# AgentArk: Distilling Multi-Agent Intelligence into a Single LLM Agent

Yinyi Luo[1], Yiqiao Jin[3], Weichen Yu[1], Mengqi Zhang[2], Srijan Kumar[3], Xiaoxiao Li[5],
Weijie Xu[4*], Xin Chen[4], Jindong Wang[2*]

[1]Carnegie Mellon University   [2]William & Mary   [3]Georgia Institute of Technology
[4]Amazon   [5]University of British Columbia

### Abstract

While large language model (LLM) multi-agent systems achieve superior reasoning performance through iterative debate, practical deployment is limited by their high computational cost and error propagation. This paper proposes **AgentArk**, a novel framework to distill multi-agent dynamics into the weights of a *single* model, effectively transforming explicit test-time interactions into implicit model capabilities. This equips a single agent with the intelligence of multi-agent systems while remaining computationally efficient. Specifically, we investigate three hierarchical distillation strategies across various models, tasks, scaling, and scenarios: reasoning-enhanced fine-tuning; trajectory-based augmentation; and process-aware distillation. By shifting the burden of computation from inference to training, the distilled models preserve the efficiency of one agent while exhibiting strong reasoning and self-correction performance of multiple agents. They further demonstrate enhanced robustness and generalization across diverse reasoning tasks. We hope this work can shed light on future research on efficient and robust multi-agent development. Our code is at https://github.com/AIFrontierLab/AgentArk.

## 1 Introduction

Multi-agent Systems (MAS), where multiple models interact through debate (Du et al., 2023; Eo et al., 2025; Hegazy, 2024), critique (Lan et al., 2024; Yu et al., 2025), and consensus (Chen et al., 2024a), have demonstrated remarkable success in complex reasoning tasks (Guo et al., 2024). By structuring reasoning as a multi-turn and multi-role dialogue, MAS can explore diverse hypotheses, uncover logical errors (Wang et al., 2024b), and iteratively refine solutions (Wan et al., 2025). However, this collaborative power is a double-edged sword that introduces systemic risks including *computational overhead* and *vulnerability amplification*. First, the reliance on multi-role and multi-turn dialogue causes inference latency and computational overhead to grow rapidly (Kim et al., 2025a; Wang et al., 2025e). In densely connected networks, computation can grow quadratically with the number of agents, making MAS prohibitively expensive for real-time (Ignise & Vahi, 2024; Kim et al., 2025b). Second, while MAS can correct errors, they can also escalate them. In high-density interactions, individual biases or hallucinations can propagate and amplify across the group, leading to collective failures in robustness and safety (He et al., 2025; Nguyen et al., 2025).

These challenges naturally raise a fundamental question: *Can a single model internalize the reasoning benefits of MAS without their high inference-time cost and collaborative vulnerabilities?*

The collective power of MAS, i.e., *inference-time compute*, suggests that the gains can possibly be "shifted forward"—internalized by a single model during offline learning (Chen et al., 2025; Liu et al., 2025c; Wang et al., 2025b). While prior work has shown that single models can absorb certain MAS reasoning benefits (Chen et al., 2024b; Li et al., 2025; Li, 2026; Zhou et al., 2025), they are often limited to imitating final answers (Han et al., 2024) or exhibiting shallow interaction traces (Li et al., 2023; Lin et al., 2025), failing to reproduce the core *iterative conflict-and-refinement dynamics* that underlie MAS reasoning (Li et al., 2024). Notably, recent evidence suggests that the structural design of multi-agent systems may play a secondary role in their observed gains. Kim et al. (2025b) demonstrate that removing or perturbing explicit agent structures leads to only marginal performance degradation, revealing that the essential contribution of MAS lies in the *reasoning dynamics* they induce, rather than in the interaction schema itself. Similarly, Ke et al. (2026) systematically studies orchestration strategies under controlled benchmarks,

---

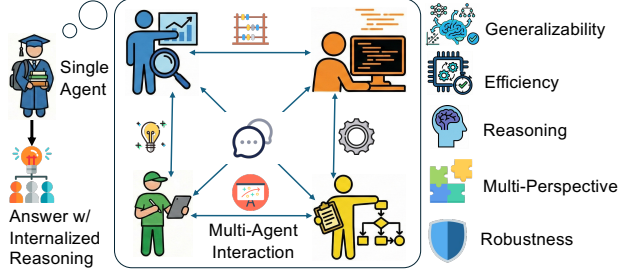arXiv:2602.03955v1 [cs.AI] 3 Feb 2026

Figure 1: AgentArk distills the reasoning capability of multi-agent systems into one single agent, such that this single unit can imitate the thinking process with boosted performance.

showing that performance is largely driven by induced reasoning behaviors rather than specific agent topologies or coordination protocols.

Our pilot study has shown that fine-tuning only on final answers can lead to overfitting on task-specific mapping with limited generalization (Appendix G.1). Therefore, a single agent should learn to internalize the reasoning dynamics, allowing for generating, evaluating, and refinement with one forward pass. This perspective also resonates with human cognition: individuals can internalize group reasoning strategies and reproduce collective wisdom independently (Curşeu et al., 2015; Navajas et al., 2018; Toyokawa et al., 2019). In this paper, we present a comprehensive investigation of multi-agent distillation for reasoning tasks. Particularly, we propose **AgentArk** (Figure 1 and 2), a general and scalable distillation paradigm that transfers MAS reasoning dynamics into a single model, without relying on handcrafted interaction or task-specific supervision.[1]

AgentArk integrates three progressively deeper distillation strategies: 1) Reasoning-Enhanced Outcome–Based Supervision *(R-SFT)*: Training the model on the final consensus reached by the agent group while leveraging reasoning trajectories to ensure the single agent can consistently reach high-quality conclusions; 2) Reasoning Trajectory-based Data Augmentation *(DA)*: Distilling the diverse reasoning chains derived from collective interaction, allowing the model to learn a variety of logical strategies and ways of thinking; and 3) Process-Aware Distillation *(PAD)*: Using process reward model to train the agent to internalize the critique-and-revision dynamics, enabling a single agent to emulate the dialectical reasoning of multi-agent debates within a single forward pass. We conduct extensive experiments by varying LLM backbones, teacher-student roles, datasets and various tasks, scaling, and evaluation settings. The following are our key findings:

1. **AgentArk enables a single agent to acquire multi-agent reasoning ability.** Our extensive experiments show that all three reasoning-centric distillation methods can boost the performance of single agents. Combination of approaches can achieve further improvement. (§4.2)
2. **PRM capacity matters more than student model size, while student capacity bounds multi-agent gains.** High-capacity PRMs enable strong improvements even for small students, whereas weak PRMs limit gains. Scaling teacher agents mainly benefits larger students, with diminishing returns for smaller ones. (§4.2, §4.3)
3. **Reasoning quality outweighs quantity.** Simply adding more reasoning trajectories does not improve performance, while PAD's high-signal process supervision yields stable gains. (§4.3)
4. **Process-aware distillation improves reasoning behavior, not just accuracy.** PAD models exhibit better step decomposition, self-checking, and error correction than RSFT and RA. (§4.4)
5. **AgentArk improves generalization and robustness.** Distilled models transfer reliably to unseen and robustness benchmarks. (§4.5)
6. **AgentArk can extend to multimodal LLMs.** (§4.6)

**Contributions.** (1) To our best knowledge, AgentArk is the *first* comprehensive framework to explore various strategies for MAS distillation. (2) We construct a scalable distillation data generation pipeline and framework agnostic to MAS strategies, which will be released to facilitate future research on reasoning distillation. (3) We perform extensive evaluation of MAS distillation from various perspectives, providing insights for future research.

## 2  Related Work

**Multi-Agent Systems.** MASs have emerged as an effective paradigm for enhancing LLM reasoning by enabling multiple agents to interact (Chen et al., 2024a; Du et al., 2023; Eo et al., 2025; Wei et al., 2026; Yuan & Xie, 2025).

---

[1]This paper only focuses on reasoning tasks; tool use (Qin et al., 2024) and memory management (Zhang et al., 2025) are left for future work. We consider the popular homogeneous setting where all agents share the same LLM backbone (Chen et al., 2023; Eo et al., 2025).
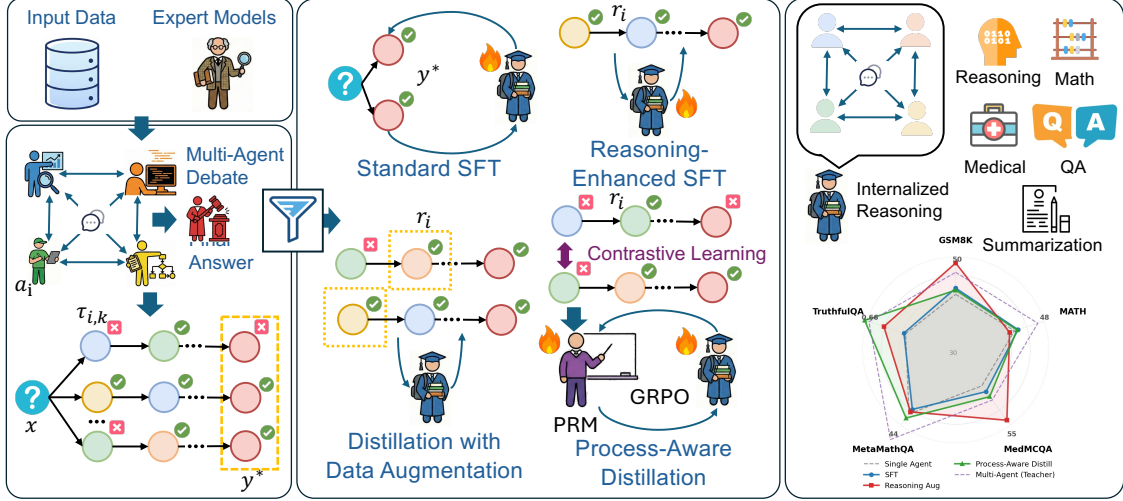
Figure 2: **Overview of AgentArk.** The pipeline proceeds through three stages: (1) **Data Generation Through Multi-Agent Debate** to produce diverse reasoning trajectories; (2) **Knowledge Extraction** to filters for high-quality corrective traces; and (3) **Distillation** utilizing Standard SFT, Reasoning-enhanced SFT, Distillation with Data Augmentation, and Process-Aware Distillation (PRM + GRPO). The resulting student model achieves optimized, low-latency reasoning that generalizes across diverse task domains.

By organizing reasoning as explicit multi-turn interactions, they can explore diverse solution paths, detect errors, and iteratively refine predictions, leading to strong performance on complex reasoning tasks (Wang et al., 2024b, 2025f, 2026). However, MAS rely on inference-time coordination among multiple agents, which incurs substantial computational cost and latency (Ma et al., 2026; Shi et al., 2026; Wang et al., 2025c,d). In addition, agent roles, interaction protocols, and evaluation criteria are typically designed for specific tasks, limiting the applicability in resource-constrained or real-time settings. Beyond efficiency concerns, recent work (Kim et al., 2025b) has examined the structural sensitivity of MAS, showing that performance is often robust to substantial simplifications or perturbations of agent structures.

