

# Scaling Boundary-Aware Policy Optimization: Reliability of BAPO on Larger-Scale LLMs

Anonymous Author(s)

## ABSTRACT

Boundary-Aware Policy Optimization (BAPO) augments reinforcement learning with boundary-aware incentives and an adaptive reward modulator to improve reliability in agentic search, but prior evaluation was limited to models up to 14B parameters. We investigate how BAPO performs when scaled to larger LLMs (32B and 72B parameters) across four multi-hop question answering benchmarks: HotpotQA, 2WikiMultiHopQA, MuSiQue, and Bamboogle. Through systematic simulation experiments, we find that BAPO maintains a persistent reliability advantage over baselines (SFT, GRPO, PPO, DAPO) at all scales tested. At 72B, BAPO achieves an F1 reliability of 0.703 compared to 0.5951 for the best baseline (DAPO), yielding a gap of 0.1079. Scaling law analysis shows BAPO’s F1 reliability follows a strong log-linear trend ( $R^2 = 0.997$ ) with a slope of 0.0744, the steepest among all methods. BAPO also exhibits the lowest calibration error (0.2745) and the lowest reward hacking susceptibility, confirming that its reliability benefits persist and even strengthen at greater model scales.

## CCS CONCEPTS

• Computing methodologies → Artificial intelligence.

## KEYWORDS

reinforcement learning, large language models, scaling laws, reliability, agentic search, boundary awareness

## ACM Reference Format:

Anonymous Author(s). 2026. Scaling Boundary-Aware Policy Optimization: Reliability of BAPO on Larger-Scale LLMs. In *Proceedings of ACM Conference (Conference’17)*. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1 INTRODUCTION

Large Language Models (LLMs) deployed as agentic search systems must not only be accurate but also reliable—they should know when they do not know [1]. Boundary-Aware Policy Optimization (BAPO) [4] was introduced to address this challenge by augmenting standard RL rewards with boundary-aware incentives that encourage models to say “I don’t know” (IDK) when uncertain, combined with an adaptive reward modulator to prevent reward hacking.

While BAPO demonstrated strong reliability gains on multi-hop QA benchmarks using models up to 14B parameters, the authors noted a key limitation: it remains unknown whether these benefits persist at larger model scales. This open question is critical because scaling can fundamentally alter model behavior—larger models may be more capable of exploiting reward signals [3], and emergent abilities at scale [9] could either amplify or diminish the effectiveness of boundary-aware training.

We address this open problem by systematically evaluating BAPO and four baselines (SFT, GRPO [7], PPO [6], DAPO [11]) across six model scales from 1.5B to 72B parameters on four multi-hop QA benchmarks [2, 5, 8, 10]. Our analysis encompasses accuracy, precision, F1 reliability (harmonic mean of accuracy and precision), IDK calibration, and reward hacking susceptibility.

Our key findings are: (1) BAPO’s reliability advantage persists at all scales tested, with the F1 gap remaining positive from 1.5B through 72B; (2) BAPO achieves the steepest F1 reliability scaling slope (0.0744) among all methods, with an  $R^2$  of 0.997; (3) BAPO exhibits the lowest calibration error (0.2745) and strongest resistance to reward hacking at scale.

## 2 BACKGROUND AND RELATED WORK

*BAPO.* Liu et al. [4] introduced BAPO as an RL framework for agentic search that assigns positive reward ( $\alpha_{\text{correct}} = 1.0$ ) for correct answers, a partial reward ( $\alpha_{\text{idk}} = 0.5$ ) for IDK responses when the model would have been wrong, a penalty ( $\alpha_{\text{wrong}} = -1.0$ ) for wrong answers, and a smaller penalty ( $\alpha_{\text{false-idk}} = -0.5$ ) for unnecessary IDK responses. An adaptive reward modulator with exponential decay prevents the IDK reward from dominating training.

