Ren, Xinyun Chen, Heng-		
821	Tze Cheng, Quoc V Le, Ed Chi, Denny Zhou, Swa-		
822	roop Mishra, and Huaixiu Steven Zheng. 2024. Self-		
823	discover: Large language models self-compose rea-		
824	soning structures. <i>Advances in Neural Information</i>		
825	<i>Processing Systems</i> , 37:126032–126058.		

**HumanEval (Chen, 2021)** A code generation dataset comprising 164 handwritten programming problems. Each problem includes a function signature, docstring, body, and several unit tests.

- **Task Type:** Code Generation
- **Data Source:** The full original test set.
- **Data Split:** 33 instances for the validation set and the remaining 131 for the test set, following prior work.

**MBPP (Mostly Basic Python Problems) (Austin et al., 2021)** An entry-level Python programming dataset consisting of around 500 crowd-sourced tasks, each with a short description, a code solution, and three test cases.

- **Task Type:** Code Generation
- **Data Source:** The full original test set.
- **Data Split:** 86 instances for the validation set and the remaining 341 for the test set.

**GSM8K (Cobbe et al., 2021)** A dataset of high-quality, grade-school math word problems that require a sequence of elementary arithmetic operations to solve.

- **Task Type:** Math Word Problems
- **Data Source:** The full original test set.
- **Data Split:** 264 instances for the validation set and 1055 for the test set.

**MATH (Hendrycks et al., 2021)** A challenging dataset of mathematics competition problems. Our experiments focus on a curated subset of 605 problems from four specific domains-Combinatorics & Probability, Number Theory, Pre-algebra, and Pre-calculus-all at difficulty level 5.

- **Task Type:** Advanced Mathematics
- **Data Source:** A curated subset of 605 problems at difficulty level 5.
- **Data Split:** 119 instances for the validation set and 486 for the test set.

## A.2 Inference Details

We use different models for the optimization and execution phases. Specifically, we employ Claude-3-7-Sonnet (Anthropic, 2025a) as the optimizer and GPT-4o-mini (OpenAI, 2024) as the executor. All models are accessed via APIs with a temperature of 0. For our method and AFlow, we set the maximum number of optimization iterations to 20. This setting follows the standard configuration in prior work and ensures a fair comparison under limited optimization budgets. For the manually designed prompting methods, we use the same executor model (GPT-4o-mini) to ensure a fair comparison. This design choice is intentional: we leverage a more powerful model for the optimizer, which is called infrequently, and a more cost-effective model for the executor, which is invoked for every task step. This strategy allows us to achieve high-quality workflow discovery while controlling the overall inference cost.

## A.3 API Cost Analysis

This section provides a breakdown of the total API costs incurred for running the complete FUSIONFLOW and AFlow optimization process (20 iterations) on each of the six benchmark datasets. The costs, summarized in Table 4, are reported in US Dollars (\$) and reflect the expenditure for the entire workflow discovery phase.

Table 4: Total API cost comparison between FUSIONFLOW and AFlow.

Dataset	FUSIONFLOW (\$)	AFlow (\$)
HotpotQA	5.8	6.7
DROP	6.5	4.2
HumanEval	3.0	3.3
MBPP	4.7	4.1
GSM8K	14.8	15.7
MATH	18.4	16.0

It is important to note that these figures represent the one-time cost for the workflow discovery and optimization phase. Once the optimal workflow is generated, the inference cost for solving individual tasks using the more economical executor model (GPT-4o-mini) is substantially lower.

## B Theoretical Analysis

This appendix provides theoretical insights that motivate the design of FUSIONFLOW. Our analysis

clarifies (i) why incremental optimization alone is insufficient under finite budgets, (ii) why workflow fusion induces different search dynamics, and (iii) why scheduling and candidate selection are necessarily post-hoc and heuristic in nature.

### B.1 Problem Structure and Search Space

Let  $\mathbb{S}$  denote the space of all executable agentic workflows. Each workflow  $\mathcal{W} \in \mathbb{S}$  is represented as a directed control-flow graph composed of atomic operators and explicit coordination logic. We denote by  $d(\mathcal{W})$  the **structural depth** of a workflow, defined as the number of effective modifications it contains. Each workflow is associated with a modification-node set  $\mathcal{N}(\mathcal{W})$ , and its depth is  $d(\mathcal{W}) = |\mathcal{N}(\mathcal{W})|$ . The formal recursive definition and detailed analysis are provided in Appendix D.

The automated workflow generation problem can be viewed as a search process over  $\mathbb{S}$  under a finite iteration budget  $T$ . At each iteration, the search applies a transformation operator that maps one or more existing workflows to new candidates. Crucially, the evaluation of a workflow’s utility is only observable through execution feedback on a finite dataset, making the search process partially observable and non-stationary.

### B.2 Limitations of Incremental Optimization

We first analyze the limitations of purely incremental optimization. Optimization operators modify a workflow through local edits, such as prompt refinement, parameter tuning, or small structural adjustments. Let  $\mathcal{O}_{\text{opt}}$  denote such an operator.

**Observation 1 (Locality of Optimization).** Each application of  $\mathcal{O}_{\text{opt}}$  induces only a bounded structural change to the workflow graph. Consequently, under a finite budget  $T$ , the structural depth achievable through repeated optimization is at most linear in  $T$ .

This implies that even if optimization consistently improves performance, it explores the workflow space along a single trajectory, accumulating depth gradually. In practice, this leads to early saturation: once local refinements exhaust easily accessible improvements, further iterations yield diminishing returns.

### B.3 Differentiation as Width Expansion

Differentiation operators address this limitation by explicitly increasing exploration width. Given a

parent workflow, differentiation generates specialized variants that target different task subspaces.

From a search perspective, differentiation expands the frontier of exploration by creating multiple distinct trajectories. However, differentiation alone does not increase structural depth: each differentiated workflow remains confined to a similar depth regime as its parent.

**Observation 2 (Width without Depth).** Differentiation increases diversity across trajectories but does not, by itself, induce super-linear depth growth. Under finite budgets, it primarily redistributes exploration rather than enabling deeper compositions.

Thus, while differentiation is essential for discovering complementary behaviors, it cannot overcome depth limitations on its own.

### B.4 Fusion as a Non-local Search Operator

We now turn to fusion, which constitutes the central mechanism of **FUSIONFLOW**. Let  $\mathcal{O}_{\text{fuse}}$  denote the fusion operator that maps a group of workflows  $\mathcal{S} = \{\mathcal{W}_1, \dots, \mathcal{W}_m\}$  to a new workflow  $\mathcal{W}_{\text{fuse}}$ .

**Observation 3 (Non-local Transition).** Fusion induces a non-local transition in the workflow space by composing multiple independently evolved substructures into a single executable graph.

Specifically, if each  $\mathcal{W}_i$  contains a mature subgraph discovered along a distinct optimization trajectory, fusion materializes these subgraphs simultaneously. Under our depth definition in Appendix D.1, each workflow  $\mathcal{W}$  is associated with a modification-node set  $\mathcal{N}(\mathcal{W})$  and  $d(\mathcal{W}) = |\mathcal{N}(\mathcal{W})|$ . For a fusion executed at round  $t$ , the recursive construction gives  $\mathcal{N}(\mathcal{W}_{\text{fuse}}) = \bigcup_i \mathcal{N}(\mathcal{W}_i) \cup \{t\}$ . Therefore,

$$\begin{aligned} d(\mathcal{W}_{\text{fuse}}) &= |\mathcal{N}(\mathcal{W}_{\text{fuse}})| = \left| \bigcup_i \mathcal{N}(\mathcal{W}_i) \cup \{t\} \right| \\ &\geq \max_i |\mathcal{N}(\mathcal{W}_i)| = \max_i d(\mathcal{W}_i). \end{aligned} \tag{7}$$

Moreover, the fusion step introduces new coordination logic (captured by the fresh fusion node  $\{t\}$  in Appendix D.1), and the union can be strictly larger than any single parent when the parents contribute complementary modifications, in which case  $d(\mathcal{W}_{\text{fuse}})$  strictly exceeds all parent depths.

**Implication.** Under a fixed iteration budget, fusion enables discrete depth expansion that cannot

1041 be achieved through sequential local edits alone. 1089  
 1042 It effectively collapses multiple exploration trajec- 1090  
 1043 tories into a single step, bypassing intermediate 1091  
 1044 states that would otherwise require many rounds of 1092  
 1045 optimization. 1093

1046 This property distinguishes fusion from both op- 1094  
 1047 timization and differentiation. 1095

### 1048 **B.5 Why Fusion Is Feasible but Not Trivial** 1096

1049 Although fusion induces powerful search transi- 1097  
 1050 tions, it is not difficult to imply from an implemen- 1098  
 1051 tation perspective. Agentic workflows are compos- 1099  
 1052 ed of atomic operators connected by explicit, 1100  
 1053 code-level control logic. This structure allows a 1101  
 1054 large language model to synthesize a composite 1102  
 1055 workflow by identifying shared components and 1103  
 1056 reconciling complementary substructures. 1104

1057 However, effective fusion is not guaranteed. Ar- 1105  
 1058 bitrarily fusing workflows may result in redun- 1106  
 1059 dancy, incompatible coordination, or degraded per- 1107  
 1060 formance. Thus, the challenge of fusion lies not in 1108  
 1061 feasibility, but in selectivity. 1109

### 1062 **B.6 Post-hoc Nature of Fusion Scheduling** 1110

1063 A central difficulty in fusion-based search is decid- 1111  
 1064 ing *when* to apply fusion and *which* workflows to 1112  
 1065 fuse. 1113

1066 The utility of fusing a candidate group depends 1114  
 1067 on latent properties such as structural compatibility 1115  
 1068 and complementary task coverage. These proper- 1116  
 1069 ties are only observable through execution feedback 1117  
 1070 and cannot be inferred reliably a priori. 1118

### 1071 **B.7 Theoretical Motivation and Practical** 1119 1072 **Approximation** 1120

1073 This section justifies why FUSIONFLOW adopts a 1121  
 1074 feedback-driven scheduling policy. Our key claim 1122  
 1075 is that workflow scheduling is inherently a sequen- 1123  
 1076 tial decision problem under partial observability: 1124  
 1077 each operator choice not only changes the current 1125  
 1078 best score, but also reshapes the future search fron- 1126  
 1079 tier (what workflows exist, what can be fused, and 1127  
 1080 what information becomes available). 1128

1081 **Scheduling as a POMDP.** Let the search process 1129  
 1082 maintain a workflow pool  $\mathcal{P}_t$  at iteration  $t$ . The 1130  
 1083 scheduler chooses an operator (and corresponding 1131  
 1084 targets) to apply, producing  $\mathcal{P}_{t+1}$  and new evalua- 1132  
 1085 tion feedback. This can be abstracted as a POMDP: 1133

$$1086 \mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{O}, \Omega, r, \gamma), \quad (8) \quad 1134$$

1087 where the latent state  $s_t \in \mathcal{S}$  summarizes the unob- 1135  
 1088 served properties of the current search frontier (e.g., 1136