**Distillation of Multi-Agent Reasoning.** To avoid the above issues, recent work has explored distilling multi-agent reasoning into a single model (Liu et al., 2025a; Pan et al., 2025; Wang et al., 2025a; Zhao et al., 2024). Early approaches supervise student models using the final outputs of agent groups or simplified interaction traces (Kang et al., 2025; Liu et al., 2025b). More recent methods transfer richer supervision signals, including graph-based interaction modeling (Chen et al., 2024b), skill selection frameworks (Li, 2026), debate-derived preference supervision (Zhou et al., 2025) and end-to-end agentic reinforcement learning (Li et al., 2025). Despite their effectiveness, they often depend on task-specific agent designs, predefined interaction structures, or environment-dependent reward functions, restricting generalization across tasks and limit scalability. In contrast, AgentArk abstracts away agent roles and interaction structures by distilling interaction-induced reasoning processes at the process level, which enables a single student model to generalize across tasks and agent configurations without task-specific agent design.

# 3 Method

## 3.1 Overview

As shown in Figure 2, AgentArk consists of three phases: (1) **Data Generation via Multi-Agent Interaction**, where an ensemble of teacher models generates diverse reasoning trajectories. Here, we leverage the iterative reflection and error-correction trajectories inherent in LLM debates (Du et al., 2023). (2) **Knowledge Extraction**, where successful and corrective traces are generated and filtered; and (3) **Hierarchical Distillation**, where the student is trained via supervised learning and process-level reinforcement learning. We explore three distillation paradigms: (1) Reasoning-enhanced Supervised Fine-Tuning (RSFT), (2) Data Augmentation via Diverse Extraction (DA), and (3) Process-Aware Distillation using PRM and GRPO. Note that while debate is adopted, AgentArk is agnostic to MAS algorithms.

## 3.2 Data Generation and Knowledge Extraction

We adopt debate (Liang et al., 2024) as the MAS mechanism for generating data that captures rich reasoning dynamics. Interactions in which individuals reflect, revise, and converge have long been recognized as hallmarks of intelligent behavior, and debate naturally surfaces processes such as self-correction, error discovery, and cross-examination. Prior work on LLM debate Du et al. (2023) demonstrates that structured disagreement can enhance both factual accuracy and reasoning quality.

**Debate Setup.** For each input $x$, we initialize $n$ debating agents $\mathcal{A} = \{a_1, a_2, \ldots, a_n\}$ sharing the same LLM. The agents engage in a $K$-round interaction. In each round, each agent $a_i$ generates a reasoning trace $\tau_{i,k}$ based on the problem $x$ and the previous traces of its peers $\{\tau_{j,k-1}\}_{j \neq i}$. This setup encourages agents to identify logical inconsistencies in others' arguments and iteratively refine their own reasoning traces. This process continues until $K$ rounds is reached or a consensus emerges. The result is a comprehensive debate log $\mathcal{L}_x$ containing a set of final reasoning traces $\{r_1, r_2, \ldots, r_n\}$ and their corresponding answers.

**Correctness-First Trajectory Selection.** While traditional distillation typically utilizes error-free paths, we prioritize **corrective trajectories**, reasoning traces where an agent initially proposed an incorrect step but recognize successfully pivoted to a correct answer $y^*$ after receiving critiques. For each task, we extract 1) the final consensus answer $y^*$, verified against ground-truth labels; 2) the intermediate reasoning traces $\{r_i\}$ that successfully lead to $y^*$.

**Knowledge Extraction.** Our pipeline extracts multiple answer-consistent yet structurally diverse reasoning trajectories from multi-agent debates, capturing both correct solutions and explicit self-correction behaviors.

## 3.3 Distillation Methods

We explore three strategies to supervise the student model $\pi_\theta$, aiming to distill not only correct solutions but also concise and structured reasoning (Zhong et al., 2025).

### 3.3.1 Reasoning-Enhanced SFT

Reasoning-enhanced SFT not only takes the final answers, but also the reasoning traces as supervision. The student model is trained to maximize the likelihood of the raw multi-agent reasoning traces. Given an input $x$ and a successful trace $r = (r_1, \ldots, r_n)$ to final answer $y^*$, the objective is:

$$\mathcal{L}_{\text{SFT}}(\theta) = -\mathbb{E}_{(x,r,y^*) \sim \mathcal{D}} \mathcal{L}_{\text{res}} + \mathcal{L}_{\text{ans}}, \tag{1}$$

where

$$\mathcal{L}_{\text{res}} = \sum_{t=1}^{|r|} \log p_\theta(r_t \mid r_{<t}, x), \tag{Reasoning}$$

$$\mathcal{L}_{\text{ans}} = \log p_\theta(y^* \mid r, x). \tag{Answer}$$

The objective consists of two components: a reasoning loss $\mathcal{L}_{\text{res}}$, which optimizes the model's ability to generate a coherent sequence of intermediate rationales $r$, and an answer loss $\mathcal{L}_{\text{ans}}$, which ensures the final prediction $y^*$ is grounded in both the input context $x$ and the preceding reasoning path. SFT assesses the student's capacity for basic imitation of the multi-agent reasoning style, testing whether it can maintain structural consistency across the entire generation sequence.

### 3.3.2 Distillation with Data Augmentation

To fully exploit the diversity and multi-perspective nature of the multi-agent debate, we implement a *Correctness-First Diverse Extraction* strategy.

**Selection and LLM-based Extraction.** For each problem $x$, we filter the agents to define a set of successful contributors $\mathcal{A}_{\text{correct}} \subseteq \mathcal{A}$. We then utilize a high-capacity teacher LLM as a "distiller" to parse the raw debate logs of these agents. The distiller is tasked with extracting $k \in \{1, 2, 3\}$ reasoning trajectories that are: (1) **Correct:** They must lead strictly to the ground-truth $y^*$. (2) **Diverse:** The teacher is prompted to select traces that employ distinct

mathematical identities, different logical heuristics, or varied starting assumptions (details in Appendix A). The student is trained on this augmented set $\mathcal{D}_{\text{aug}}$ using:

$$\mathcal{L}_{\text{Aug}}(\theta) = -\frac{1}{k} \sum_{i=1}^{k} \sum_{t=1}^{T} \log p_\theta(y_t \mid y_{<t}, r_i, x). \qquad (2)$$

This forces the model to learn multiple valid paths to the same solution, improving its robustness and generalization.

### 3.3.3 Process-Aware Distillation

The third method treats distillation as a reinforcement learning problem, using a Process Reward Model (PRM) (Lightman et al., 2023; Wang et al., 2024a) for granular, step-level supervision that reinforces the intermediate logical transitions found in multi-agent debates. Specifically, we train a PRM $R_\phi$ to predict the probability that a reasoning step is correct. To ensure the PRM reliably captures the nuances of the debate, we use a two-stage curriculum: **Stage I: Feature Alignment (Backbone Frozen).** We initialize $R_\phi$ with the weights of the student model and freeze all layers except for the final one and the reward head. This prevents the loss of pre-trained linguistic features while the reward head learns to map existing representations to the correctness labels $z_t \in \{0, 1\}$. **Stage II: Full Specialization (Backbone Unfrozen).** We unfreeze the entire backbone for end-to-end fine-tuning. This allows the model to develop specialized internal attention patterns for detecting logical fallacies.

We design the Process Reward Model (PRM) to provide step-level supervision via a **contrastive reward objective** rather than standard binary cross-entropy. This encourages the model to assign higher rewards to reasoning steps that are more consistent with the multi-agent debate consensus, reflecting relative correctness rather than absolute labels (details in Appendix B.1).

**Reinforcement Learning via GRPO.** Finally, we fine-tune the student policy $\pi_\theta$ using Group Relative Policy Optimization (GRPO) (Shao et al., 2024). Given an input $x \sim \mathcal{D}$, we sample a group of $G$ reasoning outputs $o_1, \ldots, o_G$ from a fixed behavior policy $\pi_{\text{old}}$, which is a snapshot of the student policy before the current update. GRPO optimizes the policy by comparing the rewards of these outputs within the group, removing the need for a separate value function:

$$\mathcal{J}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{o_i\} \sim \pi_{\text{old}}} \left[ \frac{1}{G} \sum_{i=1}^{G} \mathcal{L}_i(\theta) - \beta \mathbb{D}_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}}) \right] \qquad (3)$$

where $\pi_{\text{ref}}$ denotes a fixed reference policy used to regularize the update, $\beta$ is the KL penalty coefficient, and $\mathcal{L}_i(\theta)$ is the clipped surrogate objective for output $o_i$, the surrogate objective for each output $o_i$ is:

$$\mathcal{L}_i(\theta) = \min \left( \rho_i(\theta) \hat{A}_i, \text{clip}(\rho_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right). \qquad (4)$$

In this formulation, $\rho_i(\theta) = \frac{\pi_\theta(o_i|x)}{\pi_{\text{old}}(o_i|x)}$ denotes the probability ratio. The advantage $\hat{A}_i$ ensures **reward consistency** by normalizing the step-wise aggregate score $R_\phi(o_i)$ from the trained PRM:
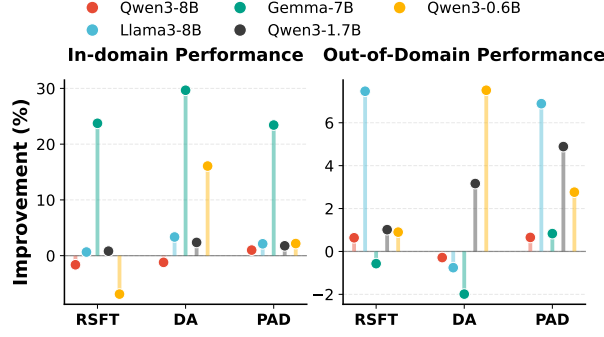
$$\hat{A}_i = \frac{R_\phi(o_i) - \mu_R}{\sigma_R}, \qquad (5)$$

where $\mu_R$ and $\sigma_R$ are the mean and standard deviation of the PRM scores within the group of $G$ sampled outputs.
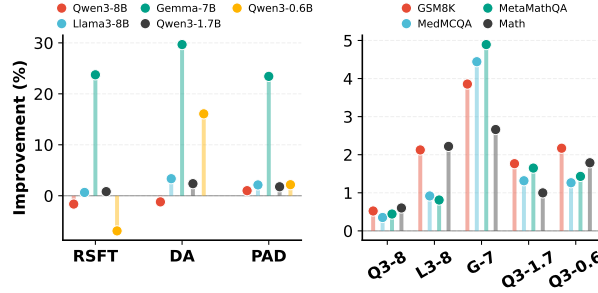
## 4 Experiments

### 4.1 Experimental Setup

**Models.** We conduct evaluation across three major model families: (1) **Qwen 3** Yang et al. (2025): Qwen3-32B, Qwen3-8B, Qwen3-1.7B, and Qwen3-0.6B; (2) **Gemma 3** Team et al. (2025): Gemma3-27B-it and Gemma-7B; and (3) **Llama 3** Dubey et al. (2024): Llama3-8B-Instruct. Particularly, we perform distillation from larger models (Qwen3-32B, and Gemma3-27B-it) to smaller ones (Qwen3-8B, Qwen3-1.7B, Qwen3-0.6B, Llama3-8B, and Gemma-7B) for comprehensive evaluation. We additionally evaluate multimodal LLM distillation as a preliminary study in §4.6.

Legend: Qwen3-8B, Gemma-7B, Qwen3-0.6B, Llama3-8B, Qwen3-1.7B

**In-domain Performance** | **Out-of-Domain Performance**

(a) Performance on in-domain (left) and OOD (right) datasets



Legend: Qwen3-8B, Gemma-7B, Qwen3-0.6B, Llama3-8B, Qwen3-1.7B; GSM8K, MetaMathQA, MedMCQA, Math

(b) Performance by datasets (left) and models (right)

Figure 3: Distillation from Qwen3-32B to different student models.