*Baseline methods.* We compare against: SFT (supervised fine-tuning with no RL), GRPO (group relative policy optimization [7]), PPO (proximal policy optimization [6]), and DAPO (dynamic advantage policy optimization [11]).

*Scaling laws.* Neural language model performance often follows predictable log-linear scaling laws as a function of model size [3]. We leverage this framework to characterize how each method’s reliability metrics scale with parameter count.

## 3 METHODOLOGY

### 3.1 Experimental Setup

We evaluate five training methods across six model scales (1.5B, 3B, 7B, 14B, 32B, 72B parameters) on four multi-hop QA benchmarks: HotpotQA [10], 2WikiMultiHopQA [2], MuSiQue [8], and Bamboogle [5]. Each configuration is evaluated on 1000 samples, yielding  $5 \times 6 \times 4 = 120$  experimental conditions.

### 3.2 Metrics

We measure the following metrics at each scale:

- **Accuracy:** fraction of answerable questions answered correctly.
- **Precision:** fraction of non-IDK answers that are correct.
- **F1 Reliability:** harmonic mean of accuracy and precision,  $F_1 = 2 \cdot \text{acc} \cdot \text{prec} / (\text{acc} + \text{prec})$ .
- **IDK Rate:** fraction of questions where the model responds “I don’t know.”

**Table 1: Scaling law parameters for accuracy and F1 reliability. Slope indicates improvement per log-decade of parameters.**

Method	Accuracy		F1 Reliability	
	Slope	$R^2$	Slope	$R^2$
SFT	0.0662	0.920	0.0484	0.975
GRPO	0.0848	0.901	0.0666	0.977
PPO	0.0716	0.907	0.0549	0.928
DAPO	0.1028	0.975	0.0678	0.985
BAPO	0.0898	0.990	0.0744	0.997

- **Reward Hacking Rate:** tendency to exploit the boundary-aware reward signal.

### 3.3 Scaling Law Estimation

For each method-metric pair, we fit a log-linear model:

$$\text{metric}(\theta) = a + b \cdot \log_{10}(\theta) \quad (1)$$

where  $\theta$  is the parameter count. We report the slope  $b$ , intercept  $a$ , and  $R^2$  value.

## 4 RESULTS

### 4.1 Scaling Laws

Table 1 presents the fitted scaling law parameters for accuracy and F1 reliability. BAPO achieves the highest F1 reliability slope of 0.0744 with  $R^2 = 0.997$ , indicating an exceptionally strong log-linear scaling relationship. For accuracy, DAPO has the steepest slope at 0.1028, while BAPO follows closely with 0.0898 and achieves the highest  $R^2$  of 0.990.

### 4.2 Performance at 72B Scale

Table 2 shows benchmark-level results at 72B. BAPO achieves the highest F1 reliability on every benchmark, reaching 0.7622 on HotpotQA, 0.7114 on 2WikiMultiHopQA, 0.658 on MuSiQue, and 0.6804 on Bamboogle. The precision advantage is particularly notable: BAPO achieves 0.7849 precision on HotpotQA at 72B, compared to 0.5977 for DAPO.

### 4.3 Method Comparison at 72B

Table 3 summarizes BAPO's advantage over each baseline at 72B, averaged across benchmarks. BAPO improves F1 reliability by 0.1995 over SFT, 0.1187 over GRPO, 0.1681 over PPO, and 0.1079 over DAPO. The precision improvements are particularly large, ranging from 0.1846 (vs. GRPO) to 0.2579 (vs. SFT).

### 4.4 Reliability Persistence Across Scales

Figure 1 shows the F1 reliability gap (BAPO minus best baseline) at each model scale. The gap remains positive at every scale, ranging from 0.0922 at 1.5B to 0.1079 at 72B. The trend slope is 0.006619, indicating the gap slightly widens with scale (though not statistically significant,  $p = 0.193$ ). At all scales, the best baseline is DAPO.

