

AgentSlimming: Towards Efficient and Cost-Aware Multi-Agent Systems

Anonymous ACL submission

Abstract

Large Language Model-based Multi-Agent Systems (MAS) have demonstrated remarkable capabilities in complex tasks. However, manually designing optimal communication topologies is labor-intensive, while automated expansion methods often result in bloated structures with redundant agents, leading to excessive token consumption. To address this problem, we introduce **AgentSlimming**, a plug-and-play compression framework for graph-structured multi-agent workflows. Motivated by the AgentPruner and AgentQuant in neural networks, AgentSlimming compresses workflows by firstly estimate the importance score of each agent with a hybrid mechanism, and then removing redundant agents or replacing them with low-cost ones, where each operation is then validated with a baseline-anchored acceptance rule to prevent performance collapse. Experiments show that AgentSlimming reduces average token cost by up to 78.9% with negligible performance degradation, and even sometimes improves accuracy, achieving a strong Pareto-optimal trade-off between cost and quality. *Our codes are available in the supplementary materials and will be released on Github.*

1 Introduction

The paradigm of Multi-Agent Systems (MAS) has shifted from manually crafting static prompts to orchestrating dynamic collaborations among specialized agents. Frameworks such as AutoGen (Wu et al., 2023) and ADAS (Hu et al., 2025) have shown that decomposing a complex problem into sub-tasks (Wei et al., 2022). Recent innovations, such as MetaGPT (Hong et al., 2024), AFlow (Zhang et al., 2025c), have revolutionized this field by reformulating workflow generation as a search problem. By leveraging Monte Carlo Tree Search (MCTS), AFlow can automatically navigate the vast space of agent interactions to discover highly effective reasoning topologies.

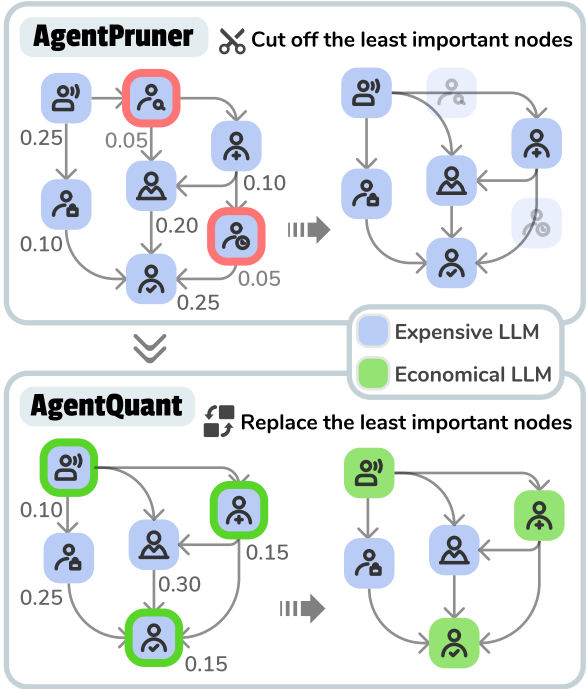


Figure 1: Overview of the pruning and quantization concepts in AgentSlimming. The process begins with an initially high-performance but computationally expensive workflow. Then, a hybrid importance evaluation mechanism is utilized to calculate each agent node’s importance score (i.e., the float number), which guides pruning and quantization in AgentSlimming.

However, this automated discovery process operates under an "unconstrained resource" assumption. The resulting workflows, while accurate, tend to be computationally bloated. They often feature dense, redundant connectivity and uniformly employ the most capable and expensive LLMs (e.g., GPT-4) as agent nodes for every sub-task (Li et al., 2023; Qian et al., 2024). This leads to two critical bottlenecks: **redundant communication**, where agents exchange repetitive or low-value information, and **excessive computational cost**, where simple sub-tasks consume high-value resources. As the number of agents and interaction turns grows, the quadratic increase in token consumption renders these systems difficult to scale.

To solve this problem, we propose AgentSlimming, a framework that can be directly applied to any MAS and workflows to reduce their execution costs. Motivated by the concepts of pruning, quantization, and finetuning in traditional neural network compression, AgentSlimming introduces AgentPruner, AgentQuant and AgentTuner as tools towards lightweight MAS. As shown in Figure 1, AgentPruner aims to identify the least important agent nodes in MAS and then remove them. In contrast, AgentQuant is designed as model substitution. When high-cost, high-precision models (e.g., GPT-4) are considered as not necessary for an unimportant role in MAS, they are dynamically replaced by cost-effective, lightweight surrogates (e.g., GPT-4o-mini) to reduce the inference costs. After AgentPruner and AgentQuant, AgentTuner is optionally employed to further optimize the MAS to recover the performance loss.

The proposition of AgentSlimming introduce a curcial problem for MAS compression, *i.e.*, how to determine the importance of each node in MAS. To tackle this challenge, we introduce a hybrid importance evaluation mechanism. Generally, a node’s value is determined by both its *structural position* and *functional contribution*. Based on this observation, we consider the three following indicators: Degree Centrality (Freeman, 1978) and Betweenness Centrality (Brandes, 2001) for topological structure in the graph of MAS, and an Approximate Shapley value (Shapley, 1953; Lundberg and Lee, 2017) for functional contribution. Specifically, we adopt a Leave-One-Out (LOO) estimation strategy to calculate the marginal contribution of each node against the complete topology, serving as a computationally efficient proxy for the exact Shapley value. These rankings are fused via Reciprocal Rank Fusion (RRF) (Cormack et al., 2009) to generate a robust importance score.

Guided by this score, AgentSlimming employs an iterative greedy strategy: we progressively prune and quantize agent nodes starting from the least important ones. After each operation, we re-evaluate the workflow against a performance threshold and re-calculate the RRF scores, ensuring the optimization dynamically adapts to the changing graph structure until the efficiency limit is reached. While the pruned and quantized workflows typically achieve a superior Pareto frontier (Deb et al., 2002; Chen et al., 2024) between cost and accuracy, AgentSlimming further incorporates a MCTS fine-tuning mechanism that optimizes the stream-

lined graph rather than accepting the compressed policy as final. Here, we report the average execution cost in USD per problem. Empirically, the resulting workflows are consistently more cost-efficient than AFlow, often achieving strictly superior performance in both accuracy and cost. Reporting the average execution cost per problem (USD), AgentSlimming demonstrates significant gains across diverse benchmarks.

For instance, on GSM8K, AgentSlimming matches AFlow’s score (95.5) while reducing the cost by **78.8%**. On code generation and complex reasoning tasks like MBPP and LiveCode, AgentSlimming achieves a "dual victory": it improves accuracy while reducing costs by **71.7%** and **78.9%**, respectively. Even in cases where AFlow retains a marginal score advantage (e.g., HotpotQA: 77.3 vs. 77.0), AgentSlimming achieves this at a substantially lower cost (a 27.4% reduction), indicating that our pipeline reliably locates the workflow at a more favorable operating point on the cost–quality Pareto frontier.

Our contributions can be summarized as follows:

- We introduce AgentSlimming, a training-free framework that integrates automated agentic workflow exploration with a novel pruning and semantic quantization pipeline.
- We first propose an iterative multi-metric optimization algorithm that integrates topological and functional importance via RRF to accurately identify redundant components.
- Extensive evaluations on eight benchmarks show that AgentSlimming achieves a striking reduction in token costs of up to **78.9%**.

2 Related Work

2.1 Agentic Workflows

LLM-based systems can be broadly characterized into two paradigms: *agentic workflows* and *autonomous agents*. The former executes tasks through predefined multi-step pipelines with repeated LLM invocations, whereas the latter emphasizes dynamic decision-making and planning under interaction and feedback. Recent work has made notable progress in language-driven decomposition and collaboration (Zhuge et al., 2025), data-science agents (Hong et al., 2025), tool-enabled mobile agent teams (Zhang et al., 2024), and open-ended embodied exploration (Wang et al., 2024a).

Compared to autonomous agents that often require environment-specific action spaces and decision patterns, workflows can more easily incorporate human domain expertise and improve through iterative refinement, making them particularly amenable to automated construction and optimization.

2.2 Multi-Agent Evolving

Many effective workflows are still developed primarily through manual discovery and engineering practice. At the general level, common transferable reasoning recipes include chain-of-thought, self-consistency, self-refine and structured self-collaboration (Wei et al., 2022; Wang et al., 2023; Madaan et al., 2023b; Wang et al., 2024b). At the domain level, multi-step procedures are specialized in reusable pipelines, such as code generation and debugging (Hong et al., 2024; Ridnik et al., 2024; Zhong et al., 2024), data analysis and visualization (Xie et al., 2024; Ye et al., 2024; Li et al., 2024a; Zhou et al., 2023), mathematical reasoning (Xu et al., 2024), and planning-style search for question answering (Nori et al., 2023; Zhou et al., 2024).

