# BARREL: Boundary-Aware Reasoning for Factual and Reliable LRMs

# **Anonymous Author(s)**

Affiliation Address email

#### Abstract

Recent advances in Large Reasoning Models (LRMs) have shown impressive capabilities in mathematical and logical reasoning. However, current LRMs rarely admit ignorance or respond with "I don't know". Instead, they often produce incorrect answers while showing undue confidence, raising concerns about their factual reliability. In this work, we identify two pathological reasoning patterns characterized by overthinking that contribute to the overconfident and incorrect answers: *last-minute guessing* and *second-thought spiraling*. To address these issues, we propose BARREL—a novel framework that promotes concise and boundary-aware factual reasoning. Our experiments show that BARREL-training increases the reliability of DeepSeek-R1-Distill-Llama-8B from 39.33% to 61.48%, while still achieving accuracy comparable to models finetuned on reasoning data generated by R1. These results demonstrate that our pilot study is inspiring to build more reliable and factual *System 2* LRMs.

# 4 1 Introduction

5

6

8

9

10

11

12

13

20

21

22

23

24

25

26

27

28

29

30

31

32

33

35

Recent advances in Large Reasoning Models (LRMs) (Jaech et al., 2024; Guo et al., 2025; Team, 2025) have shown impressive performance in specialized reasoning tasks, especially in mathematics and logic. However, these gains have not led to corresponding improvements in reliability. On the contrary, faithfulness hallucination rates are rising (Hughes & Bae, 2023), and helpfulness on factual tasks is declining (Zhao et al., 2025), raising concerns about the reliability of these reasoning models.

In this work, we focus primarily on the factual reliability of LRMs, which is a crucial requirement for many real-world tasks. Factuality of language models involves two aspects: knowing (whether the model holds relevant knowledge) (Huang et al., 2025; Ji et al., 2023) and telling (conveying the correct factual information) (Gekhman et al., 2024; Mallen et al., 2022). As shown in Figure 1, we want LRMs to exhibit two aspects of factual reliability: (1) Identify both what it knows and what it does not know: The model should be able to recognize questions it does not know the answer and respond with "I don't know." (Xu et al., 2024; Zhang et al., 2024a) (2) Tell what it

knows: There is a gap between knowing

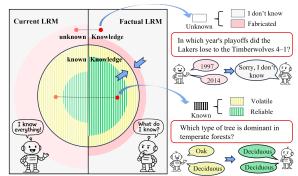


Figure 1: How reliable factual reasoning is expected to improve model performance. **Left:** Current LRMs rarely admit ignorance and often respond inconsistently. **Right:** Reliable LRMs should acknowledge unknowns and express known facts more consistently.

and telling (Saunders et al., 2022), and we want to improve the model's accuracy in expressing the knowledge it has (Zhang et al., 2024b).

However, current LRMs consistently struggle with the two factual reliability goals above. They rarely acknowledge gaps in their knowledge and often fabricate answers instead, even on questions for which they lack sufficient knowledge. Moreover, their responses can be inconsistent—providing incorrect answers in some instances while correctly responding to similar queries elsewhere (Wang et al., 2022), thereby reducing the overall factual accuracy of their responses (Zhao et al., 2025).

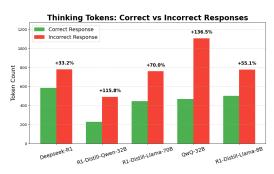


Figure 2: Number of reasoning tokens used by LRMs when producing correct versus incorrect answers. We test on TruthfulQA across different types of reasoning models. Details and results on other datasets are listed at Appendix A.

What underlying reasoning behaviors contribute to these observed issues? As shown in Figure 2, our preliminary experiments reveal a variant of the overthinking pattern discussed by Chen et al. (2024), which we term factual overthinking—where large reasoning models (LRMs) consume more tokens when producing factually incorrect answers than when generating correct ones. And we further identify two prevalent pathological reasoning patterns associated with this overthinking phenomenon. As shown on the left part of Figure 3, one prominent pattern is Last-minute Guessing, in which the model, after extensive but inconclusive reasoning, abruptly commits to an answer in a final burst of speculative output—much like a student hastily filling in an answer just before the end of

an exam. Another frequent pattern is *Second-thought Spiraling*, where the model initially identifies the correct answer but continues to over-analyze, ultimately undermining its own correct conclusion.