1089 latent structural compatibility between workflows, 1090  
 1091 potential synergy from fusion, and the “distance” to 1092  
 1093 high-performing regions). The action  $a_t \in \mathcal{A}$  cor- 1094  
 1095 responds to selecting an operator (Optimization/D- 1096  
 1097 ifferentiation/Fusion) and its operands (selected 1098  
 1099 workflow(s) or group). After executing  $a_t$ , the pro- 1100  
 1101 cess transitions according to an unknown dynamics 1102  
 1103  $\mathcal{T}$ , and produces an observation  $o_t \in \mathcal{O}$ , which in- 1104  
 1105 cludes measurable statistics such as best score  $R_t$ , 1106  
 1107 per-workflow coverage sets  $\mathbf{C}_i$ , and execution logs. 1108  
 1109 The scheduler only observes  $o_t \sim \Omega(\cdot | s_t, a_{t-1})$  1109  
 1110 and thus cannot directly compute an optimal policy 1110  
 1111 over  $s_t$ . Importantly, the consequence of an op- 1111  
 1112 erator is intrinsically multi-step: a differentiation 1112  
 1113 operator primarily increases future fusion upside by 1113  
 1114 creating complementary specialists, and a fusion 1114  
 1115 operator becomes beneficial after sufficient matu- 1115  
 1116 ration of its parents. This induces a long-horizon 1116  
 1117 credit assignment problem, where the value of an 1117  
 1118 action depends on how it alters future workflow 1118  
 1119 pools and future available compositions, which is 1119  
 1120 exactly the dependency that a POMDP formulation 1120  
 1121 makes explicit. 1121

1122 Two practical properties make the process ad- 1122  
 1123 ditionally challenging: (i) the transition is non- 1123  
 1124 stationary because the pool distribution changes 1124  
 1125 as new workflows are generated; (ii) the reward is 1125  
 1126 expensive and delayed since any operator’s effect is 1126  
 1127 only known after execution evaluation. Under these 1127  
 1128 conditions, an effective scheduler must depend on 1128  
 1129 history feedback. This motivates a feedback-driven 1129  
 1130 policy class that uses lightweight sufficient statis- 1130  
 1131 tics extracted from  $o_t$  to approximate the decision- 1131  
 1132 making signal required by the underlying POMDP. 1132

1133 **Exploration–Exploitation View: Operators as** 1123  
 1134 **Arms in a Contextual Bandit.** At a coarse level, 1124  
 1135 the scheduler repeatedly chooses among operator 1125  
 1136 types, which can be viewed as a contextual bandit 1126  
 1137 with three arms: 1127

$$1138 a_t \in \{\mathcal{O}_{\text{opt}}, \mathcal{O}_{\text{diff}}, \mathcal{O}_{\text{fuse}}\}, \quad (9) \quad 1128$$

1129 where the context is the observation  $o_t$  (recent score 1129  
 1130 trajectory, pool diversity, coverage patterns). Op- 1130  
 1131 timization corresponds to exploitation (local im- 1131  
 1132 provements), while differentiation and fusion cor- 1132  
 1133 respond to exploration in structure space. Under 1133  
 1134 finite budgets, an effective policy should allocate 1134  
 1135 more probability mass to exploration when the 1135  
 1136 marginal returns of exploitation diminish. 1136

1137 **Stagnation as a Proxy for Negative Advantage.** 1137  
 1138 In a rigorous RL setting, action selection is of- 1138

ten driven by an advantage function  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ , which indicates whether a specific action improves upon the baseline value. However, since the latent state  $s_t$  is unobservable and computing  $Q^\pi$  requires prohibitive sampling costs, we seek a computationally tractable estimator. We observe that when the marginal return of local optimization diminishes, the advantage of exploitation decreases relative to exploration. Therefore, we employ the measurable improvement rate  $\Delta_t$  (defined below) as an empirical estimator for the advantage of repeated local refinement. Let

$$R_t^{\text{recent}} = \max_{u \in [t-k+1, t]} R_u, \quad (10)$$

$$R_t^{\text{prev}} = \max_{u \in [t-2k+1, t-k]} R_u, \quad (11)$$

$$\Delta_t = R_t^{\text{recent}} - R_t^{\text{prev}}. \quad (12)$$

When  $\Delta_t$  is small or negative, the empirical advantage of repeated local refinement becomes weak, suggesting that the expected gain of  $\mathcal{O}_{\text{opt}}$  is decreasing. We therefore define a monotone transformation into a stagnation signal:

$$\pi_t = \sigma(-\kappa\Delta_t) = \frac{1}{1 + \exp(\kappa\Delta_t)}, \quad (13)$$

which can be interpreted as a *soft gating mechanism*, approximating a stochastic policy that shifts probability mass toward high-variance structural exploration (Fusion/Differentiation) when the estimated advantage of exploitation vanishes.

**Why Fusion Should be Scheduled More Aggressively than Generic Exploration.** Differentiation and fusion both expand exploration, but fusion uniquely induces non-local transitions by composing mature substructures. In the POMDP view, fusion is a high-variance, potentially high-reward action whose payoff depends on latent compatibility between workflows. It is therefore sensible to schedule fusion under two conditions: (i) *exploitation plateau* (high  $\pi_t$ ), when local search returns diminish; (ii) *candidate maturity*, when the pool contains sufficiently competent and diverse workflows that make positive fusion payoff more likely. This matches standard RL intuitions: reserve high-variance actions for states where the baseline action underperforms, and when the context indicates high upside.

**Connection to Tree Search and Mixed-Probability Selection.** Prior automated

workflow search frameworks have also adopted RL-inspired selection mechanisms to balance exploration and exploitation, e.g., Monte Carlo Tree Search variants and mixed-probability selection rules. Such designs empirically validate the necessity of stochastic, feedback-driven scheduling in large, expensive, and partially observable search spaces. (Zhang et al., 2024) uses MCTS-style selection and execution feedback to guide workflow expansion, which is conceptually aligned with our view that operator scheduling should be adaptive and history-dependent rather than fixed.

**Takeaway.** In summary, **FUSIONFLOW**'s scheduling is theoretically motivated as a tractable instantiation of decision-making under partial observability and expensive feedback: (1) model the process as a POMDP / contextual bandit; (2) use measurable score momentum  $\Delta_t$  as a proxy for exploitation advantage; and (3) increase the probability of high-upside structural actions (especially fusion) when exploitation stagnates and candidate maturity is satisfied. This provides a principled justification for our lightweight, feedback-driven scheduling mechanism.

**Observation 4 (A Posteriori Decision Problem).**

Fusion scheduling constitutes an a posteriori decision problem under partial observability. No fixed or optimal fusion policy can be defined in advance without executing the fused workflow itself.

As a result, scheduling must rely on indirect signals, such as performance stagnation, as proxies for diminishing returns from local refinement. This motivates the use of lightweight, feedback-driven scheduling mechanisms rather than explicit policy optimization.

## B.8 Rationale for Heuristic Candidate Selection

Similarly, selecting fusion or differentiation candidates involves balancing complementary objectives-coverage, consensus, and computational tractability-under incomplete information.

Attempting exhaustive or optimal selection is infeasible due to combinatorial complexity and the cost of evaluation. Instead, **FUSIONFLOW** adopts heuristic scoring functions that are grounded in observable execution statistics, such as task-level coverage sets.

**Observation 5 (Principled Heuristics).** Although heuristic, the candidate selection mechanisms in **FUSIONFLOW** are principled in the sense that they directly operationalize the structural objectives required for effective fusion and differentiation.

## B.9 Summary

Taken together, this analysis explains why **FUSIONFLOW** combines optimization, differentiation, and fusion under a dynamic, feedback-driven scheduling framework. Incremental optimization alone is insufficient for deep exploration under finite budgets. Differentiation expands exploration width but not depth. Fusion enables non-local depth expansion but requires careful, post-hoc control.

The design of **FUSIONFLOW** is therefore not ad hoc, but a direct response to the structural constraints of the workflow search problem.

## C Detailed Algorithmic Description

This appendix provides a complete and formal specification of the algorithmic components of **FUSIONFLOW**, including execution feedback representation, dynamic scheduling, workflow selection for differentiation, fusion candidate selection, and fusion execution. The notation and formulation here are fully consistent with the main text.

### C.1 Workflow Pool and Execution Feedback

At each iteration  $t$ , **FUSIONFLOW** maintains a workflow pool

$$\mathcal{P}_t = \{\mathcal{W}_1, \mathcal{W}_2, \dots, \mathcal{W}_{|\mathcal{P}_t|}\}, \quad (14)$$

where each  $\mathcal{W}_i$  is a fully executable agentic workflow.

Each workflow is evaluated on the task dataset

$$\mathcal{D} = \{(x_j, y_j)\}_{j=1}^{|\mathcal{D}|}, \quad (15)$$

using a task-specific evaluation metric  $G(\cdot)$ . The execution feedback of workflow  $\mathcal{W}_i$  is summarized by:

### Scalar Performance Score

$$R_i = \mathbb{E}_{(x,y) \sim \mathcal{D}}[G(f_{\mathcal{W}_i}(x), y)], \quad (16)$$

which reflects the overall task performance of the workflow.

### Task-level Coverage Set

$$C_i \subset \mathcal{D}, \quad (17)$$

where  $C_i$  denotes the subset of task instances in  $\mathcal{D}$  that  $\mathcal{W}_i$  successfully solves. The set  $C_i$  serves as a discrete approximation of the functional coverage of the workflow and forms the basis for fusion candidate selection.

### C.2 Global Scope of Dynamic Scheduling

The dynamic scheduling mechanism operates over all transformation operators in **FUSIONFLOW**, including optimization, differentiation, and fusion. At each iteration, the scheduler determines which operator to apply based on recent performance dynamics and operator usage history.

Since workflow fusion constitutes the core mechanism for inducing non-local structural transitions, the main text focuses on fusion-related scheduling. This appendix presents the complete scheduling formulation.

### C.3 Performance Stagnation Signal

Let

$$R_t = \max_{\mathcal{W}_i \in \mathcal{P}_t} R_i \quad (18)$$

denote the best performance achieved at iteration  $t$ .

To detect diminishing returns from local refinement, we compute two sliding-window statistics with window size  $k$ :

$$\begin{aligned} R_t^{\text{recent}} &= \max_{u \in [t-k+1, t]} R_u, \\ R_t^{\text{prev}} &= \max_{u \in [t-2k+1, t-k]} R_u. \end{aligned} \quad (19)$$

The improvement difference is defined as

$$\Delta_t = R_t^{\text{recent}} - R_t^{\text{prev}}. \quad (20)$$

We map this quantity to a continuous stagnation signal:

$$\pi_t = \frac{1}{1 + \exp(\kappa \cdot \Delta_t)}, \quad (21)$$

where  $\kappa > 0$  controls sensitivity to performance stagnation. Larger values of  $\pi_t$  indicate reduced marginal gains from local refinement.

The signal  $\pi_t$  is not treated as a decision rule, but as a soft indicator that biases the scheduler toward exploration-oriented operators.

#### 1309 C.4 Workflow Selection Mechanism

1310 Both differentiation and fusion workflows are es-  
 1311 sentially selection processes driven by distinct po-  
 1312 tential algorithms. We unify their formulations  
 1313 based on the set of correctly solved tasks. Let  
 1314  $\mathbf{C}_i \subseteq \mathcal{D}$  denote the set of task instances correctly  
 1315 processed by workflow  $\mathcal{W}_i$ .

##### 1316 C.4.1 Differentiation Candidate Selection

1317 Differentiation aims to select a parent workflow  
 1318  $\mathcal{W}_i$  that exhibits high potential for specializing in a  
 1319 specific task cluster  $k$ .