**Datasets.** We adopt diverse benchmarks: (1) MATH (Lightman et al., 2023) and GSM8K (Cobbe et al., 2021) for *mathematical reasoning*; (2) MedMCQA (Pal et al., 2022) for *domain-specific knowledge*; (3) MetaMathQA (MMQA) (Yu et al., 2023) for augmented math tasks; and (4) QASPER (Dasigi et al., 2021), HotpotQA (Yang et al., 2018) and QMSum (Zhong et al., 2021) for *multi-hop and long-form reasoning*.

**Comparison Methods.** We mainly compare three distillation methods in §3: RSFT, DA, and PAD, with single-agent and the vanilla multi-agent debate as baselines. Our primary study are based on 5 agents following existing work (Liang et al., 2024) and §4.3 presents the scaling exploration. By varying teacher-student models and training-test datasets on different distillation methods, we conducted 120 experiments in total. Accuracy on the test dataset is adopted as the primary evaluation metric while other metrics for further analysis are introduced in later sections.

## 4.2 Main Results

Figure 3 presents the results by distilling from Qwen3-32B to different student models across various datasets and the right-hand-side of Figure 2 shows the average comparison results. Other results are in Appendix G.2. Overall speaking, AgentArk significantly improves the performance of a single agent by 4.8% and only slightly worse than the vanilla MAS. More insightful findings are as follows.

**Performance across Distillation Methods.** Figure 3a show the results tested on GSM8K and MedMCQA, respectively. (1) **ID vs. OOD.** While distillation shows improvement in both settings, the gain in ID is more significant than in the OOD setting (30% vs. 7% maximum improvement and 4-6% vs. 1-3% on average). This is expected as OOD is more challenging, suggesting more future work. (2) **Methods.** Different methods show varied performance across datasets and tasks. While RSFT and DA sometimes provide improvements, their gains are inconsistent and can fluctuate by dataset. In contrast, PAD consistently yields performance improvements, demonstrating its robustness and reliable transfer of reasoning capabilities. (3) **Compatibility.** Different distillation strategies are mutually compatible and can be composed to yield consistent gains (Appendix F).

**Performance across Student Models.** We evaluate the generalization of all three distillation strategies by fixing the teacher model and varying the student model family. The results are shown in Figure 3a and 3b (left). (1) **Same-family distillation (Qwen-3)**. When both teacher and student are drawn from the Qwen-3 family, distillation yields stable but relatively moderate improvements, with smaller variants (1.7B and 0.6B) consistently benefiting more than the 8B
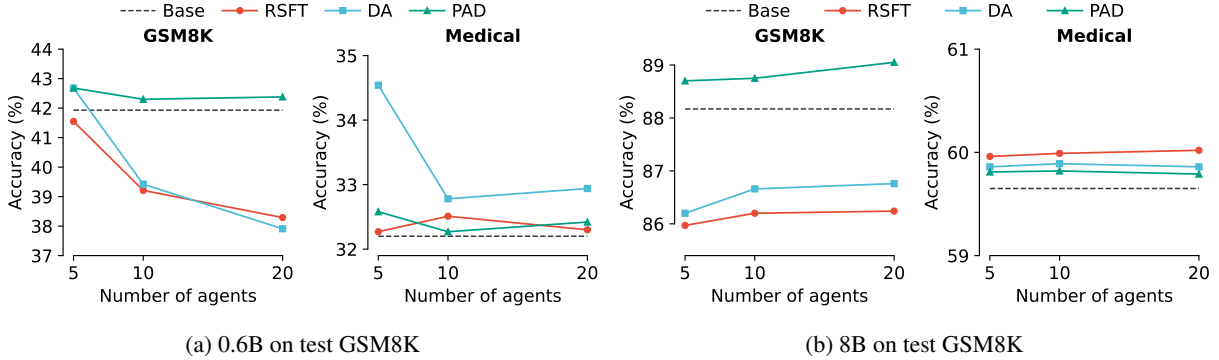
6

Figure 4: Effect of agent scale $(5, 10, 20)$ on distillation performance evaluated on GSM8K and MedMCQA.

model, indicating that same-family distillation mainly alleviates capacity constraints rather than inducing substantial representational changes. (2) **Cross-family distillation.** When distillation is performed across different model families, we observe larger and more consistent performance gains, particularly for Gemma-7B and LLaMA-3-8B, indicating that heterogeneous architectures benefit more from transferred reasoning signals. (3) **Effect of PRM supervision.** Across both large and small student models, PRM-based distillation consistently improves performance on ID and OOD tasks, demonstrating that the learned gains are not limited to scale or data distribution. Notably, the improvements persist under distribution shift, suggesting that PRM supervision primarily transfers reasoning behaviors rather than merely improving surface-level alignment to training data (Detailed results in Appendix G.2). We further conduct detailed PRM ablation studies to analyze its effects (Appendix C).

**Performance across datasets.** For a fixed student model, distillation consistently improves performance across all benchmark datasets, but the magnitude of gains varies (Figure 3b). MetaMathQA exhibits the largest improvements, followed by GSM8K, while Math shows moderate gains and MedMCQA benefits the least. This pattern suggests that the observed improvements are primarily driven by enhanced reasoning capabilities rather than dataset-specific overfitting. In particular, the strong gains on MetaMathQA and GSM8K indicate that these datasets contain reasoning-intensive problems, such as multi-step logical or arithmetic chains, which benefit most from transferred reasoning knowledge. By contrast, Math, with more formulaic or domain-specific tasks, shows moderate improvement, and MedMCQA, which relies heavily on specialized medical factual knowledge, benefits the least, implying that distillation contributes less when reasoning is minimal. Overall, this analysis highlights that our distillation strategies effectively enhance generalized reasoning skills, with larger impact on datasets that require complex reasoning rather than memorization.

## 4.3 Scaling and Data Dynamics

**Scaling the Number of Agents.** While the primary experiments employ 5 agents following existing work, we further scale to 10 and 20 agents to study how increasing teacher diversity and interaction complexity affects distillation. We perform evaluation by distilling from a fixed Qwen3-32B model into two different student models: Qwen3-8B and Qwen3-0.6B. The results are shown in Figure 4.

For the smaller student Qwen3-0.6B, scaling beyond 5 does not yield additional benefits and, in some cases, leads to performance degradation. We attribute this to the limited representational capacity of the student: as the teacher ensemble becomes more diverse and produces more complex or longer reasoning trajectories, the student is unable to faithfully absorb and generalize this information. By contrast, the larger Qwen3-8B benefits modestly from scaling. Performance improves consistently as the number of agents increases. However, the incremental gains diminish at higher scales, indicating that once the student can effectively utilize the ensemble's reasoning, additional agents offer limited benefit. These results highlight that multi-agent distillation is bounded by student capacity: smaller models saturate quickly and may even struggle with overly diverse teacher signals, whereas larger models can better exploit richer supervision from multiple agents. Overall, scaling teachers is effective only when matched to the student's representational capacity (More results in Appendix D.2).

**Data Quantity vs. Quality.** We investigate the trade-off between data quantity and quality by varying the amount of training data in Figure 5. Across datasets, increasing the training data does not lead to monotonic improvements. For both RSFT and DA, performance exhibits high variance as data scale grows: moderate data sizes can yield gains, while further scaling often leads to stagnation or even degradation, particularly for GSM8K and MedMCQA. In contrast, PAD

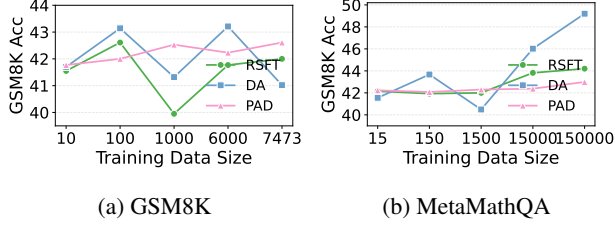|                                | (a) GSM8K | (b) MetaMathQA |
|--------------------------------|-----------|----------------|

Figure 5: Data scaling behavior of distillation from Qwen3-32B to Qwen3-0.6B, showing target model performance as a function of training data size across datasets.

Table 1: Quantitative evaluation of reasoning quality.

| Metric | Single | RSFT | DA | PAD |
|---|---|---|---|---|
| Avg NLL (↓) | 0.6529 | 0.4092 | 0.4449 | 0.5876 |
| Perplexity (↓) | 1.9211 | 1.6388 | 1.5603 | 1.7996 |
| Step Decomposition (↑) | 2.75 | 3.13 | 3.38 | 3.23 |
| Intermediate Verification (↑) | 2.41 | 3.48 | 4.04 | 4.07 |
| Error Localization (↑) | 1.97 | 2.19 | 2.91 | 2.78 |
| Reasoning Coherence (↑) | 1.88 | 2.25 | 3.07 | 3.96 |

demonstrates significantly more stable behavior across data scales. They consistently achieve competitive performance and avoid the sharp fluctuations observed with raw data scaling. This suggests that for capacity-limited students, reasoning quality rather than data volume is the primary bottleneck. Excessive or noisy supervision from large-scale MAS outputs can overwhelm the target model, whereas PRM-guided distillation preserves high-signal reasoning trajectories that enable more reliable transfer.

**Training Time.** The offline training cost is shown in Appendix D.3, which is mild and does not affect the inference efficiency of a single agent.

## 4.4 Analysis of the Reasoning Quality

To better understand why distillation works, we conduct a focused analysis on the reasoning quality by combining perplexity analysis and LLM-based qualitative evaluation with detailed case study. For each dataset, we randomly sample 100 examples.

**Perplexity Analysis.** We first measure the perplexity ((Bengio et al., 2003); details in Appendix E.1) of distilled models on held-out samples from GSM8K, focusing on reasoning tokens using Qwen3-32B and 0.6B as the teacher and student models, respectively. Perplexity measures how well a model predicts the next token in a sequence; lower values indicate the model's reasoning steps are more predictable and coherent, aligning with our goal of producing structured reasoning trajectories. The results are in Table 1. All distilled models achieve substantially lower reasoning perplexity than the single ones, indicating that distillation improves the predictability of the student's reasoning. This suggests that distillation not only helps the model produce more fluent outputs, but also encourages more structured and consistent reasoning trajectories.

**Evaluation of Reasoning Quality.** For more quantitative evaluation, we employ InternLM-2.5-20b-chat (InternLM2, 2024) as an automatic evaluator to score model outputs along four dimensions: step decomposition (Hwang et al., 2025), intermediate verification (Zheng et al., 2025), error localization (Mukherjee et al., 2025), and overall reasoning coherence (Lee & Hockenmaier, 2025). Step decomposition evaluates whether the model explicitly breaks problems into logical substeps; intermediate verification measures self-checking at each step; error localization assesses the model's ability to identify and correct mistakes; and overall reasoning coherence captures the consistency and logical flow of the solution. As shown in Table 1, across all dimensions, PAD achieves the highest scores, indicating that it most effectively preserves explicit multi-step structure, self-checking behavior, and coherent reasoning flows. In contrast, DA shows moderate improvements, particularly in intermediate verification and error localization, suggesting that it captures surface-level reasoning structure without fully inheriting reflective reasoning behaviors. RSFT only outperforms the baseline, implying that direct reasoning-level supervision alone is insufficient.

**Case Study.** We further presented detailed comparisons of reasoning states in Appendix E. This example demonstrates that the distilled single-agent model acquires the multi-agent reasoning patterns: its reasoning is more structured,

Table 2: Robustness evaluation on TruthfulQA.