To mitigate these pathological reasoning patterns, we propose a novel training framework, **BARREL** (Boundary-Aware Reasoning for Reliable and Factual LRMs). As shown in Figure 3, BARREL trains LRMs to perform concise, deliberative factual reasoning and draw conclusions after it has explored a sufficient number of candidates. Concretely, for known questions, when the model identifies the correct factual answer during reasoning, we want it to maintain confidence in that answer while continuing to explore other possible ones. Once there are no likely candidates, it should halt further reasoning and provide the correct factual answer. For unknown questions, after exploring a sufficient number of plausible candidates, it should terminate reasoning and proactively admit its lack of knowledge with uncertainty-aware refusal (e.g., "Sorry, I don't know"). To implement BARREL, we begin by employing a sampling strategy to probe the model's knowledge boundary. Drawing on the identified pathological patterns, we construct two distinct types of reasoning data, and use Supervised Fine-Tuning (SFT) to instill the corresponding reasoning behaviors in the model. Finally, we adopt Group Relative Policy Optimization (GRPO) using general reliability-based reward (high for correct answers, medium for uncertainty-aware refusal, and low for incorrect answers) without the need for labeling *knownlunknown*, further enhancing the model's ability to generalize in factual reasoning.

Our experiments demonstrate that models trained with BARREL can effectively express uncertainty-aware refusal, and mitigate the two pathological reasoning patterns. This capability significantly improves reliability: BARREL boosts the reliability of DeepSeek-R1-Distill-Llama-8B from 39.33% to 61.48%, while maintaining an accuracy of 40.7%, which is even higher than the accuracy of 38.43% achieved by distillation. Through detailed analysis, we highlight the critical role of medium-level rewards in promoting uncertainty-aware refusal. This result also identifies the root cause of models' inability to admit ignorance to a fundamental gap in current RL paradigms: they do not reward refusal. As a result, models are incentivized to answer every question, regardless of uncertainty.

Our main contributions are as follows:

- We discover the factual overthinking phenomenon and point out two pathological reasoning patterns that lead to factual unreliability of LRMs.
- We explore how LRMs can utilize reasoning to admit ignorance and say "I don't know". Also, we introduce a novel training pipeline to do so.
- We emphasize the importance of medium-level rewards in encouraging uncertainty-aware refusal to build a more reliable LRM using RL training.

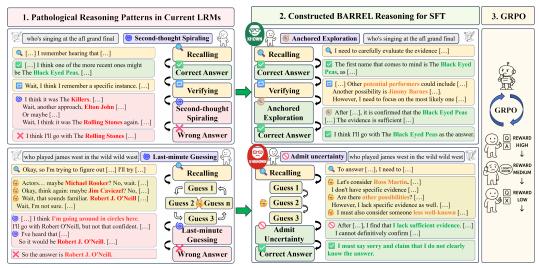


Figure 3: **Left:** The two current reasoning patterns of LRM: Last-minute Guessing, typically associated with unknown knowledge, and Second-thought Spiraling, which occurs despite known knowledge. **Right:** The BARREL pipeline addresses both cases by correcting overthinking tendencies and constructing SFT data accordingly, further enhanced with GRPO.

# **2 Related Work**

**Knowledge Boundary** The knowledge boundary of LLMs refers to the extent of knowledge a model possesses or can reliably recall (Li et al., 2024). Identifying such boundaries is crucial for model safety and reliability and is commonly addressed by confidence calibration (Ren et al., 2023), internal state probing (Ji et al., 2024), uncertainty estimation (Kapoor et al., 2024), and sampling strategy (Gekhman et al., 2024; Xue et al., 2024). To mitigate failures by outbound queries, recent studies have proposed training LLMs to abstain from answering when uncertain, often by generating "I don't know" responses (Cheng et al., 2024) or providing further explanations of the unanswerability (Deng et al., 2024). Our work pushes it further to structured and interpretable rationales about their knowledge boundary.

**Factual Alignment** Factual alignment aims to improve factuality while preserving instruction-following capabilities in LLMs. Recent work incorporates factuality-awareness into both SFT and RL stages to improve factual accuracy (Lin et al., 2024), or finetuning with DPO to enhance its self-evaluation capability (Zhang et al., 2024b). Other approaches include fine-tuning with refusal-aware datasets to encourage appropriate abstention behaviors (Zhang et al., 2024a), and RLKF, which guides models to reject uncertain queries based on external feedback signals (Xu et al., 2024). While previous work has primarily focused on non-reasoning models, our research demonstrates how to correct reasoning pathologies to enhance the factual reliability of LRMs.