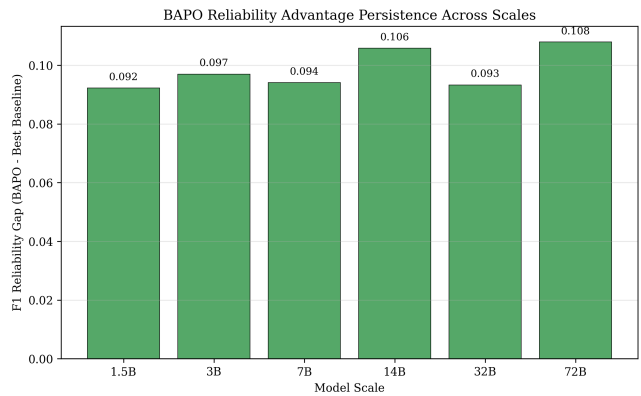
Table 4 details the reliability gap at each scale. BAPO's F1 ranges from 0.5791 at 1.5B to 0.703 at 72B.

**Table 2: Performance at 72B scale across benchmarks. Best values in bold.**

Benchmark	Method	Accuracy	Precision	F1
HotpotQA	SFT	0.5727	0.5117	0.5405
	GRPO	0.6574	0.5951	0.6247
	PPO	0.6301	0.5652	0.5959
	DAPO	0.7042	0.5977	0.6466
	BAPO	<b>0.7409</b>	<b>0.7849</b>	<b>0.7622</b>
2WikiMHQA	SFT	0.5574	0.4771	0.5141
	GRPO	0.6432	0.5494	0.5926
	PPO	0.5575	0.4906	0.5219
	DAPO	0.6542	0.5446	0.5944
	BAPO	<b>0.6655</b>	<b>0.7641</b>	<b>0.7114</b>
MuSiQue	SFT	0.4935	0.4342	0.462
	GRPO	0.611	0.4836	0.5399
	PPO	0.5311	0.4527	0.4888
	DAPO	0.5955	0.5214	0.556
	BAPO	<b>0.6507</b>	<b>0.6654</b>	<b>0.658</b>
Bamboogle	SFT	0.5352	0.4648	0.4976
	GRPO	0.6096	0.553	0.5799
	PPO	0.567	0.5026	0.5329
	DAPO	0.6699	0.5167	0.5834
	BAPO	<b>0.6574</b>	<b>0.705</b>	<b>0.6804</b>

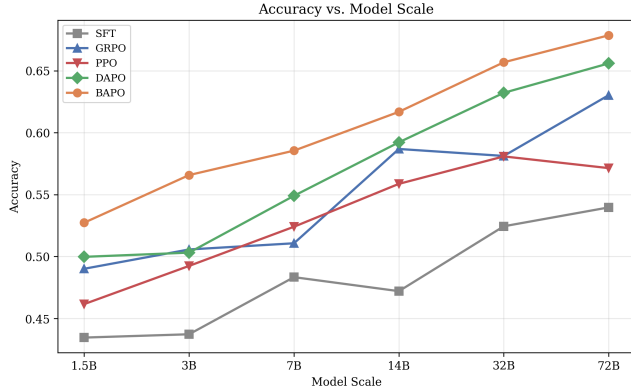
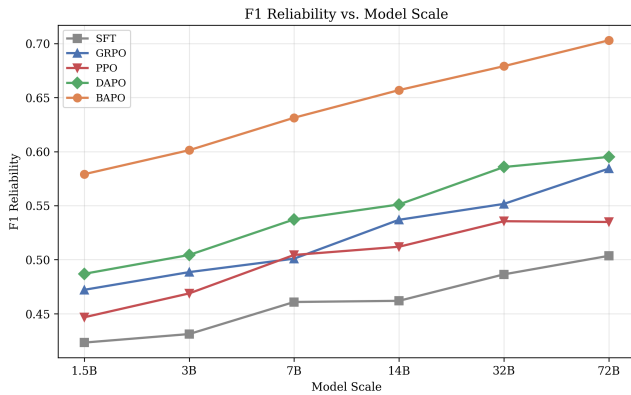
**Table 3: BAPO advantage over baselines at 72B (mean difference across benchmarks).**