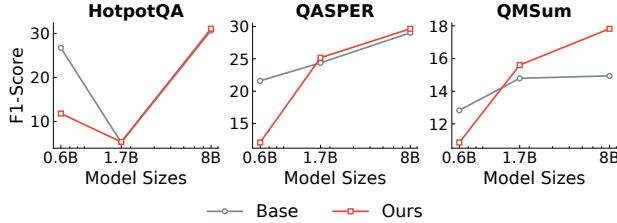| Metric | Single | PAD | DA | RSFT |
|---|---|---|---|---|
| BLEU (↑) | 0.6034 | 0.6634 | 0.6353 | 0.6059 |
| ROUGE-1 (↑) | 0.6144 | 0.6659 | 0.6389 | 0.6157 |
| ROUGE-2 (↑) | 0.5704 | 0.6414 | 0.5961 | 0.5777 |
| ROUGE-L (↑) | 0.6132 | 0.6573 | 0.6401 | 0.6157 |



Figure 6: Performance on 3 open-ended datasets.

logically coherent, and self-consistent, producing the correct answer without the repeated self-corrections seen in the baseline single-agent model.

## 4.5  Robustness and Generalization

We further analyze the robustness and generalization. Robustness is particularly important in multi-agent distillation, where the teacher ensemble may exhibit heterogeneous reasoning styles and varying levels of reliability.

**Robustness.** We evaluate the Qwen3-8B student model distilled from the Qwen3-32B teacher using the GSM8K dataset on TruthfulQA (Lin et al., 2022). It measures the model's ability to maintain factual accuracy and coherent reasoning. As shown in Table 2, all distillation methods improve performance over the base model, indicating their robustness and resilience to catastrophic forgetting. In particular, PAD achieves the highest scores across all metrics, indicating more robust reasoning and better retention of factual correctness. These results suggest that MAS distillation not only enhances average accuracy but also increases robustness, allowing the student model to generalize more reliably to unseen or challenging tasks.

**Open-ended Generalization.** To assess cross-domain generalization, we evaluate whether reasoning skills learned from mathematical problem-solving transfer to other complex reasoning tasks such as open-ended questions. We train the models on a single source dataset GSM8K (math reasoning) and evaluate them on three out-of-domain (OOD) datasets spanning diverse tasks: HotpotQA Yang et al. (2018) (multi-hop reasoning), QASPER Dasigi et al. (2021) (long-context understanding), and QMSum Zhong et al. (2021) (summarization). For comprehensive open-ended assessment, we leverage F1 scores and ROUGE-1/2/L Lin (2004) to measure lexical overlap between each prediction and the ground-truth, and BERTScore Zhang et al. (2019) for semantic similarity. Figure 6 shows that AgentArk significantly enhances cross-dataset reasoning transfer, particularly for larger models, which consistently improve performance on diverse OOD tasks such as multi-hop QA, long-context understanding, and summarization. These results indicate that AgentArk strengthens general reasoning capabilities rather than merely fitting dataset-specific pattern on open-ended tasks. More details on other metrics are in Appendix D.1.

## 4.6  Distillation of Multimodal LLMs

We finally study the transferability of MAS reasoning to **Multimodal LLMs (MLLMs)**. Specifically, we distill from Qwen2.5-VL-32B-Instruct into the smaller Qwen2.5-VL-3B-Instruct (Bai et al., 2025). Figure 7 shows that MAS distillation remains effective for MLLMs, despite being trained solely on text-only reasoning datasets. Across both benchmarks, PAD consistently achieves the strongest or near-strongest performance, indicating that distilling process-level reasoning signals generalizes better than RSFT and DA. Notably, the absolute gains are modest compared to text-only settings, which is expected since the model is not explicitly trained on MLLM reasoning. Nevertheless, it suggest that MAS distillation primarily enhances the model's internal reasoning competence, which can be reused by MLLMs. These results indicate that AgentArk captures model-agnostic reasoning patterns that are transferable beyond the training modality.
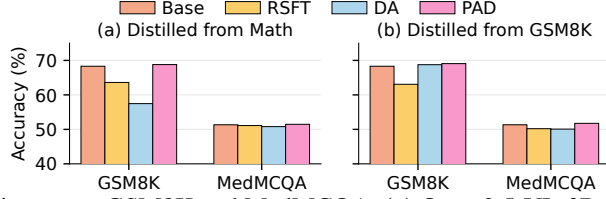
Figure 7: Multimodal distillation on on GSM8K and MedMCQA. **(a)** Qwen2.5-VL-3B-Instruct distilled from Math, **(b)** Qwen2.5-VL-3B-Instruct distilled from GSM8K.

# 5 Discussion

Our experiments reveal several insights and avenues for future research in multi-agent reasoning distillation.

*(1) Multi-agent distillation strategies:* Structured distillation, particularly PRM-guided methods, effectively transfers complex reasoning behaviors to smaller or multimodal models. Future work could explore adaptive distillation strategies that dynamically select reasoning trajectories based on task complexity, model capacity, or ensemble diversity. Additionally, methods that weigh intermediate self-corrections differently from final answers may further improve learning efficiency and reasoning fidelity.

*(2) Process modeling and policy optimization:* Larger PRMs yield stronger gains even for small target models. Future research could investigate modular or hierarchical PRMs, which allow selective guidance for different reasoning components, and alternative policy optimization techniques that balance stability and scalability beyond PPO and GRPO.

*(3) Scaling agent ensembles and data quality:* Increasing the number of debating agents benefits larger models but shows diminishing returns for smaller models. Similarly, training data volume alone is insufficient; high-quality reasoning trajectories are crucial. Future work could focus on adaptive ensemble scaling or selective trajectory sampling that optimizes the trade-off between diversity and learnability.

*(4) Transfer to multimodal models:* Distilled reasoning behaviors generalize to multimodal LLMs, suggesting potential for cross-modal knowledge transfer. Future directions include extending structured distillation to richer modalities (e.g., vision, audio) and evaluating the interplay between textual and non-textual reasoning pathways.

*(5) Implication for foundation models:* Our study suggests that in addition to developing advanced MAS algorithms, it is promising to leverage MAS to augment single LLMs, which will dramatically reduce cost with improved performance. Moreover, the distillation on small models implies that small language models can be significantly enhanced by MAS, highlighting promise for lightweight and cost-effective deployment of foundation models.

**Limitations.** This work has the following limitations. First, experiments are limited to a subset of reasoning benchmarks and multimodal models, evaluating additional tasks and modalities could help assess broader applicability. Second, we focus on a specific set of distillation pipelines. Exploring alternative or hybrid approaches may provide further insights into multi-agent knowledge transfer mechanisms. Third, while AgentArk is agnostic to MAS strategies, we focus on debate in this work. More insightful findings could be obtained by applying AgentArk to other MAS algorithms.

# 6 Conclusion

We introduced AgentArk, a framework for distilling multi-agent reasoning into a single agent. Our experiments demonstrate that structured distillation, particularly PRM-guided methods, enables smaller models to approximate complex reasoning behaviors while maintaining efficiency and generalization across diverse tasks. With the increasing attention to LLM efficiency and multi-agent systems, we hope our findings can provide insights for future research.

# Acknowledgment

## Impact Statements

AgentArk achieves the benefits of multi-agent reasoning through efficient distillation rather than expensive test-time orchestration, reducing latency and deployment cost for reasoning-heavy applications where multi-agent inference is impractical, such as on-device settings. By lowering the computational barrier, it broadens access to stronger agentic reasoning capabilities for resource-constrained users. Meanwhile, distilled students may inherit less desirable behaviors from teacher models, such as persuasive but biased or logically unsound rationales. To mitigate these risks, we recommend enforcing correctness checks for both PRMs and RL-finetuned models, and auditing distilled agents for hallucination, bias, and harmful content prior to deployment. Future work may extend AgentArk to more real-world agentic settings, including interactive tool-usage workflows or safety-critical decision support tasks.

## References

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155, 2003.

Huaben Chen, Wenkang Ji, Lufeng Xu, and Shiyu Zhao. Multi-agent consensus seeking via large language models. *arXiv preprint arXiv:2310.20151*, 2023.

Justin Chen, Swarnadeep Saha, and Mohit Bansal. Reconcile: Round-table conference improves reasoning via consensus among diverse llms. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7066–7085, 2024a.

Justin Chih-Yao Chen, Swarnadeep Saha, Elias Stengel-Eskin, and Mohit Bansal. Magdi: Structured distillation of multi-agent interaction graphs improves reasoning in smaller language models. *arXiv preprint arXiv:2402.01620*, 2024b.

Weize Chen, Jiarui Yuan, Chen Qian, Cheng Yang, Zhiyuan Liu, and Maosong Sun. Optima: Optimizing effectiveness and efficiency for llm-based multi-agent system. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 11534–11557, 2025.

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

Petru L Curşeu, Nicoleta Meslec, Helen Pluut, and Gerardus JM Lucas. Cognitive synergy in groups and group-to-individual transfer of decision-making competencies. *Frontiers in psychology*, 6:1375, 2015.

Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. A dataset of information-seeking questions and answers anchored in research papers. *arXiv preprint arXiv:2105.03011*, 2021.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2023.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv:2407.21783*, 2024.

Sugyeong Eo, Hyeonseok Moon, Evelyn Hayoon Zi, Chanjun Park, and Heuiseok Lim. Debate only when necessary: Adaptive multiagent collaboration for efficient llm reasoning. *arXiv preprint arXiv:2504.05047*, 2025.

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

Shanshan Han, Qifan Zhang, Yuhang Yao, Weizhao Jin, and Zhaozhuo Xu. Llm multi-agent systems: Challenges and open problems. *arXiv preprint arXiv:2402.03578*, 2024.

Pengfei He, Yue Xing, Shen Dong, Juanhui Li, Zhenwei Dai, Xianfeng Tang, Hui Liu, Han Xu, Zhen Xiang, and Charu C Aggarwal. Comprehensive vulnerability analysis is necessary for trustworthy llm-mas. *arXiv preprint arXiv:2506.01245*, 2025.

Mahmood Hegazy. Diversity of thought elicits stronger reasoning capabilities in multi-agent debate frameworks. *arXiv preprint arXiv:2410.12853*, 2024.

Hyeon Hwang, Yewon Cho, Chanwoong Yoon, Yein Park, Minju Song, Kyungjae Lee, Gangwoo Kim, and Jaewoo Kang. Assessing llm reasoning steps via principal knowledge grounding. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pp. 19925–19948, 2025.

Ashrey Ignise and Yashika Vahi. Applications of multi-agent systems. *Authorea Preprints*, 2024.

InternLM2. Internlm2 technical report, 2024.

Minki Kang, Jongwon Jeong, Seanie Lee, Jaewoong Cho, and Sung Ju Hwang. Distilling llm agent into small models with retrieval and code tools. *arXiv preprint arXiv:2505.17612*, 2025.

Zixuan Ke, Yifei Ming, Austin Xu, Ryan Chin, Xuan-Phi Nguyen, Prathyusha Jwalapuram, Semih Yavuz, Caiming Xiong, and Shafiq Joty. Mas-orchestra: Understanding and improving multi-agent reasoning through holistic orchestration and controlled benchmarks. *arXiv preprint arXiv:2601.14652*, 2026.

Jiin Kim, Byeongjun Shin, Jinha Chung, and Minsoo Rhu. The cost of dynamic reasoning: Demystifying ai agents and test-time scaling from an ai infrastructure perspective. *arXiv preprint arXiv:2506.04301*, 2025a.

Yubin Kim, Ken Gu, Chanwoo Park, Chunjong Park, Samuel Schmidgall, A Ali Heydari, Yao Yan, Zhihan Zhang, Yuchen Zhuang, Mark Malhotra, et al. Towards a science of scaling agent systems. *arXiv preprint arXiv:2512.08296*, 2025b.