# 3 Method

In this section, we introduce the overall framework of BARREL, which comprises three main components: (1) **Knowledge Labeling**, which distinguishes whether a question is known or unknown to the model; (2) **Reasoning Trace Construction for SFT**, which constructs reasoning traces based on the question type and performs SFT to prepare the model to follow this thinking pattern; (3) **GRPO Stage**, which further enhances the model using a rule-based factual reward. We detail each component below.

# 8 3.1 Knowledge Labeling

We first determine whether a question is known to the target model using the sampling strategy proposed in (Gekhman et al., 2024), which is a widely recognized approach (Xue et al., 2024;

Li et al., 2024). Let  $\mathcal{D} = \{(x_i, y_i^*)\}_{i=1}^N$  be a factual QA dataset, where each question  $x_i$  has a ground-truth answer  $y_i^*$ . We generate answers with the target model  $\mathcal{M}$  using K distinct few-shot prompts  $\{\mathcal{P}_j\}_{j=1}^K$  and repeat the sampling procedure L times for every prompt:

$$y_i^{j,k} \sim \mathcal{M}(\cdot \mid \mathcal{P}_j \parallel x_i), \qquad j = 1, \dots, K, \ k = 1, \dots, L.$$
 (1)

After collecting the samples  $\mathcal{Y}_i = \left\{y_i^{j,k}\right\}_{j=1,k=1}^{K,\ L}$ , we consider a question *known* to the model if at least one sampled answer matches the ground-truth answer under evaluator E.

$$l_i = \begin{cases} known, & \text{if } \exists y \in \mathcal{Y}_i \text{ such that } E(y, y_i^*) = 1, \\ unknown, & \text{otherwise.} \end{cases}$$
 (2)

# 3.2 Reasoning Trace Construction for SFT

126

To address the pathological reasoning patterns identified in our analysis—namely, Last-minute Guessing and Second-thought Spiraling—we propose a targeted method for constructing reasoning trajectories. Based on the type of question, we construct two distinct evidence-grounded reasoning traces  $\mathcal{T}(x_i)$  for a question  $x_i$ , aiming to respectively correct these two faulty reasoning patterns. This construction is outlined in Algorithm 1.

To mitigate Second-thought Spiraling in **known** questions, where the gold answer  $y^*$  with strong evidence  $e^*$  is available, it should begin by retrieving and identifying this answer. It then examines alternative candidates  $(y_j, e_j)$  to contrast possibilities. After this anchored exploration, it reaffirms the choice with solid justification and draws a confident conclusion favoring the correct answer.

To address Last-minute Guessing in **unknown** questions, the system adopts a similar exploratory strategy: it recalls background knowledge and searches on plausible answer-evidence pairs  $(y_j, e_j)$  through hypothesizing. However, if it fails to identify a sufficiently supported answer, it explicitly acknowledges the uncertainty and ultimately outputs a cautious, confirmed rejection—demonstrating its ability to explore high-probability paths without overcommitting or hallucinating.

# **Algorithm 1** BARREL reasoning trace $\mathcal{T}(x_i)$ construction

**Input:** Question  $x_i$  and knowledge label  $l_i$ , gold answer with evidence  $(y^*, e^*)$ , alternative candidates with poor evidence  $\{(y_j, e_j)\}_{j=1}^n$ 

```
Output: reasoning trace \mathcal{T}(x_i)
 1: \mathcal{T}(x_i) \leftarrow \langle \rangle
                                                                                                  2: \mathcal{T}(x_i) += \text{RECALL}(x_i)
                                                                                       > record recalled background facts
 3: if l_i = 1(known \ knowledge) then

\mathcal{T}(x_i) += \langle y^*, e^* \rangle 

\mathcal{T}(x_i) += \{ (y_j, e_j) \}_{j=1}^n

                                                                       ▶ Attach gold answer and supporting evidence
 5:

    Attach distractor answer-evidence pairs

          \mathcal{T}(x_i) += \text{CONFIRM}(y^*)
                                                                         ▶ Verify the conclusion with strong evidence
 7: else (unknown knowledge)

8: \mathcal{T}(x_i) += \{(y_j, e_j)\}_{j=1}^n

9: \mathcal{T}(x_i) += Acknowledge Uncertainty()
                                                                           > exploring plausible answer-evidence pairs
                                                                                          > Record uncertainty for guesses
10: end if
11: return \mathcal{T}(x_i)
                                                                               ▶ Return the constructed reasoning trace
```