1320 **Task Partition.** The dataset  $\mathcal{D}$  is partitioned into  
 1321 disjoint clusters  $\mathcal{K} = \{k_1, \dots, k_{|\mathcal{K}|}\}$ . For any cluster  
 1322  $k \in \mathcal{K}$ , the local correct set of workflow  $\mathcal{W}_i$  is  
 1323 defined as:

$$1324 \mathbf{C}_{i,k} = \mathbf{C}_i \cap \mathcal{D}_k. \quad (22)$$

1325 **Performance Metrics.** We define the global ac-  
 1326 curacy rate  $R_{\text{global}}$  and the local cluster recall rate  
 1327  $R_{\text{local}}$  as:

$$1328 R_{\text{global}}(\mathcal{W}_i) = \frac{|\mathbf{C}_i|}{|\mathcal{D}|}, \quad (23)$$

$$R_{\text{local}}(\mathcal{W}_i, k) = \frac{|\mathbf{C}_{i,k}|}{|\mathcal{D}_k|}.$$

1329 A workflow is a valid candidate for cluster  $k$   
 1330 only if  $R_{\text{local}}(\mathcal{W}_i, k) > R_{\text{global}}(\mathcal{W}_i)$ .

1331 **Differentiation Potential.** The potential focuses  
 1332 on the absolute volume of specialized tasks. We de-  
 1333 fine the differentiation potential  $\Phi_{\text{diff}}$  for workflow  
 1334  $\mathcal{W}_i$  as the maximum specialization gain across all  
 1335 clusters:

$$1336 \Phi_{\text{diff}}(\mathcal{W}_i) = \max_{k \in \mathcal{K}} \left[ \frac{|\mathbf{C}_{i,k}|}{|\mathcal{D}|} \cdot (R_{\text{local}}(\mathcal{W}_i, k) \right. \\ \left. - R_{\text{global}}(\mathcal{W}_i)) \right]_+, \quad (24)$$

1337 where  $[\cdot]_+ = \max(0, \cdot)$ .

##### 1338 C.4.2 Fusion Candidate Selection

1339 Fusion aims to select a group of workflows  $\mathcal{S} =$   
 1340  $\{\mathcal{W}_1, \dots, \mathcal{W}_m\}$  that maximizes the joint potential  
 1341 defined by complementarity and consensus.

1342 **Set-Based Metrics.** We evaluate the group  $\mathcal{S}$  us-  
 1343 ing two fundamental set operations on their correct  
 1344 sets:

- 1345 • **Complementarity (Union):** Represents the  
 1346 collective coverage.

$$1347 U(\mathcal{S}) = \left| \bigcup_{\mathcal{W}_i \in \mathcal{S}} \mathbf{C}_i \right|. \quad (25)$$

- **Consensus (Intersection):** Represents the be-  
 havioral consistency.

$$I(\mathcal{S}) = \left| \bigcap_{\mathcal{W}_i \in \mathcal{S}} \mathbf{C}_i \right|. \quad (26)$$

Ideally, a fusion group should maximize the union  
 (solving more problems) while maintaining a suffi-  
 cient intersection (agreement on ground truths).

**Fusion Potential.** The fusion potential  $\Phi_{\text{fuse}}$   
 combines the set-based metrics. To capture both  
 pair-wise interactions and group-wise properties,  
 we formulate:

$$\Phi_{\text{fuse}}(\mathcal{S}) = (\lambda_1 \bar{U}_{\text{pair}} + \lambda_2 U(\mathcal{S})) + \\ (\mu_1 \bar{I}_{\text{pair}} + \mu_2 I(\mathcal{S})) \cdot \eta(\mathcal{S}). \quad (27)$$

where  $\bar{U}_{\text{pair}}$  and  $\bar{I}_{\text{pair}}$  are the average pairwise  
 union and intersection sizes within  $\mathcal{S}$ ,  $\lambda, \mu$  are  
 weighting coefficients, and  $\eta(\mathcal{S})$  is a penalty term  
 for group complexity.

The group maximizing the potential is selected:

$$\mathcal{S}^* = \arg \max_{\mathcal{S} \subseteq \mathcal{P}'} \Phi_{\text{fuse}}(\mathcal{S}). \quad (28)$$

#### 1365 C.5 Fusion Execution

1366 Given a selected group  $\mathcal{S}$ , fusion is performed by  
 1367 prompting a large language model to synthesize a  
 1368 composite workflow that preserves complementary  
 1369 substructures, reconciles shared components, and  
 1370 introduces coordination logic.

The resulting workflow  $\mathcal{W}_{\text{fuse}}$  is validated and  
 added to the pool:

$$1371 \mathcal{P}_{t+1} \leftarrow \mathcal{P}_t \cup \{\mathcal{W}_{\text{fuse}}\}. \quad (29)$$

#### 1374 C.6 Computational Considerations

1375 Although both differentiation and fusion in-  
 1376 volve combinatorial elements, practical complex-  
 1377 ity is controlled through coverage-based filtering,  
 1378 bounded group size, and probabilistic selection.  
 1379 In practice, these mechanisms keep the overhead  
 1380 of structure-level operations negligible relative to  
 1381 workflow evaluation.

## 1382 D Empirical Analysis of structural depth

1383 This section provides an empirical analysis of the  
 1384 **structural depth** achieved by different workflow  
 1385 generation methods. Intuitively, depth captures  
 1386 how many *effective structural changes* have been  
 1387 accumulated into a workflow. The key empirical

claim is that **fusion** can aggregate effective changes produced in different branches, thereby overcoming the depth limitations of incremental optimization approaches.

### D.1 Depth Definition and Computation

We define the **structural depth** of a workflow  $\mathcal{W}$  as the number of effective modifications it contains. Each workflow can be represented as a set of modification nodes, where each node corresponds to an optimization or fusion operation applied during the search process.

**Definition 1** (Structural Depth). *Let  $\mathcal{N}(\mathcal{W})$  denote the set of modification nodes contained in workflow  $\mathcal{W}$ , and  $d(\mathcal{W}) = |\mathcal{N}(\mathcal{W})|$  be its depth. The node set is computed recursively:*

- **Root workflow:**  $\mathcal{N}(\mathcal{W}_{root}) = \{r\}$ , where  $r$  is the root node.
- **Optimization:**  $\mathcal{N}(\mathcal{W}_{opt}) = \mathcal{N}(\mathcal{W}_{parent}) \cup \{t\}$ , where  $t$  is the current round.
- **Differentiation:**  $\mathcal{N}(\mathcal{W}_{diff}) = \mathcal{N}(\mathcal{W}_{parent})$  (no new modification).
- **Fusion:**  $\mathcal{N}(\mathcal{W}_{fuse}) = \bigcup_i \mathcal{N}(\mathcal{W}_i) \cup \{t\}$ , where  $\{\mathcal{W}_i\}$  are parent workflows.

This definition captures the intuition that:

- Each workflow represents a tree of modifications in the search space, where each node is a chain from the root.
- Optimization adds one new modification node to the parent’s tree.
- Differentiation creates specialized variants without introducing new modifications (breadth expansion).
- Fusion merges multiple independently evolved trees into a single tree (taking the union), plus the fusion operation itself. This enables aggregation of modifications from different evolutionary branches.

### D.2 Depth as Structural Capacity under Finite Budgets

It is important to clarify the interpretation and limitations of the structural depth metric used in this analysis. The depth defined here does not directly measure the semantic quality or functional optimality of a workflow. Instead, it characterizes the

maximum structural capacity that a workflow can embody under a finite optimization budget.

Under a fixed number of optimization rounds, exploratory operations such as differentiation are unavoidable in complex search spaces, as they are necessary to discover diverse and potentially complementary solution strategies. However, such exploration incurs an opportunity cost: once the budget is consumed by branching into multiple trajectories, purely incremental methods that refine a single workflow cannot subsequently integrate the effective modifications discovered along different branches into a single executable structure. As a result, the maximum structural complexity that any individual workflow can attain is inherently bounded by the depth of a single optimization chain.

Fusion addresses this limitation by enabling ex post aggregation of independently evolved workflows. By integrating modification histories from multiple search trajectories, fusion allows the optimization process to recover and consolidate structural information that would otherwise remain fragmented across branches. In this sense, fusion does not guarantee that the resulting workflow is intrinsically more effective or semantically superior, but it expands the space of structurally realizable workflows that can be reached within a fixed iteration budget.

In **FUSIONFLOW**, this increased structural capacity is further guided by optimization, differentiation, and dynamic scheduling, which are informed by empirical observations of task structure (e.g., sub-domain specialization in the dataset). While deeper workflows are not inherently better, this coordinated process increases the likelihood that workflows with higher structural capacity can represent and support more complex reasoning behaviors when such structure is required by the task distribution.

### D.3 Depth Accumulation and Structure-Aligned Exploration

While increased structural depth alone does not guarantee improved workflow quality, the manner in which depth is accumulated plays a critical role. In purely incremental optimization, additional modifications are applied along a single trajectory, often reflecting local adjustments whose composition may be incidental rather than structurally meaningful. As a result, deeper workflows obtained through such processes do not necessarily correspond to

Table 5: **structural depth** comparison between AFlow and FUSIONFLOW on MBPP dataset. Depth is defined as the number of effective modifications. Fusion operations (marked with †) aggregate modifications from multiple branches.

Round	AFlow		FUSIONFLOW	
	Depth	Depth	Op.	Modification Set
1	1	1	Root	{1}
2	2	2	Opt	{1, 2}
3	3	3	Opt	{1, 2, 3}
4	2	2	Diff	{1, 2}
5	3	2	Diff	{1, 2}
6	2	4	Fuse†	{1, 2, 3, 6}
7	4	3	Opt	{1, 2, 7}
8	2	3	Opt	{1, 2, 8}
9	2	4	Opt	{1, 2, 8, 9}
10	3	2	Diff	{1, 2}
11	4	5	Opt	{1, 2, 8, 9, 11}
12	3	5	Opt	{1, 2, 8, 9, 12}
13	3	6	Fuse†	{1, 2, 3, 8, 9, 13}
14	3	5	Opt	{1, 2, 8, 9, 14}
15	3	7	Fuse†	{1, 2, 3, 8, 9, 14, 15}
16	3	5	Opt	{1, 2, 8, 9, 16}
17	3	2	Diff	{1, 2}
18	2	7	Fuse†	{1, 2, 3, 8, 9, 14, 18}
19	4	6	Opt	{1, 2, 8, 9, 12, 19}
20	4	5	Opt	{1, 2, 8, 9, 20}
<b>Max</b>	4	7	–	–

more effective reasoning structures.

In contrast, FUSIONFLOW accumulates depth through an explicit differentiation–fusion pattern that mirrors a natural problem-solving strategy: exploring diverse solution variants first, and subsequently integrating complementary components into a unified workflow. Differentiation encourages specialization across branches, allowing distinct substructures to emerge under different inductive biases or task sub-distributions. Fusion then selectively aggregates these independently evolved structures, resulting in deeper workflows whose modifications originate from semantically distinct yet complementary search trajectories.

Under this paradigm, depth growth reflects not merely the accumulation of more changes, but the consolidation of diverse structural insights discovered across the search space. Although such depth remains an upper-bound measure of structural capacity rather than a direct indicator of quality, aligning depth accumulation with the compositional structure of natural tasks increases the likelihood that deeper workflows correspond to more effective reasoning pipelines when such structure is required.

#### D.4 Depth Comparison: AFlow vs. FUSIONFLOW

We compare the **structural depth** trajectories of AFlow (optimization-only) and FUSIONFLOW (with fusion) on the MBPP dataset over 20 optimization rounds. The results are presented in Table 5 and visualized in Figure 5.