Tian Lan, Wenwei Zhang, Chengqi Lyu, Shuaibin Li, Chen Xu, Heyan Huang, Dahua Lin, Xian-Ling Mao, and Kai Chen. Training language models to critique with multi-agent feedback. *arXiv preprint arXiv:2410.15287*, 2024.

Jinu Lee and Julia Hockenmaier. Evaluating step-by-step reasoning traces: A survey. *arXiv preprint arXiv:2502.12289*, 2025.

Tiansi Li, Yuxuan Zhang, Yuxuan Liu, Yujia Zhang, Yujie Liu, Wayne Xin Zhao, and Ji-Rong Wen. Camel: Communicative agents for "mind" exploration. *arXiv preprint arXiv:2303.17760*, 2023.

Weizhen Li, Jianbo Lin, Zhuosong Jiang, Jingyi Cao, Xinpeng Liu, Jiayu Zhang, Zhenqiang Huang, Qianben Chen, Weichen Sun, Qiexiang Wang, et al. Chain-of-agents: End-to-end agent foundation models via multi-agent distillation and agentic rl. *arXiv preprint arXiv:2508.13167*, 2025.

Xiaoxiao Li. When single-agent with skills replace multi-agent systems and when they fail. *arXiv preprint arXiv:2601.04748*, 2026.

Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinagearth*, 1(1):9, 2024.

Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate. In *Proceedings of the 2024 conference on empirical methods in natural language processing*, pp. 17889–17904, 2024.

Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.

Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.

Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. In *Proceedings of the 60th annual meeting of the association for computational linguistics (volume 1: long papers)*, pp. 3214–3252, 2022.

Yi-Cheng Lin, Kang-Chieh Chen, Zhe-Yan Li, Tzu-Heng Wu, Tzu-Hsuan Wu, Kuan-Yu Chen, Hung-yi Lee, and Yun-Nung Chen. Creativity in llm-based multi-agent systems: A survey. *arXiv preprint arXiv:2505.21116*, 2025.

Jiaqi Liu, Chengkai Xu, Peng Hang, Jian Sun, Mingyu Ding, Wei Zhan, and Masayoshi Tomizuka. Language-driven policy distillation for cooperative driving in multi-agent reinforcement learning. *IEEE Robotics and Automation Letters*, 2025a.

Jun Liu, Zhenglun Kong, Peiyan Dong, Changdi Yang, Tianqi Li, Hao Tang, Geng Yuan, Wei Niu, Wenbin Zhang, Pu Zhao, et al. Structured agent distillation for large language model. *arXiv preprint arXiv:2505.13820*, 2025b.

Wendy Yaqiao Liu, Rui Jerry Huang, Anastasia Miin, and Lei Ding. Adaptive coopetition: Leveraging coarse verifier signals for resilient multi-agent llm reasoning. In *The 14th International Joint Conference on Natural Language Processing and The 4th Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics*, pp. 145–155, 2025c.

Tie Ma, Yixi Chen, Vaastav Anand, Alessandro Cornacchia, Amândio R Faustino, Guanheng Liu, Shan Zhang, Hongbin Luo, Suhaib A Fahmy, Zafar A Qazi, et al. Maestro: Multi-agent evaluation suite for testing, reliability, and observability. *arXiv preprint arXiv:2601.00481*, 2026.

Sagnik Mukherjee, Abhinav Chinta, Takyoung Kim, Tarun Anoop Sharma, and Dilek Hakkani-Tür. Premise-augmented reasoning chains improve error identification in math reasoning with llms. *arXiv preprint arXiv:2502.02362*, 2025.

Joaquin Navajas, Tamara Niella, Gerry Garbulsky, Bahador Bahrami, and Mariano Sigman. Aggregated knowledge from a small number of debates outperforms the wisdom of large crowds. *Nature Human Behaviour*, 2(2):126–132, 2018.

Thi-Nhung Nguyen, Linhao Luo, Thuy-Trang Vu, and Dinh Phung. The social cost of intelligence: Emergence, propagation, and amplification of stereotypical bias in multi-agent systems. *arXiv preprint arXiv:2510.10943*, 2025.

Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on health, inference, and learning*, pp. 248–260. PMLR, 2022.

Jiqun Pan, Zhenke Duan, Jiani Tu, Anzhi Cheng, and Yanqing Wang. Knowledge graph-guided multi-agent distillation for reliable industrial question answering with datasets. *arXiv preprint arXiv:2510.06240*, 2025.

Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Xuanhe Zhou, Yufei Huang, Chaojun Xiao, et al. Tool learning with foundation models. *ACM Computing Surveys*, 57(4):1–40, 2024.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv:1707.06347*, 2017.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

Xi Shi, Mengxin Zheng, and Qian Lou. Learning latency-aware orchestration for parallel multi-agent systems. *arXiv preprint arXiv:2601.10560*, 2026.

Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv:2503.19786*, 2025.

Wataru Toyokawa, Andrew Whalen, and Kevin N Laland. Social learning strategies regulate the wisdom and madness of interactive crowds. *Nature Human Behaviour*, 3(2):183–193, 2019.

David Wan, Justin Chen, Elias Stengel-Eskin, and Mohit Bansal. Mamm-refine: A recipe for improving faithfulness in generation with multi-agent collaboration. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 9882–9901, 2025.

Jinlong Wang, Yunting Wu, Xiaoyun Xiong, Yuanyuan Zhang, Zhihan Lyu, Ahmed Ghoneim, and Haoran Zhao. Fedlma: A federated learning framework integrating llm-based multi-agent reasoning with knowledge distillation. *IEEE Transactions on Consumer Electronics*, 2025a.

Kun Wang, Guibin Zhang, ManKit Ye, Xinyu Deng, Dongxia Wang, Xiaobin Hu, Jinyang Guo, Yang Liu, and Yufei Guo. Mas$^2$: Self-generative, self-configuring, self-rectifying multi-agent systems. *arXiv preprint arXiv:2509.24323*, 2025b.

Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9426–9439, 2024a.

Peng-Yuan Wang, Tian-Shuo Liu, Chenyang Wang, Ziniu Li, Yidi Wang, Shu Yan, Chengxing Jia, Xu-Hui Liu, Xinwei Chen, Jiacheng Xu, et al. A survey on large language models for mathematical reasoning. *ACM Computing Surveys*, 2025c.

Song Wang, Zhen Tan, Zihan Chen, Shuang Zhou, Tianlong Chen, and Jundong Li. Anymac: Cascading flexible multi-agent collaboration via next-agent prediction. *arXiv preprint arXiv:2506.17784*, 2025d.

Xiao Wang, Jia Wang, Yijie Wang, Pengtao Dang, Sha Cao, and Chi Zhang. Mars: toward more efficient multi-agent collaboration for llm reasoning. *arXiv preprint arXiv:2509.20502*, 2025e.

Xinchen Wang, Pengfei Gao, Xiangxin Meng, Chao Peng, Ruida Hu, Yun Lin, and Cuiyun Gao. Aegis: An agent-based framework for general bug reproduction from issue descriptions. *arXiv preprint arXiv:2411.18015*, 2024b.

Yiyang Wang, Chen Chen, Tica Lin, Vishnu Raj, Josh Kimball, Alex Cabral, and Josiah Hester. Companioncast: A multi-agent conversational ai framework with spatial audio for social co-viewing experiences. *arXiv:2512.10918*, 2025f.

Yiyang Wang, Yiqiao Jin, Alex Cabral, and Josiah Hester. Mascot: Towards multi-agent socio-collaborative companion systems. *arXiv:2601.14230*, 2026.

Tianxin Wei, Ting-Wei Li, Zhining Liu, Xuying Ning, Ze Yang, Jiaru Zou, Zhichen Zeng, Ruizhong Qiu, Xiao Lin, Dongqi Fu, et al. Agentic reasoning for large language models. *arXiv preprint arXiv:2601.12538*, 2026.

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

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pp. 2369–2380, 2018.

Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023.

Peiying Yu, Guoxin Chen, and Jingjing Wang. Table-critic: A multi-agent framework for collaborative criticism and refinement in table reasoning. *arXiv preprint arXiv:2502.11799*, 2025.

Yurun Yuan and Tengyang Xie. Reinforce LLM reasoning through multi-agent reflection. In Aarti Singh, Maryam Fazel, Daniel Hsu, Simon Lacoste-Julien, Felix Berkenkamp, Tegan Maharaj, Kiri Wagstaff, and Jerry Zhu (eds.), *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pp. 73701–73731. PMLR, 13–19 Jul 2025. URL https://proceedings.mlr.press/v267/yuan25l.html.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2019.

Zeyu Zhang, Quanyu Dai, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model-based agents. *ACM Transactions on Information Systems*, 43(6):1–47, 2025.

Zhonghan Zhao, Ke Ma, Wenhao Chai, Xuan Wang, Kewei Chen, Dongxu Guo, Yanting Zhang, Hongwei Wang, and Gaoang Wang. Do we really need a complex agent system? distill embodied agent into a single model. *arXiv preprint arXiv:2404.04619*, 2024.

Xinyi Zheng, Ningke Li, Xiaokun Luan, Kailong Wang, Ling Shi, Meng Sun, and Haoyu Wang. Beyond correctness: Exposing llm-generated logical flaws in reasoning via multi-step automated theorem proving. *arXiv preprint arXiv:2512.23511*, 2025.

Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, et al. Qmsum: A new benchmark for query-based multi-domain meeting summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5905–5921, 2021.

Yiwu Zhong, Zi-Yuan Hu, Yin Li, and Liwei Wang. Rethinking chain-of-thought reasoning for videos. *arXiv:2512.09616*, 2025.

Xiaofeng Zhou, He-Yan Huang, and Lizi Liao. Debate, reflect, and distill: Multi-agent feedback with tree-structured preference optimization for efficient language model enhancement. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 9122–9137, 2025.

<div align="center">

**Appendix**

</div>

# A   Data Generation Methods and Pipeline

This section provides detailed descriptions of the data generation pipeline, dataset composition, and statistics used in our experiments.

## A.1   Dataset Overview

We construct the multi-agent distillation dataset using four core benchmarks, focusing on training the student models:

- **Mathematical reasoning:** GSM8K and MATH, covering arithmatic, algebraic, and multi-step symbolic reasoning.

- **Augmented math reasoning:** MetaMathQA, providing additional multi-step problems and varied solution strategies.

- **Domain-specific knowledge:** MedMCQA, focusing on medical exam-style multiple-choice questions.

The remaining benchmarks, HotpotQA, QAPER, and QMSum, are reserved exclusively for zero-shot generalization evaluation and are not used during training or data augmentation.

For the training datasets, after multi-agent data generation, correctness filtering, and diversity-based selection, the final distillation dataset contains approximately 342k unique input questions and 2M reasoning trajectories.

Table 3 reports detailed statistics for each training dataset, including the number of questions, retained debates, and augmented trajectories.