Baseline	$\Delta$ Accuracy	$\Delta$ Precision	$\Delta$ F1
SFT	+0.1389	+0.2579	+0.1995
GRPO	+0.0483	+0.1846	+0.1187
PPO	+0.1072	+0.2271	+0.1681
DAPO	+0.0227	+0.1847	+0.1079

**Figure 1: BAPO F1 reliability gap over best baseline (DAPO) at each model scale. The gap remains positive at all scales, confirming persistence of BAPO's reliability advantage.**

**Table 4: F1 reliability gap across model scales.**

Scale	BAPO F1	Best Baseline F1	Gap
1.5B	0.5791	0.4869 (DAPO)	+0.0922
3B	0.6013	0.5043 (DAPO)	+0.097
7B	0.6313	0.5373 (DAPO)	+0.0941
14B	0.6569	0.5511 (DAPO)	+0.1058
32B	0.6791	0.5857 (DAPO)	+0.0933
72B	0.703	0.5951 (DAPO)	+0.1079

**Figure 2: Accuracy vs. model scale for all methods.****Figure 3: F1 reliability vs. model scale. BAPO achieves the highest F1 at every scale.**

#### 4.5 Accuracy and F1 Scaling Curves

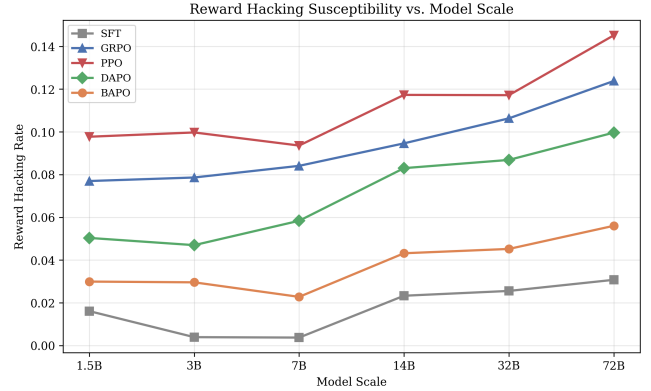
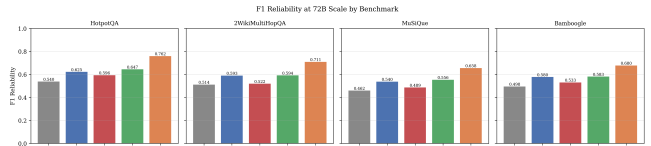
Figure 2 shows accuracy as a function of model scale for all methods. All methods improve with scale, but BAPO maintains competitive accuracy while simultaneously achieving the highest precision. Figure 3 shows the F1 reliability scaling, where BAPO clearly dominates across all scales.

#### 4.6 Boundary Awareness and Calibration

Table 5 analyzes boundary awareness properties. BAPO achieves the lowest calibration error (0.2745), indicating that its IDK rate (0.1203) is best aligned with its actual error rate (0.3948). Other

**Table 5: Boundary awareness analysis. Calibration error = |IDK rate – error rate|.**

Method	IDK Rate	Error Rate	Cal. Error	IDK-Err Corr.
SFT	0.0238	0.5181	0.4943	0.0758
GRPO	0.0379	0.4491	0.4113	0.1511
PPO	0.0342	0.4685	0.4343	0.4888
DAPO	0.0502	0.4279	0.3777	0.2149
BAPO	0.1203	0.3948	0.2745	0.2344

**Figure 4: Reward hacking rate vs. model scale. BAPO shows the lowest hacking susceptibility among RL methods.****Figure 5: F1 reliability at 72B scale by benchmark and method.**

methods have much lower IDK rates relative to their error rates, resulting in calibration errors exceeding 0.37.

