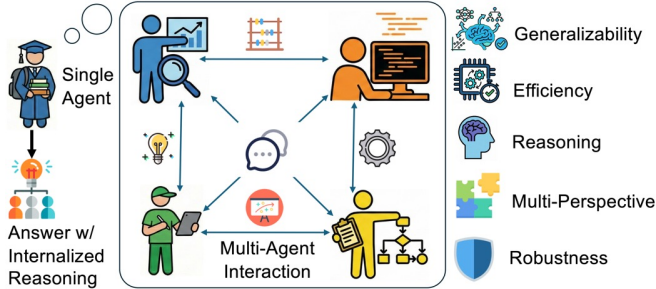




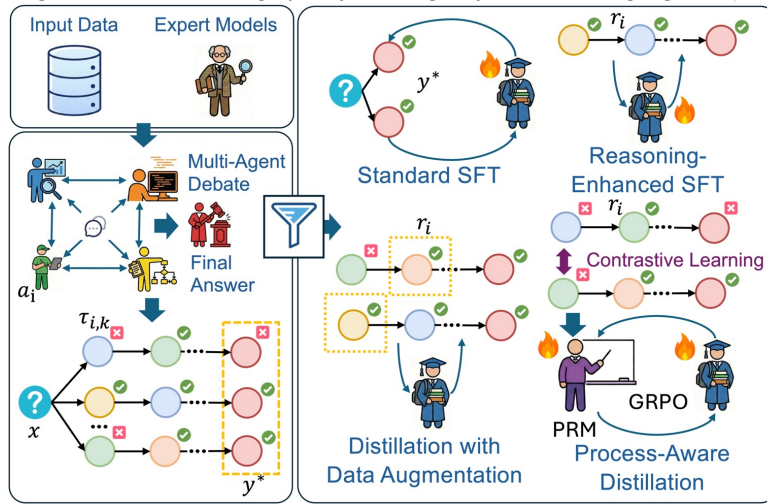
Multi-agent Systems (MASs)

Agents interact via debate, critique, and consensus, but ...

- Multi-role and multi-turn interactions can be expensive and amplify errors
 - Biases & hallucinations can spread across agents, reducing robustness and safety
- Can a single model internalize multi-agent reasoning without the inference-time cost and vulnerabilities?



AgentArk distills the reasoning capability of multi-agent systems into one single agent.



Reasoning-enhanced Supervised Fine-Tuning (RSFT).

Use both the final answers and reasoning traces as supervision.

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

$$\mathcal{L}_{\text{res}} = \sum_{t=1}^{|r|} \log p_{\theta}(r_t | r_{<t}, x) \quad \mathcal{L}_{\text{ans}} = \log p_{\theta}(y^* | r, x)$$

Data Augmentation via Diverse Extraction (DA)

Filter the agents A into a subset of successful contributors $A_{\text{correct}} \subseteq A$.

Utilize a high-capacity teacher LLM as a "distiller" to parse the raw debate logs of A_{correct}

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

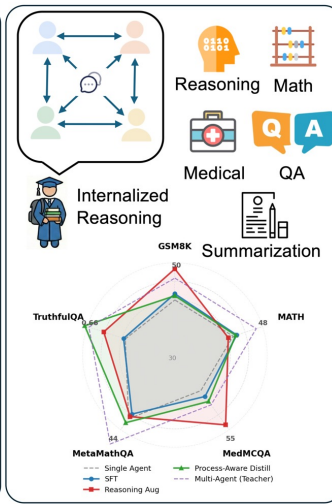
Process-Aware Distillation using PRM and GRPO.

Given an input $x \sim \mathcal{D}$, sample a group of G reasoning outputs o_1, \dots, o_G from a fixed behavior policy π_{old}

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

Key Findings

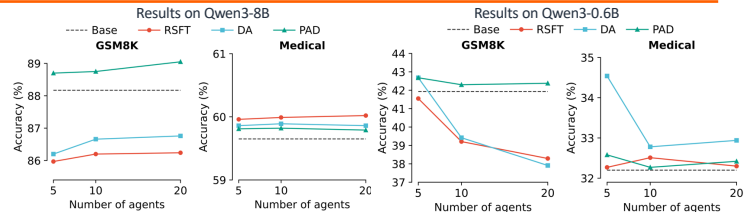
- Single-agent gains**: AgentArk gives a single model multi-agent reasoning ability; all three methods help, and combining them helps more.
- PRM matters most**: PRM capacity matters more than student size, while student capacity limits the gains from larger teacher ensembles.
- Quality > quantity**: More trajectories alone do not help; PAD's high-signal process supervision yields more stable gains.
- Better reasoning, not just accuracy**: PAD improves step decomposition, self-checking, and error correction beyond RSFT and DA.
- Stronger transfer**: AgentArk improves generalization and robustness on unseen and robustness benchmarks.



- Multi-Agent Debate**: generates diverse reasoning trajectories
- Knowledge Extraction**: filters high-quality corrective traces
- Distillation**: uses RSFT, DA, and PAD (PRM + GRPO)
- Student Model**: learns efficient, generalizable reasoning

Promising results across 3 model families

- Qwen3-8B / 1.7B / 0.6B
- Gemma3-27B / 7B
- Llama3-8B-Instruct



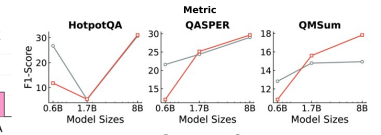
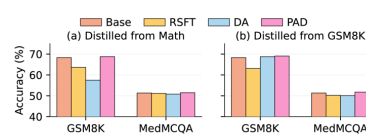
Effects of MAS scale (5, 10, 20) on distillation performance

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

Llama3-8B on OOD datasets

Dataset	Qwen3-0.6B	Qwen3-8B	Gemma3-27B	Llama3-8B	Improvement (Δ)
HotpotQA	29.46 (+0.42)	17.17 (+0.24)	29.39 (+0.41)	86.85 (+0.13)	29.46 (+0.44)
QASPER	22.44 (+0.41)	10.39 (+0.42)	19.69 (+0.39)	81.81 (+1.28)	22.60 (+0.59)
QMSum	25.72 (+0.26)	6.10 (+0.69)	16.53 (+0.75)	84.78 (+0.58)	14.92 (+1.87)

Reasoning Quality



Multimodal distillation (GSM8K & MedMCQA)

AgentArk vs. base models (OOD datasets)