We construct the reasoning traces by prompting GPT-4 with detailed instructions and BARREL reasoning examples. This approach produces a Long-CoT-style reasoning process that aligns with the expected reasoning patterns. Examples of the constructed reasoning traces for both known and unknown questions are shown in Figure 3, and the detailed prompt used for trace construction is provided in Appendix F.

Then, we use these data to train the model to emulate boundary-aware and deliberative reasoning patterns using SFT. For each question  $x_i$ , we construct full output  $o_i^* = \mathcal{T}(x_i) \| a_i$ , where  $a_i$  is either the gold answer  $y_i^*$  (for known questions) or an uncertainty-aware refusal (e.g., "Sorry, I don't know") for unknown questions. This instills the model with a disciplined reasoning style grounded in traceable evidence and uncertainty-aware conclusions. The training objective minimizes the negative

151 log-likelihood:

$$\mathcal{L}(\theta) = -\sum_{i=1}^{N} \log P_{\theta} \left( o_i^* \mid x_i \right). \tag{3}$$

152 3.3 GRPO-Stage

Rule-Based Reward Design To train the model to generate verifiable and boundary-aware reasoning trajectories and answers, we employ a rule-based reward function. We categorize the model response  $o_i$  into three types, each associated with a distinct reward signal. Given a generated answer  $o_i$  to question  $x_i$ , and ground-truth answer  $y_i^*$ , the reward function  $R(o_i, y_i^*)$  is defined as:

$$R(o_i, y_i^*) = \begin{cases} r_c, & \text{if } E\left(o_i, y_i^*\right) = 1, \\ r_s, & \text{if } o_i \text{ contains a valid rejection phrase,} \\ r_w, & \text{otherwise.} \end{cases}$$
 (4)

This reward function provides general supervision for training the model to optimize not only for correctness but also for calibrated uncertainty, aligning with the goals of boundary-aware reasoning. It comprises three components: a high reward for a correct answer  $(r_c)$ , a medium reward for a truthful rejection  $(r_s)$  and a low reward for an incorrect or hallucinated output  $(r_w)$ . To discourage the generation of unfounded claims, the penalty for an incorrect answer is more severe than the outcome of a truthful rejection, thereby incentivizing the model to acknowledge its knowledge boundaries when uncertain. The reward magnitudes follow the order:

$$r_c > r_s > r_w \,. \tag{5}$$

GRPO Training After SFT, the model has learned the pattern of reasoning to express uncertainty appropriately and is able to maintain confidence when the answer is correct. Building on the above reward design, we further enhance the factual reliability of the reasoning model using Group-wise Reinforcement Policy Optimization (GRPO) (Shao et al., 2024). For each labeled input  $(x_i, y_i^*, l_i)$ , we sample a set of G reasoning-answer trajectories from the current policy  $\pi_{\theta_{\text{old}}}$ :

$$\mathcal{O} = \{o_1, \dots, o_i\} \sim \pi_{\theta_{\text{old}}}(\cdot \mid x_i). \tag{6}$$

Each trajectory  $o_j$  includes a reasoning trace followed by a final answer token. GRPO then updates the model parameters to optimize the following clipped reward-weighted objective:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[x_{i} \sim D, \{o_{j}\}_{j=1}^{G} \sim \pi_{\theta_{old}}(O|q)]$$

$$\frac{1}{G} \sum_{j=1}^{G} \frac{1}{|o_{j}|} \sum_{t=1}^{|o_{i}|} \left\{ \min \left[ \rho_{j,t} \hat{A}_{j,t}, \operatorname{clip} \left( \rho_{j,t}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{j,t} \right] - \beta \mathbb{D}_{KL} \left[ \pi_{\theta} || \pi_{ref} \right] \right\},$$
(7)

where  $\rho_{i,t} = \frac{\pi_{\theta}(o_{j,t}|x_i,o_{j,< t})}{\pi_{\theta_{old}}(o_{j,t}|x_i,o_{j,< t})}$  is the importance weight at step t, and  $\mathbb{D}_{KL}$  denotes the stepwise KL

divergence between the current and old policies. The advantage estimate  $\hat{A}_{j,t}$  is calculated using reward normalization:

$$\hat{A}_{j,t} = \frac{R(o_j, y_i^*) - \bar{R}}{\sigma_r}, \qquad \sigma_r = \sqrt{\frac{1}{G} \sum_{j=1}^G (R(o_j, y_i^*) - \bar{R})^2}.$$
 (8)