##### Key Observations.

- Modification Aggregation:** AFlow’s depth is bounded by the maximum chain length in the search tree, reaching only depth 4. FUSIONFLOW achieves a maximum depth of 7 by aggregating modifications from multiple evolutionary branches.
- Fusion as Branch Merger:** Fusion takes the union of modification sets from parent workflows, effectively combining insights from different search trajectories. For example:
  - Round 6: Fuses workflows from rounds 2, 3, 5, merging their modification histories into {1, 2, 3, 6}.
  - Round 15: Fuses rounds 3, 9, 14, combining {1, 2, 3} ∪ {1, 2, 8, 9} ∪ {1, 2, 8, 9, 14} ∪ {15} = {1, 2, 3, 8, 9, 14, 15}.

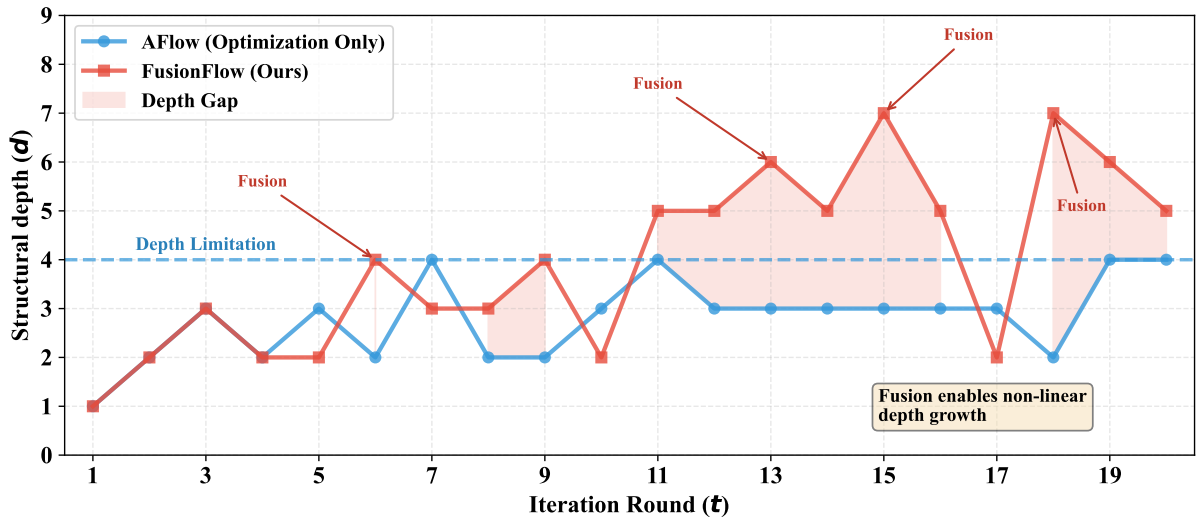


Figure 5: *Structural depth* comparison between AFlow (blue) and **FUSIONFLOW** (red) over 20 optimization rounds on the MBPP dataset. Depth is defined as the number of effective modifications contained in each workflow. Fusion operations (annotated) aggregate modifications from multiple evolutionary branches, enabling **FUSIONFLOW** to reach deeper regions of the modification space.

- Round 18: Fuses rounds 9, 14, 3, yielding {1, 2, 3, 8, 9, 14, 18}.

3. **Depth Advantage:** The maximum depth ratio is  $7/4 = 1.75\times$ , demonstrating that fusion enables access to deeper regions of the modification space by combining complementary evolutionary branches.
4. **Differentiation Preserves Diversity:** Differentiation operations (rounds 4, 5, 10, 17) create low-depth specialized variants that serve as complementary building blocks for subsequent fusion, without introducing new modifications.

## E.1 Global Dynamic Scheduling Algorithm

Table 6: Global Dynamic Scheduling Algorithm

**Algorithm 1:** Global Dynamic Scheduling Algorithm**Input:** Initial workflow pool  $\mathcal{P}_0$ , Iteration budget  $T$ , Dataset  $\mathcal{D}$ **Parameters:** Window size  $k$ , Sensitivity  $\kappa$ **Output:** Optimized workflow pool  $\mathcal{P}_{T+1}$ 


---

```

1: Initialize  $\mathcal{P}_1 \leftarrow \mathcal{P}_0$ 
2: for  $t = 1$  to  $T$  do
3:   Evaluate all  $\mathcal{W}_i \in \mathcal{P}_t$  on  $\mathcal{D}$ , obtain  $R_i$  and coverage sets  $\mathcal{C}_i$ 
4:    $R_t \leftarrow \max_{\mathcal{W}_i \in \mathcal{P}_t} R_i$  ▷ Compute best score
5:   if  $t \geq 2k$  then ▷ Check sufficient history
6:      $R_t^{\text{recent}} \leftarrow \max_{u \in [t-k+1, t]} R_u$ 
7:      $R_t^{\text{prev}} \leftarrow \max_{u \in [t-2k+1, t-k]} R_u$ 
8:      $\Delta_t \leftarrow R_t^{\text{recent}} - R_t^{\text{prev}}$  ▷ Compute performance gap
9:      $\pi_t \leftarrow \frac{1}{1 + \exp(\kappa \cdot \Delta_t)}$  ▷ Sigmoid stagnation signal
10:  else
11:     $\pi_t \leftarrow 0$ 
12:  end if
13:  Sample operator  $\mathcal{O}_t \in \{\mathcal{O}_{\text{opt}}, \mathcal{O}_{\text{diff}}, \mathcal{O}_{\text{fuse}}\}$  ▷ Sample operator
14:  if  $\mathcal{O}_t = \mathcal{O}_{\text{fuse}}$  then ▷ Apply fusion branch
15:     $\mathcal{S}^* \leftarrow \text{FUSIONWORKFLOWSELECTION}(\mathcal{P}_t)$  ▷ Algorithm 2
16:     $\mathcal{W}_{\text{fuse}} \leftarrow \mathcal{O}_{\text{fuse}}(\mathcal{S}^*)$ 
17:     $\mathcal{P}_{t+1} \leftarrow \mathcal{P}_t \cup \{\mathcal{W}_{\text{fuse}}\}$ 
18:  else if  $\mathcal{O}_t = \mathcal{O}_{\text{diff}}$  then ▷ Apply differentiation branch
19:     $(\mathcal{W}^*, k^*) \leftarrow \text{DIFFERENTIATIONWORKFLOWSELECTION}(\mathcal{P}_t, \{\mathcal{D}_k\})$  ▷ Algorithm 3
20:     $\mathcal{W}_{\text{diff}} \leftarrow \mathcal{O}_{\text{diff}}(\mathcal{W}^*, k^*)$ 
21:     $\mathcal{P}_{t+1} \leftarrow \mathcal{P}_t \cup \{\mathcal{W}_{\text{diff}}\}$ 
22:  else ▷ Apply optimization branch
23:    Select  $\mathcal{W} \in \mathcal{P}_t$  for local optimization
24:     $\mathcal{W}_{\text{opt}} \leftarrow \mathcal{O}_{\text{opt}}(\mathcal{W})$ 
25:     $\mathcal{P}_{t+1} \leftarrow \mathcal{P}_t \cup \{\mathcal{W}_{\text{opt}}\}$ 
26:  end if
27: end for

```

---

## E.2 Fusion Workflow Selection Algorithm

Table 7: Fusion Workflow Selection Algorithm

---

### Algorithm 2: Fusion Workflow Selection

---

**Input:** Current workflow pool  $\mathcal{P}_t$

**Parameters:** Selection count  $M$ , fusion group size  $m$ , weights  $\lambda_1, \lambda_2, \mu_1, \mu_2$

**Output:** Selected workflow group  $\mathcal{S}^*$

---

```

1:  $\mathcal{P}' \leftarrow \text{TOPMBYCOVERAGE}(\mathcal{P}_t, M)$  ▷ Filter top- $M$  by  $|\mathbf{C}_i|$ 
2:  $\Phi_{\text{best}} \leftarrow -\infty$ 
3:  $\mathcal{S}^* \leftarrow \emptyset$ 
4: for each candidate group  $\mathcal{S} = \{\mathcal{W}_1, \dots, \mathcal{W}_m\} \subset \mathcal{P}'$  do
5:   Retrieve coverage sets  $\mathbf{C}_1, \dots, \mathbf{C}_m$  ▷  $\mathbf{C}_i \subseteq \mathcal{D}$ : solved instances by  $\mathcal{W}_i$ 
6:    $\bar{U}_{\text{pair}} \leftarrow \frac{1}{m(m-1)} \sum_{i \neq j} |\mathbf{C}_i \cup \mathbf{C}_j|$  ▷ Pairwise complementarity (union)
7:    $U(\mathcal{S}) \leftarrow |\bigcup_{i=1}^m \mathbf{C}_i|$  ▷ Group-level complementarity
8:    $\bar{I}_{\text{pair}} \leftarrow \frac{1}{m(m-1)} \sum_{i \neq j} |\mathbf{C}_i \cap \mathbf{C}_j|$  ▷ Pairwise consensus (intersection)
9:    $I(\mathcal{S}) \leftarrow |\bigcap_{i=1}^m \mathbf{C}_i|$  ▷ Group-level consensus
10:   $\Phi_{\text{fuse}}(\mathcal{S}) \leftarrow (\lambda_1 \bar{U}_{\text{pair}} + \lambda_2 U(\mathcal{S})) + (\mu_1 \bar{I}_{\text{pair}} + \mu_2 I(\mathcal{S})) \cdot \eta(\mathcal{S})$  ▷ Fusion potential
11:  if  $\Phi_{\text{fuse}}(\mathcal{S}) > \Phi_{\text{best}}$  then
12:     $\Phi_{\text{best}} \leftarrow \Phi_{\text{fuse}}(\mathcal{S})$ 
13:     $\mathcal{S}^* \leftarrow \mathcal{S}$ 
14:  end if
15: end for
16: return  $\mathcal{S}^*$ 

```

---

Table 8: Differentiation Workflow Selection Algorithm

**Algorithm 3:** Differentiation Workflow Selection**Input:** Current workflow pool  $\mathcal{P}_t$ , dataset partition  $\{\mathcal{D}_k\}_{k \in \mathcal{K}}$ **Parameters:** Softmax temperature  $\tau$ **Output:** Selected workflow  $\mathcal{W}^*$  and target cluster  $k^*$ 