However, manual design cannot cover the combinatorial space across domains, which motivates automated agentic optimization. One line of work optimizes local instructions or components within a fixed backbone (Fernando et al., 2024; Yüksekönül et al., 2024; Yang et al., 2024; Khatatab et al., 2024) or tunes inference-time strategies (Saad-Falcon et al., 2024). Another line goes further by optimizing the end-to-end workflow structure, including automatic generation of code-represented workflows (Li et al., 2024b), viewing agent systems as optimizable graphs (Zhuge et al., 2024), and improving system designs in code space via meta-agents (Hu et al., 2025). Zhang et al. (2025c) similarly adopts code-based representations, combining finer-grained abstractions (named nodes/operators) with MCTS to leverage execution feedback and tree-structured experience for efficient workflow structure search.

2.3 Graph-Structured Orchestration, Pruning, and Cost-Aware Compression

Beyond prompt engineering, many agentic systems are best viewed as executable graphs, where nodes represent specialized modules and edges encode information flow. Recent studies focus on topology learning and pruning for improved efficiency and robustness (Zhang et al., 2025a,b; Boyi et al., 2025). These approaches also explore pruning at different

granularities, including message-level pruning under bandwidth constraints (Mao et al., 2020), dynamic agent elimination for token efficiency (Wang et al., 2025), and progressive pruning that blends heuristics with execution experience (Zhang et al., 2025d). Existing methods typically starting from (near) fully connected interactions and primarily pruning edges or communication channels. Action selection for pruning and replacement can draw on heuristic importance measures from network analysis (Freeman, 1978) or contribution-based formulations such as Shapley values (Shapley, 1953), which usually require approximation for tractability. In addition, model compression and quantization reduce inference cost while preserving quality (Frantar et al., 2022; Xiao et al., 2023; Lin et al., 2024), complementing structural optimization at the workflow level.

2.4 Search-Based Optimization for Workflows

Search provides a natural mechanism for exploring large workflow design spaces. MCTS (Coulom, 2006) and its variant UCT (Kocsis and Szepesvári, 2006) enables effective exploration through sampled evaluation and progressive expansion, and has been applied to planning and reasoning in language agents (Zhou et al., 2024) as well as workflow structure optimization (Zhang et al., 2025c). Following this trajectory, we focus on optimizing DAG workflows under dependency constraints, integrating structure search with node-level pruning and cost-driven node replacement (semantic quantization).

3 Methodology

3.1 Problem Formulation

We formulate the optimization of multi-agent systems as a discrete structural search problem over a directed graph, with an explicit trade-off between task performance and execution cost.

Workflow as a directed graph. A multi-agent workflow is represented as a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$. Each node $v \in \mathcal{V}$ denotes an executable operator (e.g., an LLM call), and each edge $e \in \mathcal{E}$ specifies an information dependency. Given a validation set \mathcal{D}_{val} , executing \mathcal{G} yields an average task score $S(\mathcal{G})$ and an average execution cost $C(\mathcal{G})$:

$$\begin{aligned} S(\mathcal{G}) &= \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{x \in \mathcal{D}_{\text{val}}} \text{score}(\mathcal{G}, x), \\ C(\mathcal{G}) &= \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{x \in \mathcal{D}_{\text{val}}} \text{cost}(\mathcal{G}, x). \end{aligned} \quad (1)$$

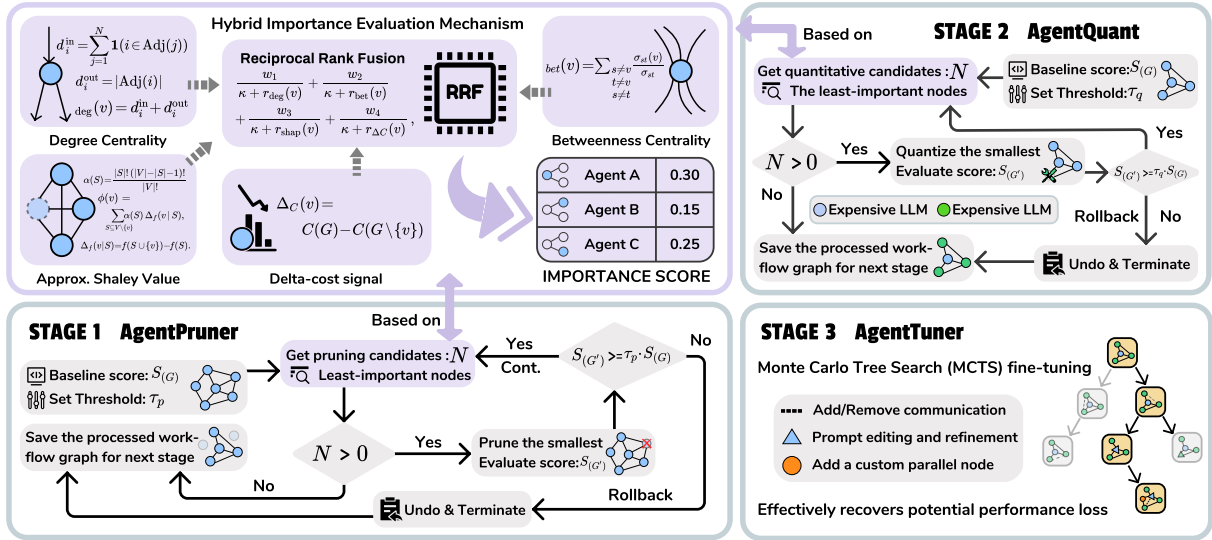


Figure 2: **Illustration of AgentSlimming.** To identify redundancy, AgentSlimming computes three distinct rankings of each agent node: *Degree Centrality*, *Betweenness Centrality*, *Costs Comparison* and *Approximate Shapley value*. Both pruning and quantization select candidate nodes based on the computed importance scores, ranking nodes from least to most important and optimizing the lowest-ranked candidate first. After each operation, a re-evaluation is performed. If the score drops below an acceptable threshold, a rollback mechanism is triggered to revoke the current operation and immediately terminate this phase, ensuring preserving superior performance.

Optimization Objective. Our goal is to maximize performance under a cost budget B . When multiple topologies achieve comparable performance, we prioritize the one with the minimal cost:

$$\max_{\mathcal{G}} S(\mathcal{G}) \quad \text{s.t.} \quad C(\mathcal{G}) \leq B. \quad (2)$$

3.2 The AgentSlimming Pipeline

As illustrated in Figure 2, our framework operates through a three-stage pipeline: *AgentPruner* (*Structural Pruning*), *AgentQuant* (*Semantic Quantization*), and *AgentTuner* (*MCTS Fine-tuning*). This pipeline progressively compresses the workflow while maintaining its reasoning capabilities.

3.2.1 Hybrid Importance Evaluation

To guide both pruning and quantization, we introduce a hybrid importance evaluation mechanism. For any node $v \in \mathcal{V}$, we compute four complementary signals using a small probe dataset $\mathcal{D}_{\text{probe}} \subset \mathcal{D}_{\text{val}}$.

1. Degree and Betweenness-Based Signals (Topological Priors). We capture the structural significance of a node using graph centrality metrics.

- **Degree-based Signal (s_{deg}):** Derived from the node’s local connectivity (in-degree and out-degree). Nodes with weaker structural connectivity are assigned larger pruning likelihood.

- **Betweenness Centrality (s_{bet}):** Measures the node’s role as an information bridge. It is calculated as:

$$s_{\text{bet}}(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}, \quad (3)$$

where σ_{st} is the total number of shortest paths from s to t , and $\sigma_{st}(v)$ is the number of those paths passing through v .

2. Approximate Shapley value signal (Functional Contribution). To quantify the semantic contribution of node v , we adopt the Shapley value concept from cooperative game theory. Since exact computation is NP-hard, we employ an efficient **Leave-One-Out (LOO)** approximation. The marginal contribution $\hat{\phi}(v)$ is estimated by the performance drop when v is removed (or replaced):

$$\hat{\phi}(v) \approx \Delta_S(v) = S(\mathcal{G}) - S(\mathcal{G} \setminus \{v\}), \quad (4)$$

where $\mathcal{G} \setminus \{v\}$ denotes the workflow after removing node v (and patching connections). A smaller $\hat{\phi}(v)$ indicates a lower contribution to task completion.

3. Delta-Cost signal (Economic Potential). We estimate the potential cost saving $\Delta_C(v)$ on $\mathcal{D}_{\text{probe}}$:

$$\Delta_C(v) = C(\mathcal{G}) - C(\mathcal{G} \setminus \{v\}). \quad (5)$$

Nodes with higher $\Delta_C(v)$ are prioritized for pruning or quantization to maximize efficiency gains.