Table 3: Dataset statistics for multi-agent distillation. Q indicates the number of unique training questions, while T denotes the number of correct reasoning trajectories extracted per dataset.

| Source Model | GSM8K(Q / T) | Math(Q / T) | MetaMathQA(Q / T) | Medmcqa(Q / T) |
|---|---|---|---|---|
| Qwen3-8B | 7473/44,838 | 500/3000 | 151k/906k | 183k/1.1M |
| Qwen3-32B | 7473/44,838 | 500/3000 | 151k/906k | 183k/1.1M |
| Gemma3-27B-It | 7473/44,838 | 500/3000 | 151k/906k | 183k/1.1M |

## A.2 Multi-Agent Data Generation

To generate rich reasoning supervision, we adopt a multi-agent debate mechanism to produce diverse solution trajectories. For each input problem $x$, we initialize a set of $n = 5$ agents, plus a agent that summarize the final answer at the end. $\mathcal{A} = \{a_1, a_2, \ldots, a_n\}$ that share the same underlying teacher language model but operate with independent generation contexts. Agents engage in a debate for up to $K = 3$ rounds. In the first round, each agent independently produces an initial reasoning trace and final answer. In subsequent rounds, each agent observes the reasoning traces generated by other agents in the previous round and is encouraged to revise its own solution by identifying potential errors, alternative solution paths, or overlooked details. Formally, in round $k$, agent $a_i$ generates a reasoning trajectory $\tau_{i,k}$ conditioned on the problem $x$ and the peer trajectories $\{\tau_{j,k-1}\}_{j \neq i}$. The complete debate log for an input $x$ is denoted as

$$\mathcal{L}_x = \{\tau_{i,k} \mid i \in [1, n], k \in [1, K]\}.$$

This interaction protocol promotes reasoning behaviors such as self-correction, hypothesis revision, and cross-verification, yielding a diverse set of candidate solutions that go beyond single-pass generation.

## A.3 Correctness Filtering

For each input problem $x$, we collect responses generated by a set of agents $\mathcal{A}$, where each agent produces a final answer along with its associated reasoning trajectory. We first apply a *correctness filtering* step to identify agents whose final answers are valid.

To ensure consistent and scalable evaluation across all tasks, we use Qwen2.5-72B-Instruct as an automatic verifier to assess whether a candidate final answer matches the ground-truth solution under task-specific evaluation rules (e.g., exact match or normalized equivalence for GSM8K-style problems). This model is prompted to perform answer verification only and does not take the reasoning trajectory into account during this stage.

An agent $a \in \mathcal{A}$ is included in the set of successful contributors $\mathcal{A}_{\text{correct}}$ if its final answer is verified as correct by the verifier model. This step ensures that all candidate responses used for distillation are answer-correct, independent of their reasoning styles.

Formally,

$$\mathcal{A}_{\text{correct}} = \{a \in \mathcal{A} \mid \text{Answer}(a, x) \text{ is verified as correct}\}. \tag{6}$$

If $|\mathcal{A}_{\text{correct}}| < 2$, the corresponding problem instance is excluded from the data augmentation process, as it does not provide sufficient diversity among correct solutions.

## A.4 Data Selection for Distillation with Data Augmentation

### A.4.1 Divergent-Reasoning Selection via LLM-based Judgment

Although all agents in $\mathcal{A}_{\text{correct}}$ arrive at the same correct answer, their reasoning processes may differ substantially. To capture this variation, we adopt a Correctness-First Diverse Extraction strategy that emphasizes reasoning divergence under answer agreement.

For each problem $x$, we consider the set of reasoning traces produced by agents in $\mathcal{A}_{\text{correct}}$. We then employ Qwen2.5-72B-Instruct as an auxiliary judge to identify responses whose reasoning patterns are meaningfully distinct while preserving correctness.

The judge model is provided with the problem $x$, the ground-truth answer, and multiple candidate reasoning traces from $\mathcal{A}_{\text{correct}}$. It is instructed to select responses that satisfy the following criteria:

1. The final answer is correct.

2. The reasoning process exhibits structural differences compared to other correct responses (e.g., different decomposition orders, intermediate representations, or solution paths).

Notably, the judge model is not used to rank responses by quality or correctness. Instead, it operates solely to assess reasoning diversity conditioned on answer correctness, which reduces bias toward specific reasoning styles.

### A.4.2 Final Augmented Dataset Construction

The selected responses are incorporated into the distillation dataset as augmented supervision signals. By construction, the resulting dataset satisfies two properties:

- **Answer consistency:** all augmented samples preserve the correct final answer.

- **Reasoning diversity:** multiple valid reasoning trajectories are retained for the same input.

This data selection procedure enables the student model to learn from a richer set of correct problem-solving behaviors without introducing incorrect or contradictory supervision.

# B    Methods

## B.1    Process Reward Model: Contrastive Loss Design

Denoting $\sigma(\cdot)$ as the sigmoid function, we design the Process Reward Model (PRM) loss as a **contrastive loss** rather than standard binary cross-entropy. This choice encourages the model to assign higher rewards to reasoning steps that are more consistent with the multi-agent debate consensus, reflecting relative correctness rather than absolute labels.

Specifically, given a positive reasoning step $r_t^+$ and a set of negative steps $\{r_t^-\}$ sampled from other agents' reasoning trajectories, the PRM loss is defined as:

$$\mathcal{L}_{\text{PRM}}(\phi) = -\sum_t \log \text{softmax}\left(\frac{\sigma(R_\phi(r_t^+)), \{\sigma(R_\phi(r^-))\}_{r^- \in \mathcal{N}_t}}{\tau}\right), \tag{7}$$

where $\tau$ is a temperature hyperparameter controlling the sharpness of the contrastive distribution. This formulation encourages the PRM to score more consistent reasoning steps higher while down-weighting less consistent or contradictory steps.

## B.2    Ablation: PPO Comparison

In addition to GRPO, we also experimented with standard Proximal Policy Optimization (PPO) (Schulman et al., 2017) to verify the benefits of group-wise relative updates. The overall setup is identical to GRPO Shao et al. (2024) training, including the use of the Process Reward Model (PRM) to provide step-level supervision. The student policy $\pi_\theta$ is fine-tuned using the conventional PPO clipped objective:

$$\mathcal{L}_{\text{PPO}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, o \sim \pi_{\text{old}}}\left[\min\left(\rho(\theta)\hat{A}, \text{clip}(\rho(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}\right) - \beta \mathbb{D}_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}})\right], \tag{8}$$

where $\rho(\theta) = \frac{\pi_\theta(o|x)}{\pi_{\text{old}}(o|x)}$ is the probability ratio and $\hat{A}$ is the advantage derived from PRM scores:

$$\hat{A} = \frac{R_\phi(o) - \mu_R}{\sigma_R}, \tag{9}$$

with $\mu_R$ and $\sigma_R$ denoting the mean and standard deviation of PRM scores across sampled outputs for a given input $x$.

Unlike GRPO, PPO treats each sample independently rather than comparing outputs within a group, which can lead to slower convergence and reduced reward consistency in multi-step reasoning tasks. The comparison between GRPO and PPO in our ablation (Table X) demonstrates that GRPO achieves more stable learning and higher final reasoning accuracy, validating the benefit of group-relative updates.

# C  Ablation Studies on PAD

## C.1  PRM Modeling and Policy Separation

We study the role of PRM by explicitly separating the model used to learn process rewards from the model used to train the final policy. In our setup, PRM is first trained independently, and its learned process-level guidance is then used to supervise the distillation of a target policy model. All experiments distill from the same Qwen3-32B, while the PRM and the policy are instantiated with different parameter scales to analyze the effect of PRM capacity under this decoupled training scheme. As shown in Table 4, using a smaller PRM model consistently leads to limited improvement, even when paired with the same policy model. In contrast, assigning a larger one yields gains across the GSM8K and MedMCQA test sets, even when the target policy remains lightweight. These results indicate that effective distillation relies more critically on the modeling size of the PRM than on the scale of the policy itself, highlighting the importance of explicit process modeling and role separation.

Table 4: Ablation on PRM and LLM role separation. All models are distilled from Qwen3-32B.

| Test set | PRM | LLM | GSM8K | MATH | MedMCQA | MMQA | Avg |
|---|---|---|---|---|---|---|---|
| GSM8K | 0.6B | 0.6B | 42.84 | 42.68 | 42.46 | 42.53 | 42.63 |
| | 8B | 0.6B | 42.84 | 42.53 | 42.91 | 42.23 | 42.63 |
| | 0.6B | 8B | 88.48 | 88.55 | 88.40 | 88.63 | 88.52 |
| | 8B | 8B | 88.63 | 88.70 | 88.48 | 88.56 | 88.59 |
| MedMCQA | 0.6B | 0.6B | 32.39 | 32.58 | 32.36 | 32.54 | 32.47 |
| | 8B | 0.6B | 32.49 | 32.56 | 32.35 | 32.20 | 32.40 |
| | 0.6B | 8B | 59.50 | 59.74 | 59.74 | 59.65 | 59.66 |
| | 8B | 8B | 59.77 | 59.81 | 59.86 | 59.77 | 59.80 |

## C.2  PPO vs. GRPO

To optimize the target policy during PRM-guided distillation, we compare Proximal Policy Optimization (PPO) (Schulman et al., 2017) (details in Section B.2) with GRPO. PRMs in both settings are trained identically in a first stage and then fixed. The two methods differ only in the second stage, where the target language model is fine-tuned using either PPO or GRPO under the same PRM-based reward signals. As shown in Table 5, PPO achieves marginally higher accuracy in most settings. This advantage is expected, as PPO employs a learned value function to provide a lower-variance, state-dependent baseline for policy updates, which can lead to more stable optimization under imperfect PRM rewards. In contrast, GRPO removes the value function and relies on group-relative comparisons, resulting in slightly noisier updates but significantly reduced computational overhead. Despite this difference, GRPO achieves performance comparable to PPO across all benchmarks, making it a scalable and effective alternative for large-scale PRM-guided distillation.

Table 5: Ablation on PPO vs GRPO.

| | Test set | GSM8K | MedMCQA | MMQA | Math | Avg |
|---|---|---|---|---|---|---|
| PPO | GSM8K | 53.37 | 53.15 | 53.24 | 52.62 | 53.10 |
| | MedMCQA | 49.56 | 50.08 | 49.77 | 49.41 | 49.71 |
| GRPO | GSM8K | 52.71 | 52.48 | 52.62 | 52.64 | 52.61 |
| | MedMCQA | 49.77 | 50.02 | 49.26 | 49.33 | 49.60 |

# D  Generalization

## D.1  Generalization on Out-of-Domain Datasets

To evaluate the robustness of AgentArk, we conduct experiments across three out-of-domain (OOD) datasets spanning diverse domains: HotpotQA Yang et al. (2018) (multi-hop reasoning), QASPER Dasigi et al. (2021) (long-context understanding), and QMSum Zhong et al. (2021) (summarization). The results in Figure 6 and Appendix Figure 8 & 9 reveal a clear scaling law for reasoning transferability.

The most substantial gains are observed in the larger 8B parameter class. Qwen3-8B and Llama3-8B demonstrate superior generalizability, effectively internalizing multi-agent thought processes to achieve consistent performance

lifts across all OOD tasks. This 'reasoning lift' is most pronounced in the complex QMSum task, where Qwen3-8B improves its F1-Score from 14.94 to 17.82 and its ROUGE-L Score from 15.72 to 17.41. Similarly, Llama3-8B sees its F1-Score increase from 13.05 to 14.92 on QMSum. These gains suggest that AgentArk is particularly potent for complex tasks such as high-level synthesis and human-centered dialog understanding, where the model must perform multi-step reasoning and aggregate information from extensive documents multiple perspectives.

Even mid-sized models like Qwen3-1.7B begin to show this positive trajectory, with F1 scores rising from 24.38 to 25.16 on QASPER and from 14.79 to 15.60 on QMSum. This indicates that once a model surpasses a critical capacity threshold, AgentArk serves as an enhancer for its inherent reasoning capabilities. In contrast, ultra-small models like Qwen3-0.6B have limited capacity to balance new reasoning patterns with its existing pre-trained knowledge base. Thus, AgentArk's multi-agent knowledge is most effective for models with sufficient cognitive capacity.