---

```

1:  $N \leftarrow |\mathcal{D}|$  ▷  $N = \sum_{k \in \mathcal{K}} |\mathcal{D}_k|$ 
2: Initialize score list  $\mathcal{Q} \leftarrow []$ 
3: for each workflow  $\mathcal{W}_i \in \mathcal{P}_t$  do
4:   Retrieve  $\mathbf{C}_i \subseteq \mathcal{D}$  ▷ Solved instances of  $\mathcal{W}_i$ 
5:    $c_i \leftarrow |\mathbf{C}_i|$ 
6:   if  $c_i = 0$  then
7:     Append  $(\mathcal{W}_i, 0)$  to  $\mathcal{Q}$ 
8:     continue
9:   end if
10:   $\text{Acc}_{\text{global}}(\mathcal{W}_i) \leftarrow \frac{|\mathbf{C}_i|}{N}$  ▷ Global accuracy
11:   $S_i^{\text{max}} \leftarrow -\infty$ 
12:  for each cluster  $k \in \mathcal{K}$  do
13:     $n_k \leftarrow |\mathcal{D}_k|$ 
14:    if  $n_k = 0$  then
15:      continue
16:    end if
17:     $\mathbf{C}_{i,k} \leftarrow \mathbf{C}_i \cap \mathcal{D}_k$ 
18:     $c_{i,k} \leftarrow |\mathbf{C}_{i,k}|$ 
19:     $\text{Recall}_{i,k} \leftarrow \frac{c_{i,k}}{n_k}$  ▷ Cluster-level recall
20:    if  $\text{Recall}_{i,k} > \text{Acc}_{\text{global}}(\mathcal{W}_i)$  then
21:       $\text{Cov}_{i,k} \leftarrow \frac{c_{i,k}}{N}$  ▷ Absolute coverage
22:       $S_{i,k} \leftarrow \text{Cov}_{i,k} \cdot (\text{Recall}_{i,k} - \text{Acc}_{\text{global}}(\mathcal{W}_i))$  ▷ Split potential on cluster  $k$ 
23:       $S_i^{\text{max}} \leftarrow \max(S_i^{\text{max}}, S_{i,k})$ 
24:    end if
25:  end for
26:   $S_i \leftarrow \max(0, S_i^{\text{max}})$  ▷ Workflow differentiation potential
27:  Append  $(\mathcal{W}_i, S_i)$  to  $\mathcal{Q}$ 
28: end for
29: Sample  $(\mathcal{W}^*, S^*)$  from  $\mathcal{Q}$  with probability  $\propto \exp(S_i/\tau)$  ▷ Softmax selection
30:  $k^* \leftarrow \arg \max_{k \in \mathcal{K}} \left[ \frac{|\mathbf{C}_{*,k}|}{N} \cdot \left( \frac{|\mathbf{C}_{*,k}|}{|\mathcal{D}_k|} - \frac{|\mathbf{C}_{*,k}|}{N} \right) \right]$  ▷ Choose specialization target
31: return  $(\mathcal{W}^*, k^*)$ 

```

---

## F Example of Adaptive Operator Scheduling

This section demonstrates the dynamic scheduling mechanism in action through a complete 20-round optimization trace on the DROP dataset, showing how stagnation signals and operator probabilities evolve throughout the process. Note that this trace is from a different experimental run than the case study in Appendix I.

### F.1 Execution Trace: 20 Rounds

Table 9 shows the actual values from a complete optimization run.

Table 9: Evolution of Stagnation Signal and Operator Selection Probabilities over 20 Optimization Rounds

Round	$\pi_t$	$P_{\text{opt}}$	$P_{\text{diff}}$	$P_{\text{fuse}}$	Selected
1	0.0000	1.0000	0.0000	0.0000	Optimization
2	0.0000	1.0000	0.0000	0.0000	Optimization
3	0.4137	0.5450	0.2482	0.2068	Differentiation
4	1.0000	0.0000	0.5206	0.4794	Differentiation
5	0.7401	0.2664	0.3636	0.3701	Fusion
6	0.5863	0.4467	0.2880	0.2653	Optimization
7	0.3323	0.6864	0.1633	0.1504	Optimization
8	0.4137	0.6096	0.2032	0.1871	Optimization
9	0.2599	0.7547	0.1277	0.1176	Differentiation
10	0.3323	0.7019	0.1477	0.1504	Optimization
11	0.8512	0.2365	0.3784	0.3851	Optimization
12	0.5863	0.4741	0.2606	0.2653	Fusion
13	0.1986	0.8305	0.0883	0.0813	Optimization
14	0.5000	0.5731	0.2222	0.2047	Fusion
15	0.5000	0.5925	0.2222	0.1852	Optimization
16	0.8512	0.3063	0.3784	0.3153	Differentiation
17	0.8512	0.3423	0.3424	0.3153	Fusion
18	0.1488	0.8903	0.0598	0.0499	Optimization
19	0.1098	0.9191	0.0442	0.0368	Optimization

### F.2 Observable Patterns

The trace reveals several characteristic behaviors:

- **Initial exploration (Rounds 1–2):**  $\pi_t = 0 \Rightarrow P_{\text{opt}} = 1.0$  (pure exploitation)
- **First plateau (Rounds 3–5):**  $\pi_t$  increases from 0.41  $\rightarrow$  1.0  $\rightarrow$  0.74, triggering Differentiation and Fusion
- **Recovery phase (Rounds 6–10):**  $\pi_t$  drops to 0.26–0.59, returning to Optimization dominance
- **Second plateau (Rounds 11–12):**  $\pi_t = 0.85$  triggers Fusion
- **Balanced adaptation (Rounds 13–19):**  $\pi_t$  varies between 0.11 and 0.85, showing continuous adjustment

Fusion operations occur at Rounds 5, 12, 14, and 17, precisely when  $\pi_t \in [0.50, 0.85]$ , demonstrating adaptive triggering of structural exploration.

### F.3 Visual Illustration

Figure 6 shows how  $\pi_t$  evolves over time, with peaks at Rounds 4–5, 11–12, and 16–17 corresponding to performance plateaus.

Figure 7 visualizes the operator probability dynamics, showing the complementary relationship between Optimization (exploitation) and Differentiation/Fusion (exploration).

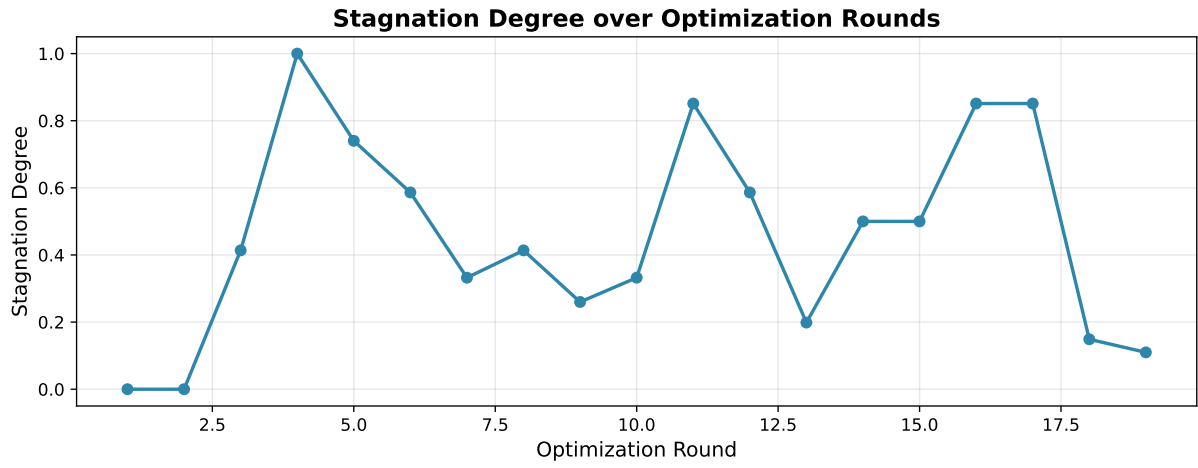


Figure 6: Stagnation signal  $\pi_t$  over 20 rounds. Peaks indicate performance plateaus.

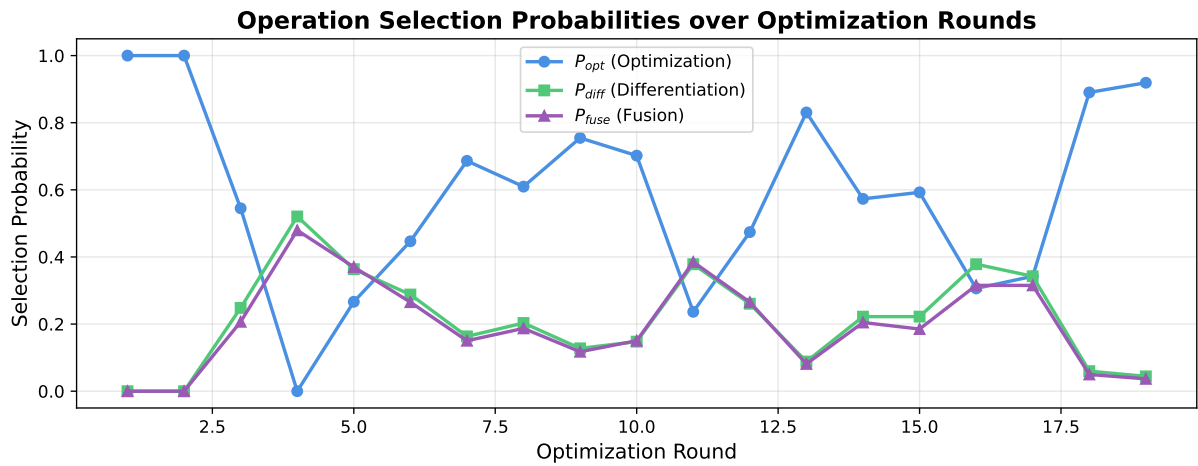


Figure 7: Operator selection probabilities. Blue: Optimization, Green: Differentiation, Purple: Fusion.

## G Framework Architecture Details

This section provides detailed examples of the core components in our **FUSIONFLOW** framework, including workflow structures, prompts, operators, and optimization metadata.

### G.1 Workflow Structure Example

A workflow in **FUSIONFLOW** consists of a directed graph that orchestrates various operators to solve specific problems. Below is a complete example of a workflow for solving GSM8K mathematical reasoning problems:

```
1581 from typing import Literal
1582 import workspace.GSM8K.workflows.template.operator as operator
1583 import workspace.GSM8K.workflows.round_1.prompt as prompt_custom
1584 from scripts.async_llm import create_llm_instance
1585 from scripts.evaluator import DatasetType
1586
1587
1588 class Workflow:
1589     def __init__(
1590         self,
1591         name: str,
1592         llm_config,
1593         dataset: DatasetType,
1594     ) -> None:
1595         self.name = name
1596         self.dataset = dataset
1597         self.llm = create_llm_instance(llm_config)
1598         self.custom = operator.Custom(self.llm)
1599
1600     async def __call__(self, problem: str):
1601         """
1602         Implementation of the workflow
1603         """
1604         solution = await self.custom(
1605             input=problem,
1606             instruction=""
1607         )
1608         return solution['response'], \
1609             self.llm.get_usage_summary()["total_cost"]
```

Listing 1: Example Workflow Graph for GSM8K Dataset

This simple workflow demonstrates the basic structure where:

- The workflow is initialized with an LLM configuration
- A Custom operator is instantiated to handle problem-solving
- The `__call__` method defines the execution flow
- The workflow returns both the solution and the cost incurred

### G.2 Node Implementation Example

Operators are the fundamental building blocks in **FUSIONFLOW**. Here we show the implementation of the Custom operator and the Programmer operator:

```
1619 from scripts.operators import Operator
1620
1621
1622 class Custom(Operator):
1623     def __init__(self, llm: AsyncLLM, name: str = "Custom"):
1624         super().__init__(llm, name)
1625
1626     async def __call__(self, input, instruction):
1627         prompt = instruction + input
```

```

response = await self._fill_node(
    GenerateOp,
    prompt,
    mode="single_fill"
)
return response

```

1628  
1629  
1630  
1631  
1632  
1633

Listing 2: Custom Operator Implementation

The Custom operator provides flexible problem-solving capability by combining custom instructions with input problems. It serves as the most versatile operator in the framework.