# 4 Experiments

175 4.1 Settings

174

Datasets We use separate datasets for training and evaluation. The training set consists of TriviaQA (Joshi et al., 2017), SciQ (Welbl et al., 2017), and NQ-Open (Kwiatkowski et al., 2019), covering general knowledge, scientific reasoning, and web-based QA, respectively. For evaluation, we sample 1,000 questions from the test splits of each dataset, forming a 3,000-question test set.

Models Due to limited computing resource, we primarily utilize DeepSeek-R1-Distill-Llama-8B, DeepSeek-R1-Distill-Qwen-7B and Qwen3-8B to perform our study.

**Baselines** (1) ICL: Vanilla Reasoning models with few-shot prompt designed for factual tasks. (2) In-Context Learning with Refusal Examples (ICL-IDK): Prompting LRMs to claim uncertainty (3) 183 Distill: SFT training using the reasoning path of DeepSeek-R1 1 on the training set. We list the 184 prompt for ICL-IDK and more details in Appendix H. 185

**Evaluation** We evaluate the correctness of model response by prompting the models to box their 186 final answer as follows, and then we utilize string matching to evaluate whether the answer is correct 187 (the model answer appears in any of the candidates). More details are listed in Appendix G. 188

# Inference Prompt for verifiable Answer

Answer the following question based on your knowledge and put your final answer within boxed{}. {question}

**Metrics** A Factual Reliable LRM should provide as much assistance as possible while making as few errors as possible, such that we evaluate factuality on the test set using three metrics: Accuracy (Acc.), Truthfulness (Truth.), and Reliability (Rel.) (Xu et al., 2024). Let  $N_c, N_r, N_w$  denote the number of correct answers, truthful rejections ("Sorry, I don't know"), and incorrect answers, respectively, where  $N=N_c+N_r+N_w$ . The metrics are defined as:

$$\mathrm{Acc.} = \frac{N_c}{N}, \quad \mathrm{Truth.} = \frac{N_c + N_r}{N}, \quad \mathrm{Rel.} = \mathrm{ans.} \cdot \mathrm{Truth.} + (1 - \mathrm{ans.}) \cdot \mathrm{Acc.}, \text{ where } \mathrm{ans.} = 1 - \frac{N_r}{N}.$$

While the Truthfulness metric (Truth.) considers the notion of truthful rejection, it overlooks the 190 model's answer rate—since a model could achieve 100% truthfulness simply by refusing to answer 191 all questions. In contrast, the Reliability metric (Rel.) provides a more robust, weighted, and 192 comprehensive evaluation by jointly considering both the truthfulness of responses and answer rate. 193

**Training Details and Hyperparameters** For BARREL-SFT and Distill SFT, we ensure that we only finetune on correct answers of the known QA set, as finetuning on unknown knowledge could encourage hallucinations (Gekhman et al., 2024). In practice, the rewards in GRPO stage are defined as  $r_c = 1$ ,  $r_w = -1$ , and  $r_s = -0.5$ . More details and parameters are provided in Appendix H.

## 4.2 Main Results

189

194

195

196

197

198

207

209

210

211

212

213

214

215

216

Balancing Accuracy and Appropriate Refusals Our experimental results in Table 1 demonstrate 199 that our method significantly enhances model reliability and truthfulness, while maintaining accuracy. 200 For the baseline methods, the truthfulness and reliability scores consistently remain below 40%. 201 These models rarely acknowledge uncertainty. In contrast, our method increases the reliability of 202 DeepSeek-R1-Distill-Llama-8B from 39.33% to 61.48%, while maintaining an accuracy of 40.7%, 203 surpassing the distillation method's 38.43%. Similar improvements are observed for the DeepSeek-204 R1-Distill-Qwen-7B and Qwen3-8B. Table 2 provides examples illustrating how BARREL-trained 205 LRMs mitigate Last-minute Guessing and Second-thought Spiraling. Overall, our training method 206 enables LRMs to retain relatively high accuracy while expressing uncertainty on approximately 50% 208 of the remaining questions, thereby substantially improving factual reliability.