Figure 8: Scaling behavior of AgentArk vs. base models. Performance comparison on three out-of-distribution (OOD) datasets across the Qwen model family demonstrate the impact of increasing model size (0.6B to 8B) on generalizability.



Figure 9: Performance of AgentArk and the base models on out-of-domain datasets

## D.2 Data Scaling

We provide additional in Figure 10 data-scaling results using the MedMCQA dataset, complementing the scaling analysis presented in Section 4.3 of the main paper.

## D.3 Computation Cost

While AgentArk significantly reduces inference-time cost by eliminating multi-agent coordination, its training pipeline, particularly Process-Aware Distillation (PAD) introduces additional computational overhead. We report concrete

Figure 10: Data scaling behavior on MedMCQA for distillation from Qwen3-32B to Qwen3-0.6B, showing student model performance as a function of the amount of training data.

training costs for different distillation strategies to provide a realistic assessment of their practicality.

**Experimental Setup.** All experiments are conducted on NVIDIA H100 GPUs (80GB). Unless otherwise noted, we are reporting the student model has 8B parameters and is trained with a global batch size of 4 and identical sequence length across methods. All distillation approaches share the same multi-agent data generation cost, differences arise solely from the supervision signal and optimization procedure applied during student training.

**Training Cost.** Table 6 summarizes the additional training cost incurred by different distillation strategies.

Table 6: Training cost for different distillation strategies (8B student model).

| Method | Additional Training Components | GPUs | Time |
|---|---|---|---|
| RSFT | Supervised fine-tuning on single reasoning traces | $1 \times$ H100 | $\sim$6 hours |
| Reasoning DA | Supervised fine-tuning on augmented multi-trajectory reasoning data | $1 \times$ H100 | $\sim$8 hours |
| PAD (PRM + GRPO) | PRM training + GRPO-based policy optimization | $8 \times$ H100 | $\sim$20 hours |
| PRM training | Step-level process reward modeling | $8 \times$ H100 | $\sim$8 hours |
| GRPO | Policy optimization with PRM reward | $8 \times$ H100 | $\sim$12 hours |

**Inference-Time Cost.** In contrast to the increased offline training cost, inference with AgentArk requires only a single auto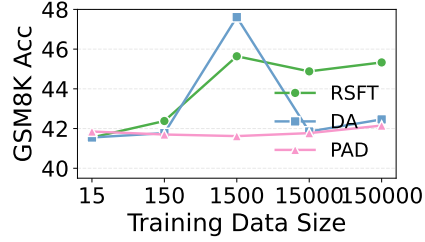regressive generation of the distilled student model. Compared to multi-agent debate, which typically involves multiple agent invocations and iterative coordination, AgentArk significantly reduces inference latency and GPU usage, making it suitable for real-time and resource-constrained deployment settings.

Overall, AgentArk trades increased offline computation for substantial reductions in inference-time cost and deployment complexity, consistent with common practices in large-scale model distillation and reinforcement learning.

# E   Case Study

> **Single Agent Reasoning**
>
> Let's think step by step... Step 1: Green eggs = 5. Step 2: Blue eggs = 2 * 5 = 10. Step 3: Pink eggs = 10 - 1 = 9. Step 4: Yellow eggs = 1/3 * 9 = 3. Total = 5 + 10 + 9 + 3 = 27.
> Wait, but the answer is 24... let me check again...
> Recalculating yellow eggs: 9 / 3 = 3. Total still 27. Something is wrong...
> Wait, but the answer is 24... let me check again...
> Recalculating yellow eggs: 9 / 3 = 3. Total still 27. Something is wrong...
> Wait, but the answer is 24... let me check again...
> Recalculating yellow eggs: 9 / 3 = 3. Total still 27. Something is wrong...
> ...

Agent 1: Step 1: Green eggs = 5 Step 2: Blue eggs = 2 * 5 = 10 Step 3: Pink eggs = 10 - 1 = 9 Step 4: Yellow eggs = 9 / 3 = 3 Total = 5 + 10 + 9 + 3 = 27

Agent 2: Step 1: Green eggs = 5 Step 2: Blue eggs = 10 Step 3: Pink eggs = 9 Step 4: Yellow eggs = 3 Total = 27 Verifies reasoning and intermediate calculations

Agent 3: Step 1: Green = 5 Step 2: Blue = 10 Step 3: Pink = 9 Step 4: Yellow = 3 Total = 27 Confirms each step independently...

Agent 4: Step 1: Green = 5 Step 2: Blue = 10 Step 3: Pink = 9 Step 4: Yellow = 3 Total = 27 Checks reasoning consistency

Agent 5: Step 1: Green = 5 Step 2: Blue = 10 Step 3: Pink = 9 Step 4: Yellow = 3 Total = 27 All calculations verified

Final Answer: 27

AgentArk:

The problem states that Cindy had 5 green eggs.

Step 1: She had twice as many blue eggs as green eggs, so she had 2 * 5 = 10 blue eggs.

Step 2: She had one fewer pink eggs than blue eggs, so she had 10 - 1 = 9 pink eggs.

Step 3: She had one-third as many yellow eggs as pink eggs, so she had 9 / 3 = 3 yellow eggs.

Step 4: The total number of eggs is 5 + 10 + 9 + 3 = 27 eggs.

Final Answer: 27

Reasoning: The problem states that Cindy had 5 green eggs. She had twice as many blue eggs as green eggs, so she had 2 * 5 = 10 blue eggs. She had one fewer pink eggs than blue eggs, so she had 10 - 1 = 9 pink eggs. She had one-third as many yellow eggs as pink eggs, so she had 9 / 3 = 3 yellow eggs. The total number of eggs is 5 + 10 + 9 + 3 = 27 eggs. The reasoning provided in the final answer is correct and matches the problem statement.

## E.1 Perplexity Calculation

**Definition.** Given a reasoning sequence $x = (x_1, \ldots, x_T)$, the perplexity (PPL) of a language model parameterized by $\theta$ is defined as:

$$\text{PPL}(x) = \exp\left(-\frac{1}{T}\sum_{t=1}^{T}\log p_\theta(x_t \mid x_{<t})\right), \tag{10}$$

where $p_\theta(x_t \mid x_{<t})$ denotes the conditional probability assigned by the model to the next token $x_t$ given all previous tokens.

**Reasoning-Token Perplexity.** Unlike standard perplexity evaluation over full responses, we restrict the computation to *reasoning tokens* only. Specifically, we identify reasoning spans corresponding to intermediate reasoning steps (e.g., chain-of-thought segments) and exclude prompt tokens, question descriptions, and final answer tokens. Let $\mathcal{R} \subseteq \{1, \ldots, T\}$ denote the index set of reasoning tokens. The reasoning perplexity is then computed as:

$$\text{PPL}_{\text{reason}}(x) = \exp\left(-\frac{1}{|\mathcal{R}|}\sum_{t\in\mathcal{R}}\log p_\theta(x_t \mid x_{<t})\right). \tag{11}$$

**Evaluation Protocol.** We evaluate perplexity on held-out GSM8K samples that are not used during distillation. For each example, the gold reasoning trajectory is tokenized using the same tokenizer as the evaluated model. The model is run in teacher-forcing mode to obtain token-level log-likelihoods. Perplexity is computed per sample and then averaged across all evaluation samples.

**Model Setup.** We compute perplexity for distilled student models and the single-agent baseline using the same teacher forcing procedure. The Qwen3-32B model serves as the teacher, while Qwen3-0.6B is used as the student architecture. All models are evaluated with identical prompts and decoding-free forward passes to ensure fair comparison.

**Interpretation.** Lower reasoning perplexity indicates that the model assigns higher likelihood to coherent and structured reasoning steps, reflecting improved internal consistency and predictability of the reasoning process. This aligns with our objective of distilling multi-agent reasoning dynamics into a single model.

# F Combined Distillation Strategies with Two-Step Training

In addition to evaluating individual distillation strategies, including reasoning-enhanced supervised fine-tuning (RSFT), reasoning data augmentation (DA), and process-aware distillation (PAD), we further investigate whether these methods can be effectively combined in a sequential manner.

## F.1 Two-Step Training Protocol

We adopt a two-step training scheme that stacks DA on top of an existing distillation method. Specifically, we consider **RSFT+DA** and **PAD+DA**, where: (i) in the first stage, the student model is trained using RSFT or PAD following the same setup as the main experiments; and (ii) in the second stage, the resulting model is further fine-tuned with reasoning data augmentation to enhance reasoning diversity.

All experiments distill from a **Qwen3-32B** source model to a **Qwen3-1.7B** student model. In the second stage, we separately use reasoning data generated by the 32B model on **GSM8K** and **MATH**, denoted as *32B_gsm8k* and *32B_math*. The *base* setting corresponds to using DA data aligned with the same task as the first-stage training.

Table 7 summarizes the results of mixed distillation strategies on GSM8K and MedMCQA.

Table 7: Performance of combined distillation strategies with two-step training.

| Method | Train / Test | Base | gsm8k | math |
|---|---|---|---|---|
| RSFT+DA | GSM8K | 68.61 | 70.27 | 69.89 |
| | MedMCQA | 42.67 | 43.92 | 44.30 |
| PAD+DA | GSM8K | 70.03 | 70.69 | 70.14 |
| | MedMCQA | 42.72 | 43.68 | 43.49 |

Overall, stacking DA on top of RSFT or PAD leads to consistent but modest gains across benchmarks. While the improvements are incremental, the results suggest that our methods are mutually compatible, as reasoning data augmentation can be applied on top of different distillation strategies without degrading performance.

# G Distillation Results

## G.1 Supervised Fine-tuning

We evaluate the performance of distillation through standard SFT trained solely with ground-truth answer supervision, following common practice. Specifically, models are fine-tuned using only input–output pairs from the ground-truth data, without any additional reasoning annotations, auxiliary losses, or trajectory-level supervision. After training, we assess the performance both in-distribution performance and out-of-distribution generation under distribution shifts that require changes in reasoning.

The results is shown in Table 8. A consistent pattern emerges across model families and scales. While SFT occasionally yields moderate gains on MedMCQA, it fails to produce consistent or reliable improvements on GSM8K, and in many cases leads to performance degradation. Specifically, MedMCQA exhibits small but repeatable improvements when models are fine-tuned on medically oriented or structurally similar datasets (e.g., +6.8 for Gemma-7B, +3.4 for Qwen3-8B). This suggests that answer-only supervision can be beneficial when the target task shares surface-level structure or domain overlap with the training data, allowing models to exploit task-specific correlations. In contrast, GSM8K performance remains largely unimproved or even deteriorates under SFT across nearly all settings. Even when trained directly on GSM8K, models frequently underperform their base counterparts. This behavior indicates that SFT struggles to induce transferable reasoning strategies required for multi-step mathematical problem solving, and instead encourages overfitting to shallow input–output mappings that do not generalize beyond the supervised distribution. Overall, these results highlight a key limitation of standard SFT: while answer-only supervision may yield localized gains on domain-specific benchmarks, it fails to support robust cross-task or reasoning-intensive generalization. This empirical evidence motivates the need for process-level supervision that exposes models to intermediate reasoning dynamics rather than only final outcomes.