Rank Fusion. To robustly combine these heterogeneous signals, we employ Reciprocal Rank Fusion (RRF). Let $r_m(v)$ be the rank of node v under metric $m \in \{\text{deg}, \text{bet}, \text{shap}, \Delta C\}$. The fused score is calculated as:

$$\text{RRF}(v) = \sum_m \frac{w_m}{\kappa + r_m(v)}, \quad (6)$$

where κ is a smoothing constant (set to $\max\{10, |\mathcal{V}|\}$) and \mathbf{w} balances the metrics. This fused score identifies nodes that are simultaneously structurally peripheral, functionally redundant, and computationally expensive.

3.2.2 Iterative Optimization Process

Stage 1: AgentPruner (Structural Pruning)

Let $\mathcal{G}_{\text{base}}$ denote the initial high-performance workflow graph. We employ an iterative greedy strategy to prune nodes. In each iteration, we rank nodes by $\text{RRF}(v)$ and evaluate the top candidates. A pruning operation is accepted if the performance degradation is within a tolerance threshold τ_p :

$$S(\mathcal{G}') \geq \tau_p \cdot S(\mathcal{G}_{\text{base}}), \quad (7)$$

where \mathcal{G}' is the pruned graph and $\tau_p \in (0, 1)$ is a predefined threshold. We evaluate the Top-1 candidate first and then the Top-k candidate if necessary. If neither candidate is accepted, the stage terminates and the workflow rolls back to the last accepted state. **Graph Surgery:** To preserve executability after removing node v , we perform edge patching to maintain the graph’s connectivity. For every pair of predecessor $s \in \text{In}(v)$ and successor $t \in \text{Out}(v)$, we add a direct edge ($s \rightarrow t$) if one does not already exist, self-loops and duplicate edges are disallowed. Each accepted pruning step is logged with the probe signals, fused ranks, and the resulting workflow artifact.

Stage 2: AgentQuant (Semantic Quantization)

Let $\mathcal{G}_{\text{base}}$ now refer to the pruned workflow graph obtained from Stage 1. The quantization process operates on this sparse structure to further reduce computational costs. We define *Semantic Quantization* as a model substitution strategy. For a selected node v , we replace its high-cost LLM with a cost-effective surrogate $\pi(v)$, yielding a quantized graph $\mathcal{G}^{(v)}$. Candidate ranking follows the same RRF mechanism, where the Shapley signal is instantiated by evaluating the performance impact of the model substitution. We also employ the greedy

strategy to quantize nodes, a quantization operation is accepted if the performance degradation is within a tolerance threshold τ_q :

$$S(\mathcal{G}') \geq \tau_q \cdot S(\mathcal{G}_{\text{base}}), \quad (8)$$

where \mathcal{G}' is the quantized graph and $\tau_q \in (0, 1)$ is a predefined threshold. This step drastically reduces token costs for non-critical reasoning steps.

Stage 3: AgentTuner (Adaptive MCTS Fine-tuning)

To address the variance in compression sensitivity across different tasks, we introduce an adaptive MCTS exploration on the compressed workflow. This phase is selectively applied and focuses on localized refinements, such as prompt adaptation and parameter tuning, to recover from potential degradation. Critically, this fine-tuned configuration is adopted only if it improves upon the initial compressed version, thereby ensuring constant and robust performance gains.

4 Experiments

4.1 Setup

Overview. We evaluate our method on eight public benchmarks. To ensure a fair comparison, we strictly adhere to AFlow’s configuration, adopting its datasets, splits, and sampling methods where applicable. Additionally, we incorporate three challenging benchmarks to test robustness.

Datasets. Our evaluation suite consists of two categories: (1) Standard Benchmarks: We utilize the complete AFlow (Zhang et al., 2025c) suite, including GSM8K (Cobbe et al., 2021), and MBPP (Austin et al., 2021) (full sets), as well as HotpotQA (Yang et al., 2018) and DROP (Dua et al., 2019) (randomly sampled 1,000-instance subsets). For MATH (Hendrycks et al., 2021), we follow the specific subsetting of 617 Level-5 problems across four categories. These datasets employ a 1:4 validation/test split. (2) High-Difficulty Benchmarks: To assess performance on complex reasoning and coding tasks, we incorporate AIME (Mathematical Association of America; Art of Problem Solving), MuSiQueAns (Trivedi et al., 2022), and LiveCode (Jain et al., 2025). These datasets are partitioned using a 3:7 validation/test split.

Baselines. We compare AgentSlimming against a diverse set of baselines: (1) *Manual Prompting Strategies*: Standard Chain-of-Thought (CoT) (Wei et al., 2022), Self-Consistency (SC-CoT)

Method	MATH		GSM8K		HotpotQA		MBPP		DROP	
	Acc.	Cost (\$/prob.)	Acc.	Cost (\$/prob.)	Acc.	Cost (\$/prob.)	Acc.	Cost (\$/prob.)	Acc.	Cost (\$/prob.)
Ours	73.9	6.88e-3 \downarrow 44.2%	95.5	9.30e-4 \downarrow 78.8%	77.0	5.55e-3 \downarrow 27.4%	77.9	7.30e-4 \downarrow 23.2%	81.7	9.50e-4 \downarrow 38.7%
AFlow	74.8	1.20e-2	95.5	4.38e-3	77.3	7.64e-3	73.3	2.58e-3	78.7	1.55e-3
LLM-Debate	74.5	1.47e-2	93.5	3.57e-3	73.4	9.47e-3	75.1	6.22e-3	75.6	2.94e-3
ADAS	69.0	1.15e-2	94.9	4.79e-3	65.4	4.72e-2	75.4	4.31e-3	78.6	3.07e-3
CoT	62.9	6.03e-4	90.5	2.48e-4	57.5	7.11e-4	73.6	2.80e-4	75.9	3.75e-4
SC-CoT	65.3	3.67e-3	90.2	1.71e-3	61.7	4.48e-3	75.6	2.12e-3	75.3	2.37e-3
Self-Refine	65.5	1.01e-3	89.8	5.60e-4	57.9	1.50e-3	74.5	5.56e-4	75.0	7.91e-4

Table 1: **Performance and inference cost comparison across standard benchmarks.** Cost denotes the average API expense (\$/prob.). **Bold** indicates the best results among **workflow-based methods** (excluding simple prompting baselines). Green percentages indicate the relative cost reduction compared to AFlow, highlighting how much inference budget is saved while maintaining competitive accuracy.

Method	AIME		LiveCode		MusiqueAns	
	Acc.	Cost (\$/prob.)	Acc.	Cost (\$/prob.)	Acc.	Cost (\$/prob.)
Ours	65.7	1.35e-2 \downarrow 19.2%	61.7	2.47e-3 \downarrow 78.9%	89.3	8.24e-3 \downarrow 23.2%
AFlow	67.1	1.67e-2	55.3	1.17e-2	84.5	1.07e-2
LLM-Debate	56.8	2.86e-2	67.8	1.67e-2	45.3	1.62e-2
Cost-aware AFlow	61.8	1.20e-2	57.9	8.41e-3	83.6	1.04e-2
ADAS	62.3	1.42e-2	52.1	9.01e-3	77.2	1.62e-2
CoT	53.4	2.99e-3	46.4	1.25e-3	73.9	1.11e-3
Self-Consistency CoT	59.7	1.41e-2	47.2	7.93e-3	74.9	6.89e-3
Self-Refine	54.1	4.08e-3	49.1	2.30e-3	72.8	2.30e-3

Table 2: **Performance and inference cost comparison across high-difficulty benchmarks.** Cost denotes the average API expense (\$/prob.). **Bold** indicates the best results among **workflow-based methods** (excluding simple prompting baselines). Green percentages indicate the relative cost reduction compared to AFlow.

(Wang et al., 2023), and Self-Refine (Madaan et al., 2023a). (II) *Automated Workflow Optimization*: ADAS (Hu et al., 2025) and AFlow (Zhang et al., 2025c). (III) *Cost-Aware AFlow*: We implement a modified AFlow variant. This baseline explicitly integrates cost comparisons into its selection mechanism during search phase, directly competing on cost efficiency.

Details. We employ GPT-4.1-mini (OpenAI, 2025) as the high-precision model and GPT-4.1-nano (OpenAI, 2025) as the cost-effective surrogate for quantization, both accessed via OpenAI API. For workflow-based optimization, we use GPT-4.1-mini to maintain computational efficiency. All models are accessed with temperature set to 0 to ensure reproducibility. Full experimental details are provided in Appendix B.