1635  
1636

```

class Programmer(Operator):
    def __init__(self, llm: AsyncLLM, name: str = "Programmer"):
        super().__init__(llm, name)

    async def exec_code(self, code, timeout=30):
        """
        Asynchronously execute code and return error if timeout.
        """
        loop = asyncio.get_running_loop()
        with concurrent.futures.ProcessPoolExecutor() as executor:
            try:
                future = loop.run_in_executor(
                    executor, run_code, code
                )
                result = await asyncio.wait_for(
                    future, timeout=timeout
                )
                return result
            except asyncio.TimeoutError:
                executor.shutdown(wait=False, cancel_futures=True)
                return "Error", "Code execution timed out"
            except Exception as e:
                return "Error", f"Unknown error: {str(e)}"

```

1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660

Listing 3: Programmer Operator with Code Execution

The Programmer operator is designed for problems requiring code execution, such as mathematical calculations or algorithmic solutions.

1662  
1663

### G.3 Prompt Template Example

1664

Prompts guide the LLM's behavior within operators. In **FUSIONFLOW**, prompts are stored separately and can be dynamically modified during optimization:

1665  
1666

```

# File: workspace/GSM8K/workflows/round_1/prompt.py

SOLVE_PROMPT = """
Please solve the following mathematical problem step by step.
Show your reasoning clearly and provide the final numerical answer.

Problem: {problem}

Solution:
"""

VERIFICATION_PROMPT = """
Verify the following solution and check if the answer is correct.
If there are errors, provide the corrected solution.

Original Problem: {problem}
Solution to Verify: {solution}

Verification:
"""

```

1667  
1668  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687

Listing 4: Prompt Template File Structure

1689 Custom prompts can be referenced in the workflow using `prompt_custom.PROMPT_NAME`.

#### 1690 G.4 Experience Metadata

1691 During optimization, **FUSIONFLOW** maintains a structured experience database that tracks the success  
1692 and failure of different modifications. Below is an example from the DROP dataset:

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727

1728

1729

1730

1731

1732

1733

1734

1735

1736

1737

1738

1739

1740

```

{
  "3": {
    "score": 0.8450,
    "success": {
      "14": {
        "modification": "3-way Workflow Fusion: Combined
          3 high-performing workflows from rounds [3, 5, 11]
          (scores: ['0.8450', '0.8450', '0.8300']). Integrated
          complementary strengths to maximize problem coverage.",
        "score": 0.8800
      }
    },
    "failure": {
      "15": {
        "modification": "Add mathematical verification and
          completeness check step using Custom operator before
          ensemble selection",
        "score": 0.8200
      }
    }
  },
  "2": {
    "score": 0.8400,
    "success": {
      "3": {
        "modification": "Add ScEnsemble operator to generate
          multiple reasoning paths and select the most consistent
          answer, improving accuracy for mathematical and logical
          reasoning problems",
        "score": 0.8450
      }
    },
    "failure": {}
  },
  "11": {
    "score": 0.8300,
    "success": {},
    "failure": {
      "17": {
        "modification": "Workflow Differentiation: Specialized
          for Mathematical Derivation but with too narrow focus",
        "score": 0.8100
      }
    }
  }
}

```

Listing 5: Processed Experience Metadata from DROP Dataset

1741

The experience metadata captures:

1742

- **Parent Round:** The round from which modifications originated

1743

- **Success/Failure Tracking:** Categorization of descendant rounds

1744

- **Modification Description:** Natural language description of the change

- **Performance Score:** Validation set accuracy

1745

This structured experience allows the optimization LLM to learn from past attempts and avoid repeating failed modifications.

1746

1747

## G.5 Fusion Metadata Example

1748

When performing workflow fusion, **FUSIONFLOW** combines multiple high-performing workflows. The fusion process generates metadata describing how workflows were combined:

1749

1750

```

{
  "fusion_round": 14,
  "parent_workflows": [
    {
      "round": 3,
      "score": 0.8450,
      "solved_problems": 169,
      "key_features": [
        "ScEnsemble for multiple reasoning paths",
        "Majority voting mechanism"
      ]
    },
    {
      "round": 5,
      "score": 0.8450,
      "solved_problems": 169,
      "key_features": [
        "Alternative ensemble configuration",
        "Enhanced problem comprehension"
      ]
    },
    {
      "round": 11,
      "score": 0.8300,
      "solved_problems": 166,
      "key_features": [
        "Specialized for mathematical derivation",
        "Step-by-step calculation verification"
      ]
    }
  ],
  "fusion_strategy": "Adaptive routing to specialized paths
    based on problem type analysis",
  "combined_coverage": 176,
  "fusion_score": 0.8800,
  "improvement_over_best_parent": 0.0350,
  "unique_problems_solved": 7
}

```

1751

1752

1753

1754

1755

1756

1757

1758

1759

1760

1761

1762

1763

1764

1765

1766

1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

1777

1778

1779

1780

1781

1782

1783

1784

1785

1786

1787

1788

1789

1790

Listing 6: Fusion Metadata Structure from DROP Dataset

## G.6 Differentiation Metadata Example

1791

Differentiation creates specialized variants of workflows by analyzing error logs:

1792

```

{
  "parent_round": 9,
  "parent_score": 0.8200,
  "error_analysis": {
    "total_errors": 36,
    "error_categories": {
      "multi_step_calculation_error": 15,
      "passage_comprehension_error": 12,
      "numerical_extraction_error": 9
    }
  },
  "example_errors": [

```

1793

1794

1795

1796

1797

1798

1799

1800

1801

1802

1803

1804

```

1805     {
1806         "problem_id": "drop-validation-127",
1807         "question": "How many more points did team A score
1808             in the first half than team B?",
1809         "passage": "Team A scored 24 in first half...
1810             Team B scored 17 in first half...",
1811         "expected": "7",
1812         "got": "24",
1813         "error_type": "multi_step_calculation_error",
1814         "root_cause": "Extracted first value instead of
1815             computing difference"
1816     }
1817 ]
1818 },
1819 "differentiation_type": "mathematical_specialization",
1820 "modification": "Enhance workflow with dedicated multi-step
1821     calculation handling and intermediate value
1822     tracking for mathematical derivation problems",
1823 "child_round": 11,
1824 "child_score": 0.8300,
1825 "improvement": 0.0100,
1826 "specialization_focus": "Complex multi-step numerical reasoning"
1827 }

```

Listing 7: Differentiation Metadata from DROP Dataset

## 1829 G.7 Log Entry Example

1830 Execution logs capture detailed information about each problem solved:

```

1831 {
1832     "question": "How many more field goals did the team make
1833         in the first half compared to the second half?",
1834     "passage": "In the first half, the team made 3 field goals.
1835         ... In the second half, they made 1 field goal.",
1836     "right_answer": "2",
1837     "model_output": "To find how many more field goals the team
1838         made in the first half:\n1. First half field goals: 3\n
1839         2. Second half field goals: 1\n3. Difference: 3 - 1 = 2\n
1840         Therefore, the team made 2 more field goals in the first
1841         half compared to the second half.",
1842     "score": 1.0,
1843     "judge_explanation": "Model correctly identified the values
1844         and computed the difference, arriving at answer '2'.",
1845     "problem_id": "drop-validation-042",
1846     "category": "Numerical Reasoning"
1847 }
1848

```

Listing 8: Example Log Entry from DROP Workflow Execution

1850 These logs are used by the differentiation operator to identify systematic errors and guide targeted  
1851 improvements.

## H Core Operator Implementations Logic

1852

This section details the prompts and implementation logic for FUSIONFLOW's three core optimization operators: Optimization, Fusion, and Differentiation.

1853

1854

### H.1 Optimization Operator

1855

The optimization operator analyzes existing workflows and proposes incremental improvements.

1856

#### H.1.1 Optimization Prompt

1857

##### Prompt for the Optimization Operator

###### [Objective]

You are building a Graph and corresponding Prompt to jointly solve {type} problems. Referring to the given graph and prompt, which forms a basic example of a {type} solution approach, please reconstruct and optimize them. You can add, modify, or delete nodes, parameters, or prompts. Include your single modification in XML tags in your reply. Ensure they are complete and correct to avoid runtime failures.

###### [Optimization Guidelines]

When optimizing, you can incorporate critical thinking methods like:

- **Review:** Add verification steps to check intermediate results
- **Revise:** Implement error correction loops
- **Ensemble:** Generate multiple answers through different/similar prompts, then vote/integrate/check the majority to obtain a final answer
- **Self-Ask:** Break down complex problems into sub-questions

Consider Python's loops (for, while, list comprehensions), conditional statements (if-elif-else, ternary operators), or machine learning techniques. The graph complexity should not exceed 10 nodes. Use logical and control flow (IF-ELSE, loops) for enhanced graphical representation.

###### [Input Format]

You will receive:

```
<sample>
  <experience>{experience}</experience>
  <modification>(e.g., add/delete/modify operator X)</modification>
  <score>{score}</score>
  <graph>{current_graph_code}</graph>
  <prompt>{current_prompts}</prompt>
  <operator_description>{available_operators}</operator_description>
</sample>
```

Error logs from current workflow:

```
{error_examples}
```

###### [Constraints]

- **Single Modification:** Only one detail point can be modified at a time
- **Code Limit:** No more than 5 lines of code may be changed per modification
- **No Placeholders:** Generated prompts must not contain any placeholders

1858

- **Format Output:** Use custom methods to restrict output format, not code
- **No None Output:** Graph must never output None for any field

### [Custom Operator Usage]

Example of using the Custom operator in your graph:

```
# Write prompt in prompt_custom and use in Custom method
response = await self.custom(
    input=problem,
    instruction=prompt_custom.XXX_PROMPT
)

# Can concatenate previous results for more context
solution = await self.generate(
    problem=f"question:{problem}, context:{response['response']}"
)
```

Note: In custom, input and instruction are directly concatenated (instruction+input), and placeholders are not supported.

### [Output Format]

Provide your optimization in XML format:

```
<modification>
Describe the single modification made
</modification>
<graph>
Complete Python code for optimized workflow
</graph>
<prompt>
All necessary prompts in prompt_custom (if needed)
</prompt>
```

## H.1.2 Optimization Implementation Logic

The optimization operator follows this process:

1. **Select Parent:** Choose high-performing workflows
2. **Load Context:** Retrieve the parent workflow's code, prompts, and error logs
3. **Format Experience:** Convert historical experience into structured format showing successful and failed modifications
4. **Generate Prompt:** Create optimization prompt with all context
5. **Call LLM:** Use optimization LLM to generate modified workflow
6. **Validate:** Check that modification is novel (not repeated from history)
7. **Evaluate:** Test the modified workflow on validation set
8. **Record:** Update experience database with success/failure

## H.2 Fusion Operator

The fusion operator combines multiple high-performing workflows into a single comprehensive workflow.

### Prompt for the Fusion Operator

#### [Objective]

You are an expert workflow designer tasked with fusing multiple high-performing workflows into a single, more comprehensive workflow that can solve {type} problems. Referring to the given graphs and prompts from the parent workflows, you must intelligently combine the strengths of each input workflow while creating robust routing logic to handle different problem types.

#### [Core Principle]

**The most crucial point:** Your task is to perform multi-way workflow fusion. This means you must identify the key components within each of the parent workflows and select only those parts that have fusion value to combine them into a new workflow. **You should not create the core components from scratch**, including graph structure and prompts.