Table 8: SFT performance

| Test \ Train | | base | gsm8k | medmcqa | metamathqa | math |
|---|---|---|---|---|---|---|
| gsm8k | llama3-8b | 75.28 | 63.61 | 73.31 | 65.78 | 75.36 |
| medmcqa | | 56.9 | 60.79 | 60.44 | 60.46 | 60.34 |
| gsm8k | qwen3-8b | 88.17 | 81.35 | 88.17 | 87.41 | 87.87 |
| medmcqa | | 59.65 | 60.05 | 63.02 | 59.98 | 59.67 |
| gsm8k | gemma-7b | 51.1 | 56.56 | 48.52 | 36.39 | 50.04 |
| medmcqa | | 49.29 | 48.98 | 57.06 | 48.15 | 48.43 |
| gsm8k | qwen3-1.7b | 69.14 | 59.59 | 70.66 | 66.34 | 67.32 |
| medmcqa | | 42.34 | 43.72 | 49.53 | 43.68 | 42.46 |
| gsm8k | qwen3-0.6b | 41.93 | 37.91 | 44.35 | 21.38 | 31.31 |
| medmcqa | | 32.2 | 32.9 | 38.9 | 34.09 | 32.27 |

## G.2   Reasoning Based Distillation Results

Results for our three proposed methods are listed in Tables 9 to 13.

Table 9: Distillation results on Qwen3-8B. Rows denote training datasets and columns denote test benchmarks (see diagonal header). Results compare RSFT, DA, and PAD across different source models. The original Qwen3-8B scores 88.17 on GSM8K and 59.65 on MedMCQA.

| source model / test | train | gsm8k | medmcqa | metamathqa | math |
|---|---|---|---|---|---|
| | | RSFT | | | |
| | gsm8k | 87.11 | 86.88 | 87.72 | 87.19 |
| | medmcqa | 58.81 | 58.88 | 59.6 | 59.77 |
| | | DA | | | |
| | gsm8k | 86.43 | 86.35 | 88.1 | 86.41 |
| qwen3-8b | medmcqa | 58.33 | 58.83 | 57.49 | 58.35 |
| | | PAD | | | |
| | gsm8k | 88.42 | 88.36 | 88.48 | 88.41 |
| | medmcqa | 59.71 | 59.81 | 59.71 | 59.73 |
| | | RSFT | | | |
| | gsm8k | 86.73 | 87.04 | 85.6 | 85.97 |
| | medmcqa | 60.03 | 59.74 | 57.95 | 59.96 |
| | | DA | | | |
| | gsm8k | 87.11 | 87.49 | 89.57 | 86.2 |
| qwen3-32b | medmcqa | 59.48 | 58.69 | 58.33 | 59.86 |
| | | PAD | | | |
| | gsm8k | 89.05 | 89.02 | 89.15 | 88.7 |
| | medmcqa | 60.04 | 63.12 | 61.53 | 61.21 |
| | | RSFT | | | |
| | gsm8k | 86.96 | 86.35 | 87.72 | 87.64 |
| | medmcqa | 59.53 | 58.43 | 59.77 | 59.74 |
| | | DA | | | |
| | gsm8k | 87.79 | 86.28 | 86.2 | 86.05 |
| gemma3-27b-it | medmcqa | 59.43 | 58.5 | 57.18 | 59.72 |
| | | PAD | | | |
| | gsm8k | 88.48 | 88.48 | 88.48 | 88.4 |
| | medmcqa | 59.96 | 59.84 | 59.86 | 59.79 |

Table 10: Distillation results on LLaMA3-8B-Instruct. Rows denote training datasets and columns denote test benchmarks (see diagonal header). Results compare RSFT, DA, and PAD across different source models. The original LLaMA3-8B-Instruct scores 75.28 on GSM8K and 56.9 on MedMCQA.

| source model | test \ train | gsm8k | medmcqa | metamathqa | math |
|---|---|---|---|---|---|
| | | RSFT | | | |
| | gsm8k | 69.22 | 75.36 | 74.98 | 76.04 |
| | medmcqa | 60.08 | 60.39 | 59.86 | 60.87 |
| | | DA | | | |
| qwen3-8b | gsm8k | 70.89 | 75.06 | 73.24 | 75.82 |
| | medmcqa | 56.44 | 58.36 | 58.59 | 56.83 |
| | | PAD | | | |
| | gsm8k | 76.42 | 75.28 | 75.55 | 75.72 |
| | medmcqa | 57.11 | 57.16 | 57.02 | 57.1 |
| | | RSFT | | | |
| | gsm8k | 75.77 | 73.54 | 76.15 | 77.18 |
| | medmcqa | 61.15 | 57.64 | 58.95 | 60.6 |
| | | DA | | | |
| qwen3-32b | gsm8k | 77.8 | 74.39 | 77.53 | 75.88 |
| | medmcqa | 56.47 | 58.12 | 58.83 | 56.77 |
| | | PAD | | | |
| | gsm8k | 76.88 | 75.97 | 75.89 | 76.95 |
| | medmcqa | 60.82 | 60.75 | 60.79 | 60.96 |
| | | RSFT | | | |
| | gsm8k | 75.26 | 72.02 | 76.55 | 74.98 |
| | medmcqa | 59.77 | 58.76 | 59.5 | 60.94 |
| | | DA | | | |
| gemma3-27b-it | gsm8k | 77.89 | 72.33 | 76.37 | 75.59 |
| | medmcqa | 56.44 | 55.3 | 58.93 | 56.92 |
| | | PAD | | | |
| | gsm8k | 76.35 | 75.97 | 77.02 | 76.27 |
| | medmcqa | 60.82 | 60.94 | 60.73 | 60.55 |

Table 11: Distillation results on Gemma-7B. Rows denote training datasets and columns denote test benchmarks (see diagonal header). Results compare RSFT, DA, and PAD across different source models. The original Gemma-7B scores 51.1 on GSM8K and 49.29 on MedMCQA.

| source model | test \ train | gsm8k | medmcqa | metamathqa | math |
|---|---|---|---|---|---|
| | | RSFT | | | |
| | gsm8k | 64.52 | 60.35 | 75.28 | 62.09 |
| | medmcqa | 47.17 | 51.37 | 45.04 | 48.51 |
| | | DA | | | |
| qwen3-8b | gsm8k | 65.73 | 62.02 | 74.75 | 62.62 |
| | medmcqa | 46.38 | 50.97 | 42.84 | 48.84 |
| | | PAD | | | |
| | gsm8k | 52.74 | 51.69 | 52.34 | 52.55 |
| | medmcqa | 49.32 | 49.66 | 49.37 | 49.42 |
| | | RSFT | | | |
| | gsm8k | 63.23 | 54.66 | 67.25 | 57.62 |
| | medmcqa | 49.01 | 49.52 | 49.76 | 48.74 |
| | | DA | | | |
| qwen3-32b | gsm8k | 66.26 | 58.23 | 73.78 | 63.46 |
| | medmcqa | 48.31 | 48.06 | 48.74 | 49.32 |
| | | PAD | | | |
| | gsm8k | 63.07 | 63.37 | 63.6 | 62.46 |
| | medmcqa | 49.7 | 50.77 | 49.61 | 50.77 |
| | | RSFT | | | |
| | gsm8k | 67.17 | 60.65 | 73.46 | 63.61 |
| | medmcqa | 48.31 | 42.03 | 45.66 | 48.46 |
| | | DA | | | |
| gemma3-27b-it | gsm8k | 68.16 | 60.27 | 71.34 | 59.59 |
| | medmcqa | 47.36 | 38.35 | 42.34 | 48.7 |
| | | PAD | | | |
| | gsm8k | 53.37 | 53.15 | 53.24 | 52.62 |
| | medmcqa | 50.62 | 50.08 | 50.77 | 49.41 |

Table 12: Distillation results on Qwen3-1.7B. Rows denote training datasets and columns denote test benchmarks (see diagonal header). Results compare RSFT, DA, and PAD across different source models. The original Qwen3-1.7B scores 69.14 on GSM8K and 42.34 on MedMCQA.

| source model | test ╲ train | gsm8k | medmcqa | metamathqa | math |
|---|---|---|---|---|---|
| | | RSFT | | | |
| | gsm8k | 67.1 | 67.55 | 69.67 | 69.98 |
| | medmcqa | 43.03 | 42.72 | 43.22 | 42.27 |
| | | DA | | | |
| qwen3-8b | gsm8k | 68.99 | 69.07 | 69.29 | 69.23 |
| | medmcqa | 48.49 | 45.18 | 42.65 | 42.41 |
| | | PAD | | | |
| | gsm8k | 69.82 | 69.41 | 69.78 | 69.76 |
| | medmcqa | 42.36 | 42.42 | 42.41 | 42.44 |
| | | RSFT | | | |
| | gsm8k | 69.7 | 65.43 | 69.85 | 68.61 |
| | medmcqa | 42.77 | 43.74 | 44.32 | 42.67 |
| | | DA | | | |
| qwen3-32b | gsm8k | 70.78 | 68.84 | 73.89 | 70.03 |
| | medmcqa | 43.68 | 42.46 | 43.03 | 42.72 |
| | | PAD | | | |
| | gsm8k | 70.36 | 70.05 | 70.28 | 69.83 |
| | medmcqa | 44.41 | 44.46 | 42.43 | 42.46 |
| | | RSFT | | | |
| | gsm8k | 61.56 | 60.8 | 69.07 | 69.14 |
| | medmcqa | 42.65 | 39.42 | 41.84 | 42.31 |
| | | DA | | | |
| gemma3-27b-it | gsm8k | 71.11 | 65.06 | 69.94 | 64.52 |
| | medmcqa | 43.99 | 40.69 | 42.58 | 42.34 |
| | | PAD | | | |
| | gsm8k | 69.67 | 70.2 | 69.98 | 69.9 |
| | medmcqa | 42.36 | 42.53 | 42.46 | 42.55 |

Table 13: Distillation results on Qwen3-0.6B. Rows denote training datasets and columns denote test benchmarks (see diagonal header). Results compare RSFT, DA, and PAD across different source models. The original Qwen3-0.6B scores 41.93 on GSM8K and 32.2 on MedMCQA.

| source model | test \ train | gsm8k | medmcqa | metamathqa | math |
|---|---|---|---|---|---|
| **qwen3-8b** | | RSFT | | | |
| | gsm8k | 44.05 | 41.89 | 45.03 | 41.39 |
| | medmcqa | 31.08 | 32.9 | 28.19 | 31.7 |
| | | DA | | | |
| | gsm8k | 43.67 | 42.23 | 51.78 | 41.17 |
| | medmcqa | 28.76 | 30.29 | 30.27 | 30.72 |
| | | PAD | | | |
| | gsm8k | 42.57 | 42.25 | 42.51 | 42.53 |
| | medmcqa | 32.43 | 32.33 | 32.38 | 32.27 |
| **qwen3-32b** | | RSFT | | | |
| | gsm8k | 39.04 | 38.51 | 48.14 | 38.74 |
| | medmcqa | 32.49 | 32.85 | 34.4 | 32.66 |
| | | DA | | | |
| | gsm8k | 48.67 | 40.94 | 53.37 | 41.55 |
| | medmcqa | 34.62 | 32.06 | 34.38 | 34.54 |
| | | PAD | | | |
| | gsm8k | 42.84 | 42.46 | 42.53 | 42.68 |
| | medmcqa | 33.09 | 34.36 | 34.54 | 32.58 |
| **gemma3-27b-it** | | RSFT | | | |
| | gsm8k | 40.33 | 32.37 | 49.73 | 40.79 |
| | medmcqa | 33.28 | 30.5 | 34.4 | 32.01 |
| | | DA | | | |
| | gsm8k | 41.02 | 42.46 | 49.2 | 40.94 |
| | medmcqa | 32.27 | 32.94 | 35.43 | 33.23 |
| | | PAD | | | |
| | gsm8k | 44.61 | 43.52 | 43.72 | 42.17 |
| | medmcqa | 33.01 | 34.51 | 33.27 | 32.92 |