4.2 Main Results

Superior Cost-Performance Efficiency. The main experimental results are summarized in Table 1 and Table 2. Overall, AgentSlimming consistently discovers workflows that are more cost-

efficient than AFlow. In many cases, our method achieves *strict dominance*, surpassing the baseline in both accuracy and execution cost. Here, we report the average cost in USD per problem.

On standard benchmarks, AgentSlimming maintains competitive performance while significantly reducing computational overhead. For instance, on GSM8K, it matches AFlow’s accuracy (95.5) but reduces the average cost from 4.38×10^{-3} to 9.30×10^{-4} , a **78.8%** cost reduction. On MBPP, AgentSlimming improves the score from 73.3 to 77.9 while simultaneously lowering the cost from 2.58×10^{-3} to 7.30×10^{-4} (**71.7%** reduction). On datasets where AFlow attains a slightly higher raw score, such as HotpotQA, AgentSlimming achieves comparable performance at a substantially lower cost, reducing the average cost by 27.4%.

The advantage of AgentSlimming is even more pronounced on high-difficulty benchmarks. On LiveCode, it increases the score from 55.3 to 61.7 while reducing cost by nearly an order of magnitude (from 1.17×10^{-2} to 2.47×10^{-3} , a **78.9%** reduction). Similarly, on MuSiQueAns, AgentSlim-

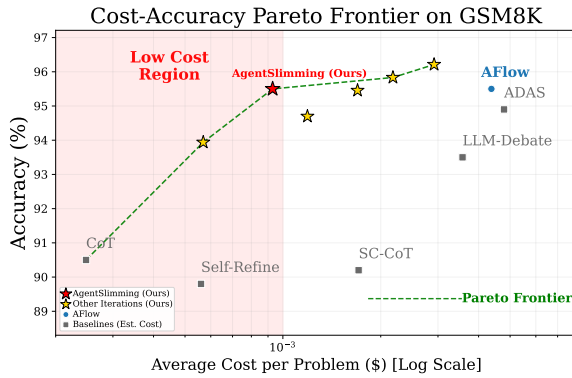


Figure 3: **Cost-Accuracy Pareto Frontier on the GSM8K dataset.** The x-axis represents the average inference cost per problem (USD), and the y-axis denotes accuracy. Data points for our method correspond to varying iteration rounds. These configurations consistently align with the Pareto frontier, demonstrating optimal cost-quality trade-offs compared to baselines.

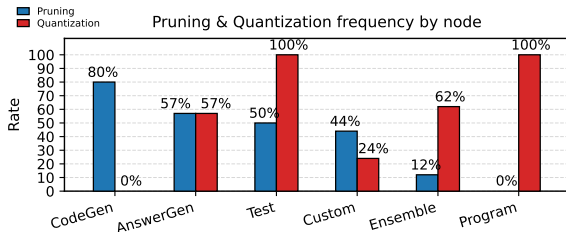


Figure 4: **Node-type specific compression rates.** We report the ratio of nodes processed by each compression strategy across six categories. Blue bars indicate structural removal (Pruning), while red bars indicate model substitution (Quantization). The variance across categories highlights the adaptive nature of our hybrid importance evaluation.

442 ming improves accuracy from 84.5 to 89.3 with
 443 a **23.2%** reduction. These results indicate that
 444 AgentSlimming reliably positions the discovered
 445 workflows at a more favorable operating points
 446 on the cost-quality Pareto frontier, particularly for
 447 complex reasoning and programming tasks.

448 4.3 Analysis

449 We conducted a comprehensive analysis to inter-
 450 pret the efficiency sources of AgentSlimming and
 451 validate its practical feasibility.

452 **Denosing via Pruning** In the initial heavy work-
 453 flows, we observed redundant parallel generation
 454 components, particularly multiple CodeGen nodes
 455 and AnswerGen nodes. Our method consistently
 456 deprioritizes these redundant nodes due to their
 457 limited marginal benefit, pruning them to yield a

simplified graph structure.

458 **Strategic Heterogeneity** When parallel struc-
 459 tures remain after pruning, AgentSlimming ap-
 460 plies non-uniform compression and retains a mixed-
 461 precision configuration. Specifically, it preserves a
 462 small number of high-capability AnswerGen nodes
 463 at full precision and quantizes the remaining par-
 464 allel auxiliary nodes to lower-cost models.
 465

466 Cost-Effective Arbitration via Functional Align- 467 ment

468 The quantization frequencies indicate that
 469 AgentSlimming primarily compresses nodes whose
 470 functionality is procedure-oriented rather than
 471 generation-intensive. ScEnsemble and Program
 472 nodes are infrequently pruned but are frequently
 473 quantized, while Test nodes are consistently
 474 quantized. Although these nodes can be topologically
 475 central, they mainly implement selection, consen-
 476 sus aggregation, and execution-based validation.
 477 By selectively quantizing them while retaining
 478 high-capability generators, AgentSlimming pre-
 479 serves accuracy and strategically directs the com-
 480 putational budget toward core reasoning rather than
 481 auxiliary validation.

482 Total Evaluation Budget Analysis.

483 To translate the structural efficiency gains into tangi-
 484 ble economic implications, we estimate the total end-
 485 to-end evaluation budget for the test sets of each
 486 benchmark by multiplying the optimized per-
 487 problem cost by the effective sample size. Re-
 488 markably, the computational burden is minimal:
 489 standard benchmarks like MBPP and GSM8K re-
 490 quire less than **\$2.00** to evaluate the entire test set
 491 (\$1.63 and \$1.84, respectively). Even for computa-
 492 tionally intensive, high-difficulty benchmarks such
 493 as AIME and HotpotQA, the total costs remain
 494 strictly affordable at \$26.06 and \$14.76.

495 Crucially, this cost-effectiveness extends to the
 496 optimization process itself. As a fully training-free
 497 framework, AgentSlimming avoids the prohibitive
 498 overhead of gradient-based updates and expensive
 499 parameter tuning. Furthermore, by leveraging topo-
 500 logical priors (e.g., degree and betweenness cen-
 501 trality) for importance scoring, we ensure that the
 502 pruning and quantization decisions remain compu-
 503 tationally efficient and lightweight.

504 4.4 Ablation Study

505 We investigate the impact of the node-importance
 506 ranking strategy used in pruning and quantization.
 Our full method employs a weighted Reciprocal

Method	DROP		MATH		MBPP	
	Acc.	Cost (\$/prob.)	Acc.	Cost (\$/prob.)	Acc.	Cost (\$/prob.)
Degree-only	66.5 ↓15.2%	3.2e-4 ↓28.9%	73.1 ↓2.3%	5.28e-3 ↑107.9%	73.3 ↓8.6%	1.56e-3 ↑119.7%
Betweenness-only	77.8 ↓0.8%	4.5e-4 ↓0.0%	74.0 ↓1.1%	4.17e-3 ↑64.2%	79.1 ↓1.4%	7.1e-4 ↓0.0%
Shapley-only	65.7 ↓16.2%	3.1e-4 ↓31.1%	71.4 ↓4.5%	5.25e-3 ↑106.7%	74.4 ↓7.2%	1.50e-3 ↑111.3%
RRF	78.4	4.5e-4	74.8	<u>2.54e-3</u>	80.2	<u>7.1e-4</u>

Table 3: **Ablation study of the importance ranking strategies.** We compare topology-based signals (degree-only, betweenness-only) and a functional signal (Shapley-only), as well as their fusion via Reciprocal Rank Fusion (RRF), on DROP, MATH, and MBPP. Results report task accuracy (Acc) and average inference cost per problem. **Bold** denotes the best accuracy on each dataset, while underline highlights the lowest average cost.

Rank Fusion (RRF) scheme to synthesize multiple signals into a unified candidate ranking. In contrast, we evaluate ablated variants that rely on a single metric for ranking: (i) Degree-only, (ii) Betweenness-only, and (iii) Shapley-only. Experimental settings are consistent across stages: for pruning, all variants initialize from the original workflow. For quantization, each variant inherits the best pruned graph derived from its corresponding pruning strategy and applies the same metric for node selection. Both stages maintain identical validation subset sizes and baseline-anchored acceptance thresholds. This design isolates the effect of ranking quality from confounders such as search budget, initialization, and stopping criteria. Table 3 details the specific performance metrics recorded after both stages.

We observe that single-metric strategies exhibit divergent cost-quality trade-offs across tasks. These complementary behaviors further underscore the necessity of multi-signal fusion via RRF, which consistently achieves a more robust balance between effectiveness and efficiency.