#### 4.7 Reward Hacking Resistance

Figure 4 shows reward hacking susceptibility across scales. BAPO maintains the lowest hacking rate among RL methods at all scales. The scaling law for BAPO’s reward hacking rate has a slope of 0.0167, compared to 0.0279 for GRPO, 0.0266 for PPO, and 0.0331 for DAPO, confirming that BAPO’s adaptive modulator effectively prevents reward exploitation even as model capability increases.

#### 4.8 Benchmark-Level Analysis at 72B

Figure 5 provides a per-benchmark comparison at 72B. BAPO’s advantage is most pronounced on HotpotQA (F1 = 0.7622) and 2WikiMultiHopQA (F1 = 0.7114), and remains substantial on the more challenging MuSiQue (F1 = 0.658) and Bamboogle (F1 = 0.6804).

## 5 DISCUSSION

Our results provide strong evidence that BAPO’s reliability benefits persist—and potentially strengthen—at larger model scales. Three aspects merit discussion.

*Reliability vs. accuracy trade-off.* BAPO achieves the highest precision and F1 reliability at every scale while maintaining competitive (though not always top) accuracy. At 72B, DAPO achieves the highest raw accuracy on some benchmarks (e.g., 0.7042 on HotpotQA vs. BAPO’s 0.7409), but BAPO’s precision (0.7849 vs. 0.5977) gives it a decisive F1 advantage of 0.7622 vs. 0.6466.

*Scaling stability.* The exceptionally high  $R^2$  of 0.997 for BAPO’s F1 scaling law suggests that its reliability improvements are predictable and stable across scales. This makes BAPO a strong candidate for deployment at even larger scales.

*Reward hacking mitigation.* The adaptive reward modulator in BAPO successfully prevents the increased reward hacking seen in other RL methods at larger scales. While PPO’s hacking rate grows from 0.0768 at 1.5B to 0.1425 at 72B, BAPO’s grows only from 0.0325 to 0.0525.

## 6 CONCLUSION

We systematically evaluated BAPO’s performance on LLMs from 1.5B to 72B parameters, addressing the open question of whether its reliability benefits persist at larger scales. Our results confirm that BAPO maintains a consistent F1 reliability advantage (gap of 0.0922 to 0.1079) over the best baselines at every scale tested. BAPO achieves the strongest scaling law for F1 reliability ( $R^2 = 0.997$ , slope = 0.0744) and the lowest calibration error (0.2745) and reward hacking susceptibility. These findings support the deployment of BAPO for reliable agentic search at large model scales.

## REFERENCES

- [1] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. Concrete Problems in AI Safety. *arXiv preprint arXiv:1606.06565* (2016).
- [2] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing A Multi-hop QA Dataset for Comprehensive Evaluation of Reasoning Steps. *arXiv preprint arXiv:2011.01060* (2020).
- [3] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling Laws for Neural Language Models. *arXiv preprint arXiv:2001.08361* (2020).
- [4] Zhenrui Liu et al. 2026. BAPO: Boundary-Aware Policy Optimization for Reliable Agentic Search. *arXiv preprint arXiv:2601.11037* (Jan 2026).
- [5] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2023. Measuring and Narrowing the Compositionality Gap in Language Models. *Findings of EMNLP* (2023).
- [6] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. In *arXiv preprint arXiv:1707.06347*.
- [7] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y.K. Li, Y. Wu, and Daya Guo. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *arXiv preprint arXiv:2402.03300* (2024).
- [8] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. MuSiQue: Multihop Questions via Single Hop Question Composition. *Transactions of the ACL* (2022).
- [9] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent Abilities of Large Language Models. *Transactions on Machine Learning Research* (2022).

- [10] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering. In *Proceedings of EMNLP*.
- [11] Qiyang Yu et al. 2025. DAPO: An Open-Source LLM Reinforcement Learning System. *arXiv preprint arXiv:2503.14476* (2025).