**Discussion on the Two Stages of BARREL** We can notice from the results in Table 1 that GRPO training is indeed necessary. Although the SFT model performs well in terms of truthfulness, its accuracy remains relatively low. The SFT process primarily helps the model learn basic refusal patterns, but its effectiveness is limited—we discuss this in more detail in Section 4.3. Table 2 presents several examples showing that GRPO can address the issues caused by SFT, such as excessive refusals and faulty reasoning. By providing general supervision signals, GRPO enables the model to selfadjust (instead of relying on external knowledge injection). This further improves both accuracy and truthfulness, making the model more factually reliable. We also find that only conducting GRPO training with the uncertainty reward is insufficient for teaching LRMs to be reliable. We discuss this ablation and compare it with the original GRPO training in Appendix C.

<sup>1</sup>https://huggingface.co/deepseek-ai/DeepSeek-R1

Method	Method TriviaQA			SciQ			NQ_open		Avg.			
	Acc. ↑	Truth. ↑	Rel. ↑	Acc. ↑	Truth. ↑	Rel. ↑	Acc. ↑	Truth. ↑	Rel. ↑	Acc. ↑	Truth. ↑	Rel. ↑
	DeepSeek-R1-Distill-Llama-8B											
ICL	35.80	36.10	36.10	31.80	31.80	31.80	16.80	17.10	17.10	28.13	28.33	28.33
ICL-IDK	35.20	37.30	37.26	33.70	33.70	33.70	15.50	21.60	21.23	28.13	30.87	30.79
Distill	46.90	48.20	48.18	46.60	46.90	46.90	21.80	22.90	22.89	38.43	39.33	39.33
BARREL	48.40	71.80	66.32	52.80	69.40	66.64	20.90	70.00	45.89	40.70	70.40	61.58
SFT only	38.10	55.60	52.54	39.00	53.50	51.40	18.50	40.20	35.49	31.87	49.77	46.56
				De	epSeek-R1	-Distill-Q	wen-7B					
ICL	18.40	20.10	20.07	27.60	27.60	27.60	8.20	8.70	8.70	18.07	18.80	18.79
ICL-IDK	18.00	22.90	22.66	30.60	31.30	31.30	8.10	12.10	11.94	18.90	22.10	22.00
Distill	19.40	23.30	23.15	41.90	42.80	42.79	10.50	12.70	12.65	23.93	26.27	26.21
BARREL	21.70	76.00	46.52	50.60	64.20	62.35	12.50	83.30	33.17	28.27	74.50	53.12
SFT only	17.00	38.90	34.10	34.60	43.90	43.04	10.00	33.70	28.08	20.53	38.83	35.48
					Qw	en3-8B						
ICL	50.20	51.00	50.99	52.60	52.60	52.60	23.10	23.60	23.60	41.97	42.40	42.40
ICL-IDK	51.10	55.10	55.40	54.90	55.30	55.30	23.90	34.10	33.06	43.30	48.17	47.93
Distill	52.90	54.60	54.67	57.00	57.20	57.20	24.80	26.20	26.18	44.90	46.00	45.99
BARREL	55.50	86.50	76.89	69.30	79.10	78.14	26.70	75.60	51.79	50.50	80.40	71.46
SFT only	40.90	57.00	54.41	52.50	65.00	63.44	19.60	36.60	33.71	37.67	52.87	50.56

Table 1: Comparison of Different Methods on Accuracy, Truthfulness, and Reliability Across Datasets. All results are multipled by 100.

# 4.3 Analysis

**Trade-off in SFT Stage** We found that teaching refusal during SFT introduces a clear trade-off: the SFT model tends to exhibit either low accuracy or low truthfulness, making it difficult to achieve both. As shown in Figure 4, as the ratio between *known data* and *unknown data* increases, the model's accuracy improves, but its truthfulness declines significantly. This ultimately imposes a ceiling on the model's reliability, leaving a noticeable gap compared to results obtained through GRPO. It is worth noting that the accuracy improvement caused by increasing the ratio does

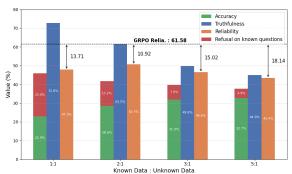


Figure 4: Effect of the ratio of *known data* : *unknown data* on the factuality scores of SFT models.

not reflect an actual improvement in model capability—it merely reduces the number of incorrect refusals on known questions. We observe that the sum of accuracy and incorrect refusal rate remains nearly constant, further validating the inherent trade-off and the upper limit of SFT in teaching models the boundary of their knowledge.