4.5 Case Study

We analyze the pruning process of the workflow on MATH to illustrate structural simplification. Starting from a complex graph with parallel reasoning paths and refinement steps, AgentSlimming identified redundancy in the dual Answer Generate nodes and the marginal utility of the Code Refine node. Concretely, probe-based importance ranking suggests that one generator path consistently dominates in contribution, while refinement rarely corrects errors relative to its token cost. Notably, pruning the parallel path diminished the utility of the downstream ScEnsemble node, leading to its removal. The transformation of this stage is detailed in Figure 5. The result is a streamlined pipeline

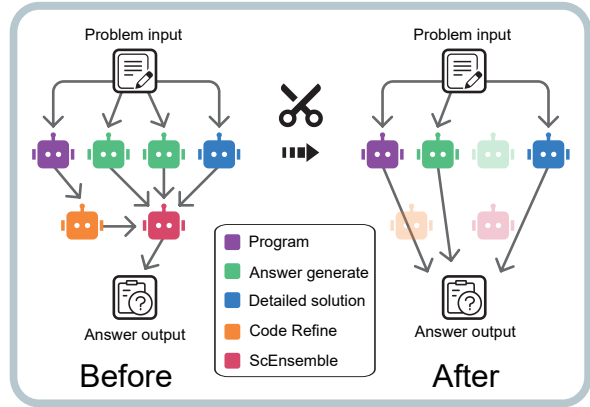


Figure 5: **Visualizing the pruning process of MATH workflow.** (Before) The initial workflow features complex parallel reasoning paths and an ensemble mechanism. (After) The streamlined workflow after AgentSlimming’s pruning stage. It effectively identifies the core reasoning backbone, pruning redundant parallel paths and refinement steps to reduce cost.

that preserves core reasoning logic while achieving substantial cost reductions, empirically validating our topo-functional optimization strategy. The complete details of this case study is provided in Appendix C.

5 Conclusion

We introduced AgentSlimming, a unified framework that optimizes multi-agent workflows via structural pruning and semantic quantization. Guided by a novel hybrid metric combining topological and functional signals, AgentSlimming achieves up to **79.8%** cost reduction across diverse benchmarks. Importantly, this efficiency is attained with negligible impact on reasoning quality, matching or even exceeding their performance. Our method shows that sparse, heterogeneous topologies can effectively replace computationally redundant dense multi-agent systems, establishing a new Pareto frontier for scalable agentic collaboration.

564 Limitations

565 We only evaluate GPT-4.1-series models, and do
566 not benchmark newer frontier models such as Ope-
567 nAI GPT-5.2, Anthropic Claude 4.5, Google Gem-
568 ini 3 Pro, Meta Llama 4.

569 Ethics Statement

570 We acknowledge that all authors are informed
571 about and adhere to the ACL ARR Code of Ethics
572 and the Code of Conduct.

573 **Risks** Our benchmarks are sourced from publicly
574 available datasets. We cannot guarantee that they
575 are free of socially harmful, biased, or toxic con-
576 tent. In addition, AgentSlimming optimizes cost by
577 pruning and replacing nodes, which may alter the
578 behavior of the original workflow in unexpected
579 ways; when applied to safety-critical domains, such
580 changes could amplify errors or reduce reliability.
581 We use ChatGPT to correct grammatical errors in
582 this paper.

583 References

584 Art of Problem Solving. Aime problems and solutions.
585 [https://artofproblemsolving.com/wiki/
586 index.php/AIME_Problems_and_Solutions](https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions).

587 Jacob Austin, Augustus Odena, Maxwell I. Nye,
588 Maarten Bosma, Henryk Michalewski, David Dohan,
589 Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le,
590 and Charles Sutton. 2021. [Program synthesis with
591 large language models](#). *CoRR*, abs/2108.07732.

592 Li Boyi, Zhonghan Zhao, Der-Hong Lee, and Gaoang
593 Wang. 2025. [Adaptive graph pruning for multi-agent
594 communication](#). *CoRR*, abs/2506.02951.

595 Ulrik Brandes. 2001. [A faster algorithm for between-
596 ness centrality](#). *Journal of Mathematical Sociology*.

597 Lingjiao Chen, Matei Zaharia, and James Zou. 2024.
598 [Frugalgpt: How to use large language models while
599 reducing cost and improving performance](#). *Trans.
600 Mach. Learn. Res.*

601 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
602 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
603 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
604 Nakano, Christopher Hesse, and John Schulman.
605 2021. [Training verifiers to solve math word prob-
606 lems](#). *CoRR*, abs/2110.14168.

607 Gordon V. Cormack, Charles L. A. Clarke, and Stefan
608 Büttcher. 2009. [Reciprocal rank fusion outperforms
609 condorcet and individual rank learning methods](#). In
610 *Proceedings of the 32nd Annual International ACM
611 SIGIR Conference on Research and Development
612 in Information Retrieval, SIGIR 2009, Boston, MA,
613 USA, July 19-23, 2009*.

Rémi Coulom. 2006. [Efficient selectivity and backup
operators in monte-carlo tree search](#). In *Computers
and Games, 5th International Conference, CG 2006,
Turin, Italy, May 29-31, 2006. Revised Papers*, vol-
ume 4630 of *Lecture Notes in Computer Science*,
pages 72–83. Springer. 614 615 616 617 618 619

Kalyanmoy Deb, Samir Agrawal, Amrit Pratap, and
T. Meyarivan. 2002. [A fast and elitist multiobjec-
tive genetic algorithm: NSGA-II](#). *IEEE Trans. Evol.
Comput.*, 6(2):182–197. 620 621 622 623

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel
Stanovsky, Sameer Singh, and Matt Gardner. 2019.
[DROP: A reading comprehension benchmark requir-
ing discrete reasoning over paragraphs](#). In *Proceed-
ings of the 2019 Conference of the North American
Chapter of the Association for Computational Lin-
guistics: Human Language Technologies, NAACL-
HLT 2019, Minneapolis, MN, USA, June 2-7, 2019,
Volume 1 (Long and Short Papers)*, pages 2368–2378.
Association for Computational Linguistics. 624 625 626 627 628 629 630 631 632 633

Chrisantha Fernando, Dylan Banarse, Henryk
Michalewski, Simon Osindero, and Tim Rock-
täschel. 2024. [Promptbreeder: Self-referential
self-improvement via prompt evolution](#). In *Forty-
first International Conference on Machine Learning,
ICML 2024, Vienna, Austria, July 21-27, 2024*.
OpenReview.net. 634 635 636 637 638 639 640

Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and
Dan Alistarh. 2022. [GPTQ: accurate post-training
quantization for generative pre-trained transformers](#).
CoRR, abs/2210.17323. 641 642 643 644

Linton C. Freeman. 1978. [Centrality in social networks
conceptual clarification](#). *Social Networks*, 1. 645 646

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul
Arora, Steven Basart, Eric Tang, Dawn Song, and Ja-
cob Steinhardt. 2021. [Measuring mathematical prob-
lem solving with the MATH dataset](#). In *Proceedings
of the Neural Information Processing Systems Track
on Datasets and Benchmarks 1, NeurIPS Datasets
and Benchmarks 2021, December 2021, virtual*. 647 648 649 650 651 652 653

Sirui Hong, Yizhang Lin, Bang Liu, and 1 others. 2025.
[Data interpreter: An LLM agent for data science](#).
In *Findings of the Association for Computational
Linguistics, ACL 2025, Vienna, Austria, July 27 -
August 1, 2025*, pages 19796–19821. Association for
Computational Linguistics. 654 655 656 657 658 659

Sirui Hong, Mingchen Zhuge, Jonathan Chen, and 1
others. 2024. [Metagpt: Meta programming for A
multi-agent collaborative framework](#). 660 661 662

Shengran Hu, Cong Lu, and Jeff Clune. 2025. [Au-
tomated design of agentic systems](#). In *The Thir-
teenth International Conference on Learning Repre-
sentations, ICLR 2025, Singapore, April 24-28, 2025*.
OpenReview.net. 663 664 665 666 667