#### [Input Format]

You will be given multiple high-performing workflows ( $K \geq 2$ ):

```
<workflow_1>
  <round>{round_number}</round>
  <score>{performance_score}</score>
  <solved_problems>{count} problems</solved_problems>
  <prompt>{workflow_prompts}</prompt>
  <graph>{workflow_code}</graph>
</workflow_1>
... (workflow_2, workflow_3, ... similarly)
```

```
<operator_description>
{available_operators}
</operator_description>
```

#### [Fusion Strategy]

Analyze each workflow's strengths and design a fusion strategy that:

- **Maximizes Coverage:** The fused workflow should solve problems that any of the parent workflows can solve
- **Leverages Scores:** Workflows with higher scores can serve as the foundation
- **Intelligent Routing:** Create decision logic to route problems to appropriate solution paths
- **Combines Strengths:** Integrate complementary capabilities from each parent

#### [Design Patterns]

Common fusion patterns include:

##### 1. Problem Classification + Routing:

```
problem_type = await self.classify(problem)
if problem_type == "calculation":
    return await workflow_1_logic(problem)
elif problem_type == "word_problem":
    return await workflow_2_logic(problem)
```

## 2. Sequential Enhancement:

```
# Use workflow_1 for initial solution
initial = await workflow_1_solve(problem)
# Use workflow_2 for verification
verified = await workflow_2_verify(initial)
# Use workflow_3 for refinement if needed
if not verified:
    final = await workflow_3_refine(initial)
```

## 3. Parallel Ensemble:

```
results = await asyncio.gather(
    workflow_1_solve(problem),
    workflow_2_solve(problem),
    workflow_3_solve(problem)
)
final = self.vote_or_combine(results)
```

### [Constraints]

- The fused workflow must be different from any single input workflow
- Generated prompts must not contain placeholders
- Use custom methods to restrict output format, not code

### [Output Format]

Provide your fusion result in XML format:

```
<modification>
Description of the fusion strategy and how the parent workflows
are combined
</modification>
<graph>
Complete Python code for the fused workflow
</graph>
<prompt>
All necessary prompts in prompt_custom (if needed)
</prompt>
```

## H.2.2 Fusion Implementation Logic

The fusion operator executes the following workflow:

1. **Select Candidate Workflows:** Choose high-performing workflows based on the fusion potential score
2. **Analyze Diversity:** Compute solved problem overlap to ensure complementarity
3. **Load Components:** Extract code, prompts, and metadata from selected workflows

4. **Format Input:** Structure the parent workflows into fusion prompt format 1882
5. **Generate Fusion:** Call LLM to create combined workflow 1883
6. **Validate Integration:** Verify that parent workflows contributed to the fusion 1884
7. **Evaluate Performance:** Test on validation set 1885
8. **Update Metadata:** Record fusion genealogy and parent information 1886

### H.3 Differentiation Operator 1887

The differentiation operator creates specialized variants of workflows targeting particular task subspaces (e.g., algebraic problems, geometric reasoning, combinatorial counting). Unlike optimization which refines general performance, differentiation tailors workflows to excel on specific problem categories, typically through prompt specialization and targeted structural adjustments. 1888  
1889  
1890  
1891

#### H.3.1 Differentiation Prompt 1892

##### Prompt for the Differentiation Operator

##### [Objective]

You are an expert workflow designer tasked with differentiating an existing workflow to create a specialized variant for {type} problems.

##### [Core Principle - What is Differentiation]

Differentiation means creating a specialized variant that excels on a particular problem category:

- Focus on making the workflow highly effective for the target task subspace
- Specialize prompts with domain-specific instructions and reasoning patterns
- Structural adjustments are allowed when beneficial for specialization

##### [What You ARE Doing]

- Converting “analyze problem” → “analyze [specific problem type] characteristics”
- Converting “solve” → “apply [specific technique] to solve”
- Converting “verify” → “verify using [specific domain] validation rules”
- Making prompts more specific and detailed for the target domain

##### [Input Format]

You will receive:

```
<differentiation_data>
  <dataset>{dataset}</dataset>
  <target_round>{target_round}</target_round>
  <specialization_direction>{direction}</specialization_direction>
  <specialization_focus>{specialization_focus}</specialization_focus>
  <parent_workflow>
    <score>{validation_score}</score>
    <graph>{workflow_code}</graph>
    <prompt>{workflow_prompts}</prompt>
  </parent_workflow>
  <operator_description>{available_operators}</operator_description>
</differentiation_data>
```

1893

For targeted differentiation, a <target\_category> section with example problems may also be provided.

#### [How to Convert Abstract → Concrete]

Parent: "analyze the problem"

Child: "analyze {specialization\_focus} characteristics:  
identify key variables, constraints, and target"

Parent: "solve the problem"

Child: "apply {specialization\_focus} techniques:  
[specific method 1], [specific method 2]"

Parent: "verify the solution"

Child: "verify using {specialization\_focus} rules:  
check [specific constraint 1], [specific constraint 2]"

#### [Example Differentiation]

If parent workflow is:

```
analysis = await self.custom(input=problem,  
    instruction=prompt_custom.ANALYSIS_PROMPT)  
solution = await self.custom(input=analysis['response'],  
    instruction=prompt_custom.SOLUTION_PROMPT)
```

Differentiated workflow should be:

```
analysis = await self.custom(input=problem,  
    instruction=prompt_custom.ALGEBRAIC_ANALYSIS_PROMPT)  
solution = await self.custom(input=analysis['response'],  
    instruction=prompt_custom.ALGEBRAIC_SOLUTION_PROMPT)
```

The prompts should contain {specialization\_focus}-specific instructions (e.g., "identify variables and equations", "apply substitution or elimination methods"), but the structure remains identical.

#### [Output Format]

```
<modification>  
Specialized for {specialization_focus}: {description}  
</modification>  
<graph>  
Differentiated workflow code (same structure, specific prompts)  
</graph>  
<prompt>  
Category-specific prompts  
</prompt>
```

1894

1895

### H.3.2 Differentiation Modes

1896

The differentiation operator supports two modes:

1897

**Exploratory Differentiation:** When no specific target category is provided, the system randomly selects a specialization focus from the predefined categories based on the problem type (math/code/qa). This enables broad exploration of the specialization space.

1898

1899

1900

1901

**Targeted Differentiation:** When a specific target\_category is provided (along with optional description and example problems), the differentiation becomes directed toward that particular task subspace.

The prompt includes additional context: 1902

```

<target_category> 1903
  <name>{category_name}</name> 1904
  <description>{category_description}</description> 1905
  <examples> 1906
    <example_1>{problem_text_1}</example_1> 1907
    <example_2>{problem_text_2}</example_2> 1908
    ... 1909
  </examples> 1910
</target_category> 1911

```

### H.3.3 Differentiation Implementation Logic 1912

The DifferentiationPromptGenerator class implements the following process: 1913

1. **Load Parent Workflow:** Retrieve a high-performing general workflow (graph and prompts) 1914
2. **Determine Specialization Focus:** 1915
  - If target\_category is provided: use it directly for targeted differentiation 1916
  - Otherwise: randomly select from predefined specializations based on question\_type 1917
3. **Gather Category Examples** (for targeted mode): Collect up to 3 representative problems, truncating if longer than 500 characters 1918  
1919
4. **Generate Specialization Prompt:** Combine the differentiation system prompt, input data, and custom operator usage instructions 1920  
1921
5. **Call LLM:** Request the model to convert abstract steps into category-specific steps while preserving workflow structure 1922  
1923
6. **Parse Output:** Extract the modified graph and category-specific prompts 1924
7. **Evaluate:** Test the differentiated workflow on the validation set 1925

## I Case Study: Complete Optimization Cycle

This section presents a detailed walkthrough of a complete optimization cycle on the DROP dataset, demonstrating how FUSIONFLOW progressively improves workflow performance through its three core operators. DROP (Discrete Reasoning Over Paragraphs) requires complex numerical reasoning over text passages, making it an ideal showcase for our framework’s capabilities.

**Important Note:** All scores reported in this case study are **validation set scores** (200 problems), which guide the optimization process. In our experimental protocol, we select the workflow with the best validation performance and evaluate it on the held-out test set (800 problems) for final reporting. The validation scores shown here may differ from the final test set scores reported in the main paper’s experiments.

### I.1 Initial Setup

**Dataset:** DROP (Discrete Reasoning Over Paragraphs)

**Problem Type:** Numerical reasoning and question answering

**Initial Workflow:** Simple Custom operator with no instructions

**Validation Set Size:** 200 problems

**Test Set Size:** 800 problems

**Challenge:** Requires multi-step mathematical reasoning combined with passage comprehension

#### I.1.1 Round 1: Baseline Workflow

The initial workflow starts with a minimal structure:

```
class Workflow:
    def __init__(self, name: str, llm_config, dataset: DatasetType):
        self.name = name
        self.dataset = dataset
        self.llm = create_llm_instance(llm_config)
        self.custom = operator.Custom(self.llm)

    async def __call__(self, problem: str):
        solution = await self.custom(input=problem, instruction="")
        return solution['response'], \
            self.llm.get_usage_summary()["total_cost"]
```

Listing 9: Round 1: Initial Baseline Workflow

#### Performance:

- Validation Score: 0.8200 (82.00%)
- Solved Problems: 164/200

**Error Analysis:** Identified issues with multi-step reasoning, mathematical calculations, and extracting the correct numerical values from passages.

### I.2 Optimization Operator Applied

#### I.2.1 Round 2: Adding Structured Reasoning

The optimization operator analyzes Round 1 and proposes adding structured step-by-step reasoning:

**Modification Description:** “Add AnswerGenerate operator to provide structured step-by-step reasoning before generating the final solution, which should help with mathematical calculations and logical comparisons”

```
class Workflow:
    def __init__(self, name: str, llm_config, dataset: DatasetType):
        self.name = name
        self.dataset = dataset
```

```

self.llm = create_llm_instance(llm_config)
self.custom = operator.Custom(self.llm)
self.answer_generate = operator.AnswerGenerate(self.llm)

async def __call__(self, problem: str):
    # Use AnswerGenerate for structured reasoning
    solution = await self.answer_generate(
        problem=problem,
        instruction=prompt_custom.REASONING_PROMPT
    )
    return solution['response'], \
        self.llm.get_usage_summary()["total_cost"]

```

Listing 10: Round 2: Workflow with Structured Reasoning

Associated prompt in prompt\_custom:

```

REASONING_PROMPT = """
Carefully read the passage and question.
Break down the problem into steps:
1. Identify relevant information from the passage
2. Determine what calculation is needed
3. Perform the calculation step by step
4. Provide the final numerical answer
"""

```

### Performance:

- Validation Score: 0.8400 (84.00%)
- Improvement: +0.0200 (+2.00%)
- Solved Problems: 168/200

**Impact:** Modest improvement through structured reasoning, helping with some multi-step problems but still struggling with complex cases.

### I.2.2 Round 3: Adding Ensemble for Robustness

Building on Round 2's success, another optimization adds ensemble voting:

**Modification:** "Add ScEnsemble operator to generate multiple reasoning paths and select the most consistent answer, improving accuracy for mathematical and logical reasoning problems"

```

class Workflow:
    def __init__(self, name: str, llm_config, dataset: DatasetType):
        self.name = name
        self.dataset = dataset
        self.llm = create_llm_instance(llm_config)
        self.custom = operator.Custom(self.llm)
        self.answer_generate = operator.AnswerGenerate(self.llm)
        self.sc_ensemble = operator.ScEnsemble(self.llm)

    async def __call__(self, problem: str):
        # Generate multiple reasoning paths
        solutions = await self.sc_ensemble(
            problem=problem,
            instruction=prompt_custom.REASONING_PROMPT,
            num_samples=3
        )
        return solutions['response'], \
            self.llm.get_usage_summary()["total_cost"]

```

Listing 11: Round 3: Workflow with Ensemble Voting

2028 **Performance:**

- 2029 • Validation Score: 0.8450 (84.50%)
- 2030 • Improvement over Round 2: +0.0050 (+0.50%)
- 2031 • Solved Problems: 169/200

2032 **Impact:** Ensemble voting provides marginal improvement, reducing some calculation errors but  
2033 introducing additional cost without dramatic performance gains.