668	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2025. Live-codebench: Holistic and contamination free evaluation of large language models for code . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	724
669		725
670		726
671		727
672		728
673		729
674		730
675		731
676	Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, and 1 others. 2024. Dspy: Compiling declarative language model calls into state-of-the-art pipelines . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	732
677		733
678		734
679		735
680		736
681		737
682	Levente Kocsis and Csaba Szepesvári. 2006. Bandit based monte-carlo planning . In <i>Machine Learning: ECML 2006, 17th European Conference on Machine Learning, Berlin, Germany, September 18-22, 2006, Proceedings</i> , volume 4212 of <i>Lecture Notes in Computer Science</i> , pages 282–293. Springer.	738
683		739
684		740
685		741
686		742
687		743
688	Boyan Li, Yuyu Luo, Chengliang Chai, Guoliang Li, and Nan Tang. 2024a. The dawn of natural language to SQL: are we fully ready? [experiment, analysis & benchmark] . <i>Proc. VLDB Endow.</i> , 17(11):3318–3331.	744
689		745
690		746
691		747
692		748
693	Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. CAMEL: Communicative agents for "mind" exploration of large language model society . In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	749
694		750
695		751
696		752
697		753
698	Zelong Li, Shuyuan Xu, Kai Mei, and 1 others. 2024b. Autoflow: Automated workflow generation for large language model agents . <i>CoRR</i> , abs/2407.12821.	754
699		755
700		756
701	Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Guangxuan Xiao, and Song Han. 2024. AWQ: activation-aware weight quantization for on-device LLM compression and acceleration . <i>GetMobile Mob. Comput. Commun.</i> , 28(4):12–17.	757
702		758
703		759
704		760
705		761
706	Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions . In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	762
707		763
708		764
709		765
710	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, and 1 others. 2023a. Self-refine: Iterative refinement with self-feedback . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	766
711		767
712		768
713		769
714		770
715		771
716		772
717	Aman Madaan, Niket Tandon, Prakhar Gupta, and 1 others. 2023b. Self-refine: Iterative refinement with self-feedback . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	773
718		774
719		775
720		776
721		777
722		778
723		
	Hangyu Mao, Zhengchao Zhang, Zhen Xiao, Zhibo Gong, and Yan Ni. 2020. Learning agent communication under limited bandwidth by message pruning . In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 5142–5149. AAAI Press.	
	Mathematical Association of America. Maa invitational competitions: American invitational mathematics examination (aime) . https://maa.org/maa-invitational-competitions/ .	
	Harsha Nori, Yin Tat Lee, Sheng Zhang, and 1 others. 2023. Can generalist foundation models outcompete special-purpose tuning? case study in medicine . <i>CoRR</i> , abs/2311.16452.	
	OpenAI. 2025. Introducing gpt-4.1 in the api: A new series of gpt models featuring major improvements on coding, instruction following, and long context—plus our first-ever nano model .	
	Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. 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> .	
	Tal Ridnik, Dedy Kredo, and Itamar Friedman. 2024. Code generation with alphacodium: From prompt engineering to flow engineering . <i>CoRR</i> , abs/2401.08500.	
	Jon Saad-Falcon and 1 others. 2024. Archon: An architecture search framework for inference-time techniques . <i>CoRR</i> , abs/2409.15254.	
	Lloyd S. Shapley. 1953. A value for n-person games . In <i>Contributions to the Theory of Games II</i> , pages 307–317. Princeton University Press.	
	Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. 9835 musique: Multi-hop questions via single-hop question composition . <i>Trans. Assoc. Comput. Linguistics</i> , 10:539–554.	
	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2024a. Voyager: An open-ended embodied agent with large language models . <i>Trans. Mach. Learn. Res.</i> , 2024.	
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	

Appendix

A Algorithm

Algorithm 1 Prune–Quantize–Fine-tune Pipeline

```
1: Input: initial graph  $G_0$ , validation set  $D_{\text{val}}$ ,  
   probe size  $m$ , budget  $B$ .  
2:   Top- $k$  size  $k$ , thresholds  $\tau, \tau_q$ .  
3:   RRF weights  $w$ , low-cost mapping  $\pi$ .  
4: Output: optimized graph  $G^*$ .  
5: Step 1: AgentPruner.  
6: Set baseline  $S_0 \leftarrow S(G)$ .  
7: while pruning budget remains do  
8:   Rank nodes by probe signals; fuse by  
   weighted RRF; get Top- $k$   $\{v_1, v_2, \dots, v_k\}$ .  
  
9:   for  $v \in \{v_1, v_2, \dots, v_k\}$  do  
10:     $G' \leftarrow \text{Surgery}(G, v)$ ; evaluate  $S(G')$ .  
11:    if  $S(G') \geq \tau \cdot S_0$  then  
12:       $G \leftarrow G'$ ; break  
13:    end if  
14:  end for  
15:  Stop if neither candidate is accepted.  
16: end while  
17: Step 2: AgentQuant.  
18: Set baseline  $S_0 \leftarrow S(G)$ .  
19: while quantization budget remains do  
20:   Rank LLM nodes by probe replacement  
   signals; fuse by weighted RRF; get Top- $k$   
    $\{v_1, v_2, \dots, v_k\}$ .  
21:   for  $v \in \{v_1, v_2, \dots, v_k\}$  do  
22:     $G^{(v)} \leftarrow \text{Replace}(G, v, \pi(v))$ ; evaluate  
     $S(G^{(v)})$ .  
23:    if  $S(G^{(v)}) \geq \tau_q \cdot S_0$  then  
24:       $G \leftarrow G^{(v)}$ ; break  
25:    end if  
26:  end for  
27:  Stop if neither candidate is accepted.  
28: end while  
29: Step 3: AgentTuner.  
30: Run MCTS on  $G$  with atomic edits; update  $G$ .  
31: return  $G^* \leftarrow \arg \max_{H \in \mathcal{K}} (S(H), -C(H))$ .  
   s.t.  $C(G^*) \leq B$ .
```

B Experimental Details

We employ different models for workflow optimization and execution. For execution, we primarily use a two-tier configuration with GPT-4.1-mini (OpenAI, 2025) and GPT-4.1-nano (OpenAI, 2025). Specifically, we employ GPT-4.1-mini by

default during the pruning stage, while adopting the lightweight GPT-4.1-nano for quantized nodes to minimize inference costs. To ensure a fair comparison, GPT-4.1-mini serves as the optimizer for both our method and all workflow-based baselines (AFlow (Zhang et al., 2025c), ADAS (Hu et al., 2025)). All models are accessed via APIs with a temperature of 0. We set the iteration rounds to 10 for AFlow and 15 for ADAS. Notably, as iteration rounds increase, workflows become more complex—a regime in which our method demonstrates superior cost reduction.

C Case Study

C.1 Workflow Overview and Legend Mapping

Figure 5 visualizes the pruning process on the MATH (Hendrycks et al., 2021) workflow. To match the legend used in the figure, we map the implementation node IDs to the paper terms as follows: (1) GenerateSolutionA/GenerateSolutionB correspond to the dual Answer Generate nodes. (2) RefineWithCode corresponds to the Code Refine node. (3) ScEnsembler corresponds to the ScEnsemble aggregation node. (4) AnswerFormatter formats the final output to satisfy benchmark constraints.

C.2 Pruning Rationale and Structural Simplification

The original workflow (Figure 5, left) contains multiple parallel reasoning paths: a program-execution path (Programmer \rightarrow RefineWithCode), a detailed chain-of-thought path (DetailedSolution), and two additional solution generation paths (GenerateSolutionA and GenerateSolutionB). All candidate answers are routed into ScEnsembler for self-consistency selection before being formatted by AnswerFormatter.

AgentSlimming identifies two main sources of redundancy: (1) The two Answer Generate branches are highly overlapping; keeping one branch provides most of the diversity gain. (2) The marginal utility of the Code Refine step is limited for this workflow; its benefits do not justify the extra LLM calls.

C.3 Complete Code for Reproducibility

We provide the complete executable code of the original MATH workflow and its pruned counterpart below.

```

1 from typing import Literal, Dict, Any
2 from src.core.graphflow import GraphFlow
3 from src.core.edge import Edge
4 from src.core.nodes.custom_node import CustomNode
5 from src.core.nodes.input_node import InputNode
6 from src.core.nodes.program_node import ProgramNode
7 from src.core.nodes.sc_ensemble_node import ScEnsembleNode
8 from src.core.executor import GraphExecutor
9 from src.utils.cost_tracker import LLMCostTracker
10 from src.core.nodes.answer_format_node import AnswerFormatNode
11 from . import prompt as prompt_custom
12
13 DatasetType = Literal["MATH"]
14
15 class Workflow:
16
17     def __init__(self, name: str, dataset: DatasetType, llm_config: Dict[str, Any] = None) → None:
18         self.name = name
19         self.dataset = dataset
20
21         if isinstance(llm_config, dict):
22             self.llm_config = llm_config
23         elif llm_config is None:
24             self.llm_config = {"model": "gpt-4.1-mini", "temperature": 0.0}
25         else:
26             self.llm_config = {
27                 "model": getattr(llm_config, "model", "gpt-4.1-mini"),
28                 "temperature": getattr(llm_config, "temperature", 0.0),
29             }
30
31         self.llm_cost_tracker = LLMCostTracker()
32
33         Input = InputNode(
34             node_id="Input",
35             node_llm_config=self.llm_config,
36             node_description="Graph input entry"
37         )
38
39         Programmer = ProgramNode(
40             node_id="Programmer",
41             node_llm_config=self.llm_config,
42             node_description="Generate and execute Python code for the problem"
43         )
44
45         RefineWithCode = CustomNode(
46             node_id="RefineWithCode",
47             node_prompt=prompt_custom.REFINE_ANSWER_PROMPT,
48             node_llm_config=self.llm_config,
49             node_description="Refine and format the solution using program output"
50         )
51
52         DetailedSolution = CustomNode(
53             node_id="DetailedSolution",
54             node_prompt=prompt_custom.DETAILED_SOLUTION_PROMPT,
55             node_llm_config=self.llm_config,
56             node_description="Provide a comprehensive, step-by-step solution"
57         )

```

Figure 6: Code of initial workflow on MATH

```

1      GenerateSolutionA = CustomNode(
2          node_id="GenerateSolutionA",
3          node_prompt=prompt_custom.GENERATE_SOLUTION_PROMPT,
4          node_llm_config=self.llm_config,
5          node_description="Generate an additional solution (A)"
6      )
7
8      GenerateSolutionB = CustomNode(
9          node_id="GenerateSolutionB",
10         node_prompt=prompt_custom.GENERATE_SOLUTION_PROMPT,
11         node_llm_config=self.llm_config,
12         node_description="Generate an additional solution (B)"
13     )
14
15     ScEnsembler = ScEnsembleNode(
16         node_id="ScEnsembler",
17         node_llm_config=self.llm_config,
18         node_description="Select the best solution using self-consistency",
19     )
20
21     AnswerFormatter = AnswerFormatNode(
22         node_id="AnswerFormatter",
23         dataset=self.dataset,
24         node_llm_config=self.llm_config,
25         node_description=f"Format the final answer for {self.dataset} benchmark compatibility.",
26     )
27
28     self.workflow_graph = GraphFlow(
29         nodes=[Input, Programmer, RefineWithCode, DetailedSolution,
30              GenerateSolutionA, GenerateSolutionB, ScEnsembler, AnswerFormatter],
31         edges=[
32             Edge(source="Input", target="Programmer"),
33             Edge(source="Input", target="DetailedSolution"),
34             Edge(source="Input", target="GenerateSolutionA"),
35             Edge(source="Input", target="GenerateSolutionB"),
36             Edge(source="Programmer", target="RefineWithCode"),
37             Edge(source="RefineWithCode", target="ScEnsembler"),
38             Edge(source="DetailedSolution", target="ScEnsembler"),
39             Edge(source="GenerateSolutionA", target="ScEnsembler"),
40             Edge(source="GenerateSolutionB", target="ScEnsembler"),
41             Edge(source="ScEnsembler", target="AnswerFormatter")
42         ],
43         entry_node_ids=["Input"],
44         final_node_id=AnswerFormatter.node_id,
45         description="Workflow for MATH problem solving"
46     )
47
48     self.executor = GraphExecutor(workflow_graph=self.workflow_graph)
49
50     async def __call__(self, problem: str):
51
52         result, execution_time, cost_summary = await self.executor.run(problem=problem)
53
54         final_output = (
55             result.get("output")
56             or str(result)
57         )
58
59         total_cost = cost_summary.get("total_cost_usd", 0.0)
60
61         return final_output, total_cost

```

Figure 7: Code of initial workflow on MATH (cont.)

```

1 from typing import Literal, Dict, Any
2 from src.core.graphflow import GraphFlow
3 from src.core.edge import Edge
4 from src.core.nodes.custom_node import CustomNode
5 from src.core.nodes.input_node import InputNode
6 from src.core.nodes.program_node import ProgramNode
7 from src.core.nodes.sc_ensemble_node import ScEnsembleNode
8 from src.core.executor import GraphExecutor
9 from src.utils.cost_tracker import LLMCostTracker
10 from src.core.nodes.answer_format_node import AnswerFormatNode
11 from . import prompt as prompt_custom
12
13 DatasetType = Literal["MATH"]
14
15 class Workflow:
16
17     def __init__(self, name: str, dataset: DatasetType, llm_config: Dict[str, Any] = None) → None:
18         self.name = name
19         self.dataset = dataset
20
21         if isinstance(llm_config, dict):
22             self.llm_config = llm_config
23         elif llm_config is None:
24             self.llm_config = {"model": "gpt-4.1-mini", "temperature": 0.0}
25         else:
26             self.llm_config = {
27                 "model": getattr(llm_config, "model", "gpt-4.1-mini"),
28                 "temperature": getattr(llm_config, "temperature", 0.0),
29             }
30
31         self.llm_cost_tracker = LLMCostTracker()
32
33         Input = InputNode(
34             node_id="Input",
35             node_llm_config=self.llm_config,
36             node_description="Graph input entry"
37         )
38
39         Programmer = ProgramNode(
40             node_id="Programmer",
41             node_llm_config=self.llm_config,
42             node_description="Generate and execute Python code for the problem"
43         )
44
45         DetailedSolution = CustomNode(
46             node_id="DetailedSolution",
47             node_prompt=prompt_custom.DETAILED_SOLUTION_PROMPT,
48             node_llm_config=self.llm_config,
49             node_description="Provide a comprehensive, step-by-step solution"
50         )

```

Figure 8: Code of pruned workflow on MATH

```

1     GenerateSolutionA = CustomNode(
2         node_id="GenerateSolutionA",
3         node_prompt=prompt_custom.GENERATE_SOLUTION_PROMPT,
4         node_llm_config=self.llm_config,
5         node_description="Generate an additional solution "
6     )
7
8     AnswerFormatter = AnswerFormatNode(
9         node_id="AnswerFormatter",
10        dataset=self.dataset,
11        node_llm_config=self.llm_config,
12        node_description=f"Format the final answer for {self.dataset} benchmark compatibility.",
13    )
14    # Graph
15    self.workflow_graph = GraphFlow(
16        nodes=[Input, Programmer, RefineWithCode, DetailedSolution,
17              GenerateSolutionA, GenerateSolutionB, ScEnsembler, AnswerFormatter],
18        edges=[
19            Edge(source="Input", target="Programmer"),
20            Edge(source="Input", target="DetailedSolution"),
21            Edge(source="Input", target="GenerateSolutionA"),
22            Edge(source="Programmer", target="AnswerFormatter"),
23            Edge(source="DetailedSolution", target="AnswerFormatter"),
24            Edge(source="GenerateSolutionA", target="AnswerFormatter"),
25        ],
26        entry_node_ids=["Input"],
27        final_node_id=AnswerFormatter.node_id,
28        description="Workflow for MATH problem solving"
29    )
30
31    self.executor = GraphExecutor(workflow_graph=self.workflow_graph)
32
33    async def __call__(self, problem: str):
34
35        result, execution_time, cost_summary = await self.executor.run(problem=problem)
36
37        final_output = (
38            result.get("output")
39            or str(result)
40        )
41
42        total_cost = cost_summary.get("total_cost_usd", 0.0)
43
44        return final_output, total_cost

```

Figure 9: Code of pruned workflow on MATH (cont.)

```

1 from typing import List, Dict
2 from collections import defaultdict
3 from src.core.pruners.base import Pruner
4 from src.core.graphflow import GraphFlow
5 from typing import Tuple
6 from src.utils.logs import logger
7
8 class DegreePruner(Pruner):
9
10     def get_pruned(self) → Tuple[Dict[str, GraphFlow], GraphFlow]:
11
12         pruning_candidates = self.identify_pruning_candidates()
13         one_pruned_graphs: Dict[str, GraphFlow] = {}
14
15         for node_to_remove in pruning_candidates.keys():
16             new_nodes = [
17                 node for node in self.original_graph.nodes
18                 if node.node_id ≠ node_to_remove
19             ]
20             new_edges = [
21                 edge for edge in self.original_graph.edges
22                 if edge.source ≠ node_to_remove and edge.target ≠ node_to_remove
23             ]
24             new_graph = GraphFlow(
25                 nodes=new_nodes,
26                 edges=new_edges,
27                 entry_node_ids=self.original_graph.entry_node_ids,
28                 final_node_id=self.original_graph.final_node_id,
29                 description=f"Pruned from {self.original_graph.description} (removed: {node_to_remove})"
30             )
31
32             one_pruned_graphs[str(node_to_remove)] = new_graph
33
34         new_nodes = [
35             node for node in self.original_graph.nodes
36             if node.node_id not in pruning_candidates.keys()
37         ]
38         new_edges = [
39             edge for edge in self.original_graph.edges
40             if edge.source not in pruning_candidates.keys() and edge.target not in pruning_candidates.keys()
41         ]
42
43         all_pruned_graph = GraphFlow(
44             nodes=new_nodes,
45             edges=new_edges,
46             entry_node_ids=self.original_graph.entry_node_ids,
47             final_node_id=self.original_graph.final_node_id,
48             description=f"Pruned from {self.original_graph.description} (Pruned all candidates)"
49         )
50
51         return one_pruned_graphs, all_pruned_graph