2034 **I.2.3 Round 11: Differentiation-Based Specialization**

2035 The differentiation operator creates a specialized variant by converting abstract workflow steps into  
2036 concrete, domain-specific steps for a particular task subspace:

2037 **Modification:** “Workflow Differentiation: Specialized for Mathematical Derivation problems. The par-  
2038 ent workflow’s general prompts are converted into concrete, derivation-specific prompts while preserving  
2039 the exact same workflow structure.”

2040 **Differentiation Process:** The operator takes a high-performing parent workflow and specializes it for  
2041 the “mathematical derivation” category by:

- 2042 • Keeping the same operator structure (single Custom operator)
- 2043 • Converting abstract instructions (e.g., “solve the problem”) into concrete, domain-specific instructions  
2044 (e.g., “apply step-by-step mathematical derivation with explicit variable tracking”)
- 2045 • Adding category-specific reasoning patterns for numerical calculations

2046 **Performance:**

- 2047 • Validation Score: 0.8300 (83.00%)
- 2048 • Solved Problems: 166/200 (different subset than Round 3)
- 2049 • Specialization: Strong on problems requiring complex multi-step calculations

2050 **I.3 Fusion Operator Applied: The Breakthrough**

2051 **I.3.1 Round 14: Three-Way Workflow Fusion**

2052 After multiple rounds of optimization and differentiation, the fusion operator combines three complemen-  
2053 tary workflows to achieve a significant performance leap:

2054 **Parent Workflows:**

- 2055 • **Round 3** (Score: 0.8450): Strong ensemble-based reasoning
- 2056 • **Round 5** (Score: 0.8450): Alternative ensemble configuration
- 2057 • **Round 11** (Score: 0.8300): Specialized for mathematical derivation

2058 **Key Insight:** While Rounds 3 and 5 have similar overall scores, they solve different subsets of problems.  
2059 Round 11’s specialization complements both by handling complex calculation problems that the ensemble  
2060 approaches miss.

2061 **Fusion Strategy:** Create an adaptive workflow that routes problems based on their characteristics and  
2062 combines the strengths of all three parents.

```
2063 class Workflow:  
2064     def __init__(self, name: str, llm_config, dataset: DatasetType):  
2065         self.name = name  
2066         self.dataset = dataset  
2067         self.llm = create_llm_instance(llm_config)  
2068         self.custom = operator.Custom(self.llm)  
2069
```

```

self.answer_generate = operator.AnswerGenerate(self.llm)
self.sc_ensemble = operator.ScEnsemble(self.llm)

async def __call__(self, problem: str):
    # Step 1: Analyze problem complexity
    analysis = await self.custom(
        input=problem,
        instruction=prompt_custom.ANALYZE_PROBLEM_PROMPT
    )

    problem_type = analysis['response'].lower()

    # Step 2: Intelligent routing based on problem type
    if "multi-step calculation" in problem_type or \
        "complex derivation" in problem_type:
        # Use Round 11's specialized approach
        solution = await self._solve_with_derivation(problem)
    elif "comparison" in problem_type or \
        "multiple values" in problem_type:
        # Use Round 5's ensemble configuration
        solution = await self._solve_with_enhanced_ensemble(
            problem
        )
    else:
        # Use Round 3's standard ensemble approach
        solution = await self._solve_with_standard_ensemble(
            problem
        )

    # Step 3: Final verification and formatting
    final_answer = await self.custom(
        input=f"Question: {problem}\nAnswer: {solution}",
        instruction=prompt_custom.VERIFY_FORMAT_PROMPT
    )

    return final_answer['response'], \
        self.llm.get_usage_summary()["total_cost"]

async def _solve_with_derivation(self, problem: str):
    """Round 11's specialized mathematical derivation logic"""
    solution = await self.answer_generate(
        problem=problem,
        instruction=prompt_custom.DERIVATION_PROMPT
    )
    return solution['response']

async def _solve_with_enhanced_ensemble(self, problem: str):
    """Round 5's ensemble configuration"""
    solutions = await self.sc_ensemble(
        problem=problem,
        instruction=prompt_custom.ENHANCED_REASONING_PROMPT,
        num_samples=5
    )
    return solutions['response']

async def _solve_with_standard_ensemble(self, problem: str):
    """Round 3's standard ensemble approach"""
    solutions = await self.sc_ensemble(
        problem=problem,
        instruction=prompt_custom.REASONING_PROMPT,
        num_samples=3
    )
    return solutions['response']

```

2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105  
2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133

Listing 12: Round 14: Fused Workflow with Intelligent Routing

## Performance - The Breakthrough:

- Validation Score: 0.8800 (88.00%)

2134

2135

- **Improvement over best parent: +0.0350 (+3.50%)**
- **Solved Problems: 176/200**
- **New Problems Solved: 7 problems that none of the three parents could solve**
- **Coverage Enhancement: Union of parent coverage (169 + 169 + 166) → 176 unique problems through intelligent routing**

#### Analysis of the Breakthrough:

The fusion’s success demonstrates two key principles:

1. **Complementary Coverage:** Each parent excels on different problem types:
  - Round 3: Strong on standard reasoning problems (optimization result)
  - Round 5: Effective on comparison and multi-value problems (optimization result)
  - Round 11: Superior on complex mathematical derivations (differentiation result—specialized for the “mathematical derivation” task subspace)
2. **Intelligent Integration:** The fused workflow doesn’t simply combine approaches randomly, but routes each problem to the most suitable strategy, achieving **3.5% improvement** - significantly larger than the marginal gains from pure optimization (+0.5% from Round 2 to 3). This demonstrates the power of **deep trajectory crossing**.
3. **Differentiation Enables Fusion:** Round 11’s task subspace specialization is crucial—without differentiation creating a workflow specialized for mathematical derivations, the fusion would lack the diversity needed to achieve complementary coverage.

#### I.4 Performance Evolution Summary

Table 10 summarizes the performance evolution across the optimization cycle, highlighting the breakthrough achieved by fusion:

Table 10: Performance Evolution on DROP Dataset Across Optimization Rounds (Validation Set). \* Round 11’s lower overall score reflects task subspace specialization - it excels on mathematical derivation problems while trading off performance on other problem types

Round	Operator	Modification	Val Score	$\Delta$
1	Baseline	Initial workflow	0.8200	-
2	Optimization	Add structured reasoning	0.8400	+0.0200
3	Optimization	Add ensemble voting	0.8450	+0.0050
5	Optimization	Alternative ensemble	0.8450	+0.0000
11	Differentiation	Math specialization	0.8300	-0.0150*
<b>14</b>	<b>Fusion</b>	<b>3-way fusion [3+5+11]</b>	<b>0.8800</b>	<b>+0.0350</b>
<b>Total Improvement from Baseline</b>				<b>+0.0600</b>
<b>Fusion Contribution (beyond best parent)</b>				<b>+0.0350</b>
<b>Pure Optimization Gains (R1→R3)</b>				<b>+0.0250</b>

#### Key Observations:

- **Diminishing Returns from Optimization:** Pure optimization shows diminishing marginal returns (2.0% → 0.5% → 0.0%)
- **Differentiation Creates Diversity:** Round 11’s differentiation produces a specialized workflow with lower overall score (0.8300) but superior performance on mathematical derivation problems, creating valuable diversity for fusion

• <b>Fusion Breakthrough:</b> The fusion operator achieves a <b>3.5% improvement - 7x larger</b> than the gain from Round 2 to Round 3 (0.5%)	2164 2165
• <b>Deep Trajectory Crossing:</b> Fusion solves 7 additional problems that <i>none</i> of the three parents could solve individually, demonstrating the power of combining diverse optimization trajectories	2166 2167
• <b>Validation vs Test:</b> These scores are on the 200-problem validation set used to guide optimization. The final experiment selects the best validation workflow (Round 14) and evaluates it on the 800-problem test set for unbiased performance reporting.	2168 2169 2170
<b>I.5 Failed Attempts and Lessons Learned</b>	2171
Not all modifications succeed. Here are notable failures from the DROP case study:	2172
<b>I.5.1 Failed Attempt 1: Over-Verification (Round 15)</b>	2173
<b>Modification:</b> “Add a mathematical verification and completeness check step using Custom operator before the ensemble selection to validate calculations and ensure all qualifying answers are included”	2174 2175
<b>Result:</b> Score dropped to 0.8200 (regression of -0.0250 from parent Round 3)	2176
<b>Root Cause:</b> The verification step introduced excessive complexity and second-guessing, causing the model to reject correct answers or overcomplicate simple problems	2177 2178
<b>Lesson:</b> Verification should be selective, not universal. Pre-ensemble verification can interfere with the diversity needed for effective voting.	2179 2180
<b>I.5.2 Failed Attempt 2: Excessive Answer Refinement (Round 16)</b>	2181
<b>Modification:</b> “Add a dedicated answer refinement step using Custom operator after ensemble selection to standardize answer formats and verify calculations, ensuring consistency with expected answer patterns”	2182 2183
<b>Result:</b> Score: 0.7600 (severe regression of -0.0850)	2184
<b>Root Cause:</b> The post-processing step incorrectly reformatted many correct numerical answers, changing "175" to "one hundred seventy-five" when the expected format was numerical	2185 2186
<b>Lesson:</b> Answer formatting should preserve the content’s semantic meaning. Format standardization must be carefully designed to avoid introducing new errors.	2187 2188
<b>I.5.3 Failed Attempt 3: Misguided Ensemble Expansion (Round 18)</b>	2189
<b>Modification:</b> “Modify the existing Custom operator call to include mathematical verification and multi-answer handling within the formatting step, and increase the number of solution generations from 3 to 5 for better ensemble coverage”	2190 2191 2192
<b>Result:</b> Score: 0.8450 (no improvement over parent, just increased cost)	2193
<b>Root Cause:</b> Simply increasing ensemble size without improving the underlying reasoning doesn’t help. The additional samples added cost without providing new insights.	2194 2195
<b>Lesson:</b> Ensemble size should be tuned based on problem diversity. More samples help when solutions disagree, but don’t improve performance when the model consistently makes the same type of error.	2196 2197
<b>I.5.4 Failed Attempt 4: Suboptimal Fusion (Round 19 - Alternative)</b>	2198
<b>Modification:</b> “3-way Workflow Fusion: Combined 3 high-performing workflows from rounds [14, 18, 11]”	2199 2200
<b>Result:</b> Score: 0.8300 (below parent Round 14’s 0.8800)	2201
<b>Root Cause:</b> Round 18 (no improvement over its parent) didn’t contribute meaningful diversity. Fusing workflows with overlapping strengths doesn’t create complementary coverage.	2202 2203
<b>Lesson:</b> Successful fusion requires <b>diverse</b> parents with <b>complementary</b> strengths. The parent selection strategy should maximize coverage diversity, not just select the highest-scoring workflows.	2204 2205