```

Figure 10: Code of Degree Centrality Algorithm

```

1  def identify_pruning_candidates(self) → Dict[str, float]:
2
3      pruning_candidates = {}
4      in_degree = defaultdict(int)
5      out_degree = defaultdict(int)
6
7      for node in self.original_graph.nodes:
8          in_degree[node.node_id] = 0
9          out_degree[node.node_id] = 0
10
11     for edge in self.original_graph.edges:
12         out_degree[edge.source] += 1
13         in_degree[edge.target] += 1
14
15     for node in self.original_graph.nodes:
16         node_id = node.node_id
17         in_deg = in_degree[node_id]
18         out_deg = out_degree[node_id]
19
20         if node_id in self.original_graph.entry_node_ids:
21             continue
22
23         if node_id == self.original_graph.final_node_id:
24             continue
25
26         if is_low_importance_node(in_deg, out_deg):
27             pruning_candidates[node.node_id] = 1.0
28
29     return pruning_candidates
30
31
32 def is_low_importance_node(in_degree: int, out_degree: int) → bool:
33
34     if in_degree + out_degree ≤ 2:
35         return True
36
37     return False

```

Figure 11: Code of Degree Centrality Algorithm (cont.)

```

1 from typing import List, Dict, Tuple
2 from collections import defaultdict, deque
3 from src.core.pruners.base import Pruner
4 from src.core.graphflow import GraphFlow
5 from src.utils.logs import logger
6
7
8 class BetweennessPruner(Pruner):
9
10     def __init__(self, original_graph: GraphFlow):
11
12         super().__init__(original_graph)
13         self.betweenness_scores: Dict[str, float] = {}
14
15     def compute_betweenness centrality(self) → Dict[str, float]:
16
17         betweenness = defaultdict(float)
18
19         for node in self.original_graph.nodes:
20             if node.node_id not in self.original_graph.entry_node_ids and node.node_id ≠
self.original_graph.final_node_id:
21                 betweenness[node.node_id] = 0.0
22
23         graph_adj = defaultdict(list)
24         for edge in self.original_graph.edges:
25             graph_adj[edge.source].append(edge.target)
26
27         all_nodes = [node.node_id for node in self.original_graph.nodes]
28
29         for source in all_nodes:
30
31             stack = []
32             paths = defaultdict(list)
33             sigma = defaultdict(int)
34             dist = defaultdict(lambda: -1)
35             delta = defaultdict(float)
36
37             sigma[source] = 1
38             dist[source] = 0
39             queue = deque([source])
40
41             while queue:
42                 v = queue.popleft()
43                 stack.append(v)
44
45                 for w in graph_adj[v]:
46                     if dist[w] < 0:
47                         queue.append(w)
48                         dist[w] = dist[v] + 1
49                     if dist[w] == dist[v] + 1:
50                         sigma[w] += sigma[v]
51                         paths[w].append(v)
52
53             while stack:
54                 w = stack.pop()
55                 for v in paths[w]:
56                     delta[v] += (sigma[v] / sigma[w]) * (1 + delta[w])
57
58                 if w ≠ source and w not in self.original_graph.entry_node_ids and w ≠
self.original_graph.final_node_id:
59                     betweenness[w] += delta[w]
60
61         n = len(all_nodes)
62         if n > 2:
63             normalization = (n - 1) * (n - 2)
64             for node_id in betweenness:
65                 betweenness[node_id] /= normalization
66
67         return dict(betweenness)

```

Figure 12: Code of Betweenness Centrality Algorithm

```

1  def identify_pruning_candidates(self) → Dict[str, float]:
2
3      self.betweenness_scores = self.compute_betweenness_centrality()
4
5      pruning_candidates = {}
6      for node in self.original_graph.nodes:
7          node_id = node.node_id
8
9          if node_id in self.original_graph.entry_node_ids:
10             continue
11
12             if node_id == self.original_graph.final_node_id:
13                 continue
14
15                 pruning_candidates[node_id] = self.betweenness_scores.get(node_id, 0)
16
17         return pruning_candidates
18
19     def get_pruned(self) → Tuple[Dict[str, GraphFlow], GraphFlow]:
20
21         pruning_candidates = self.identify_pruning_candidates()
22
23         one_pruned_graphs: Dict[str, GraphFlow] = {}
24
25         for node_to_remove in pruning_candidates.keys():
26
27             new_nodes = [
28                 node for node in self.original_graph.nodes
29                 if node.node_id ≠ node_to_remove
30             ]
31             new_edges = [
32                 edge for edge in self.original_graph.edges
33                 if edge.source ≠ node_to_remove and edge.target ≠ node_to_remove
34             ]
35
36             new_graph = GraphFlow(
37                 nodes=new_nodes,
38                 edges=new_edges,
39                 entry_node_ids=self.original_graph.entry_node_ids,
40                 final_node_id=self.original_graph.final_node_id,
41                 description=f"Pruned from {self.original_graph.description} (removed: {node_to_remove},
betweenness: {self.betweenness_scores.get(node_to_remove, 0):.4f})"
42             )
43
44             one_pruned_graphs[str(node_to_remove)] = new_graph
45
46             new_nodes = [
47                 node for node in self.original_graph.nodes
48                 if node.node_id not in pruning_candidates.keys()
49             ]
50             new_edges = [
51                 edge for edge in self.original_graph.edges
52                 if edge.source not in pruning_candidates.keys() and edge.target not in pruning_candidates.keys()
53             ]
54
55             all_pruned_graph = GraphFlow(
56                 nodes=new_nodes,
57                 edges=new_edges,
58                 entry_node_ids=self.original_graph.entry_node_ids,
59                 final_node_id=self.original_graph.final_node_id,
60                 description=f"Pruned all candidates"
61             )
62
63         return one_pruned_graphs, all_pruned_graph

```

Figure 13: Code of Betweenness Centrality Algorithm (cont.)

```

1 from typing import Tuple, Dict, List
2 from copy import copy
3 from src.core.graphflow import GraphFlow
4 from src.utils.logs import logger
5 from src.core.pruners.base import Pruner
6 from src.core.graphflow import Edge
7
8 class ShapleyPruner(Pruner):
9
10     def __init__(self, original_graph: GraphFlow):
11         super().__init__(original_graph)
12
13     def identify_pruning_candidates(self) → Dict[str, float]:
14         pruning_candidates = {}
15         for node in self.original_graph.nodes:
16             if node.node_id ≠ self.original_graph.final_node_id and node.node_id not in
self.original_graph.entry_node_ids:
17                 pruning_candidates[node.node_id] = 1.0
18
19         return pruning_candidates
20
21     def get_pruned(self) → Tuple[Dict[str, GraphFlow], GraphFlow]:
22
23         pruning_candidates = self.identify_pruning_candidates()
24         pruned_graphs: Dict[str, GraphFlow] = {}
25
26         for node_id in pruning_candidates.keys():
27             pruned_graph = self._prune_node(self.original_graph, node_id)
28             if pruned_graph:
29                 pruned_graphs[node_id] = pruned_graph
30
31         all_pruned_graph = None
32
33         return pruned_graphs, all_pruned_graph
34
35     def prune_node(self, graph: GraphFlow, node_id: str):
36
37         new_nodes = [copy(n) for n in graph.nodes if n.node_id ≠ node_id]
38         new_edges = [copy(e) for e in graph.edges if e.source ≠ node_id and e.target ≠ node_id]
39
40         predecessors = [e.source for e in graph.edges if e.target = node_id]
41         successors = [e.target for e in graph.edges if e.source = node_id]
42         existing_edges = {(e.source, e.target) for e in new_edges}
43
44         for p in predecessors:
45             for s in successors:
46                 if (p, s) not in existing_edges and p ≠ s:
47                     old_edge = next((e for e in graph.edges if e.source = p and e.target = node_id), None)
48                     if old_edge:
49                         new_edge = copy(old_edge)
50                         new_edge.source = p
51                         new_edge.target = s
52                     else:
53                         new_edge = Edge(source=p, target=s)
54                     new_edges.append(new_edge)
55                     existing_edges.add((p, s))
56
57         try:
58             new_graph = GraphFlow(
59                 nodes=new_nodes,
60                 edges=new_edges,
61                 entry_node_ids=[id for id in graph.entry_node_ids if id ≠ node_id],
62                 final_node_id=graph.final_node_id,
63                 description=f"Graph which has pruned node:{node_id} ",
64             )
65             return new_graph
66         except Exception as e:
67             logger.error(f"Fail to prune node {node_id} failed: {e}")
68             return None

```

Figure 14: Code of Approximate Shapley value Algorithm