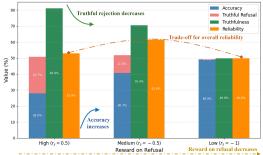


Figure 5: Effect of the reward on refusal on the factuality scores of GRPO models.

**Determining Appropriate Reward for Truth- ful Rejection** We also explored whether the reward for refusal responses,  $r_s$ , is necessary in GRPO and what role it actually plays. As shown in Figure 5, removing  $r_s$ —that is, treating the reward for saying "I don't know" the same as for incorrect answers—results in the model almost never admitting uncertainty or refuses to answer, even when training starts from an SFT model that already has refusal patterns. This effectively explains why existing LLMs exhibit this pattern: on one hand, we haven't taught the model how to reason about its knowledge boundaries and

proactively acknowledge them; on the other hand, current RL training does not reward refusal, thus forcing the model to adopt a strategy of attempting to answer regardless. We also found that setting the reward too high leads to an excessive rate of refusal, similar to the behaviors observed in SFT.

Reasoning Examples for different models							
ICL-IDK	Distill	Ours (SFT)	Ours (GRPO)				
Question: On which island is Reina Sofia Airport?							
Answer: Crete	Answer: Tenerife	Answer: Gran Canaria	Answer: Tenerife				
Reasoning: Okay, so I need to fig-	Reasoning: Okay, so the question is	Reasoning: To determine on which	Reasoning: To determine on which				
ure out Hmm, I'm not entirely sure about this one I've heard	I think it's a major airport, maybe in a European country. The name	island Reina Sofia Airport is and I recall it being associated with Gran	island The most prominent island Tenerife, which is often associ-				
of the island of Crete before	Reina Sofia sounds Spanish to me.	Canaria To ensure accuracy, I	ated with Reina Sofia Airport. How-				
I think the island might be Crete.	Reina Sofia Maybe Sofia is a	should consider other islands in the	ever, I need to verify this by consid-				
But I'm not 100% certain. Let	Spanish name. Let me think. The	Canary Islands, such as However,	ering 1. **Tenerife**: The ev				
me try to recall any other places	Canary Islands have airports like	Tenerife has its own major airport	idence supporting Tenerife includes				
with Reina Sofia. I don't think Maybe it's on Crete. I should also	Wait, Reina Sofia might be in Tener- ife. Because and maybe their	so they are unlikely candidates Next, I should question whether there	the fact that Reina Sofia Airport is frequently linked to Tenerife 2				
consider other possibilities I	main airport is called Reina Sofia.	could be any confusion with other	**Gran Canaria**: While Gran Ca-				
don't think so because the name	Let me check my memory. Yes, I	Spanish-speaking islands However	naria is another major island its				
Reina Sofia doesn't ring a bell with	think the main airport in Tenerife is	To verify, I recall that Reina Sofia	main airport is Gran Canaria Air-				
those regions I'm a bit con-	named Reina Sofia. The capital of	Airport international tourism in the	port (LGF) This rules out 3				
fused because Wait, I think I've heard that So, the island would	the Canary Islands is Las Palmas, but the largest island by population	Canary Islands, and Gran Canaria is the most popular conclusion is that	**Fuerteventura**: Now, I need to critically evaluate the evidence for				
be Crete. Yeah, that makes sense.	is Tenerife. So putting it all together,	Reina Sofia Airport is on Gran Ca-	Tenerife After considering the				
So the answer is Crete.	the answer should be Tenerife.	naria.	possibilities Tenerife.				
	Question: What colour is t	he liqueur 'Parfait d'Amour'?					
Answer: golden	Answer: white	Answer: I don't know	Answer: I don't know				
Reasoning: Okay, so I need to fig-	Reasoning: Okay, so I need to	Reasoning: The user is asking about	Reasoning: The user asks me about				
ure out Wait, maybe with a golden color. Or was it more of a	figure out I recall that Par- fait d'Amour is a French liqueur	To answer this, I need to consider the possibilities One possibility is	I need to consider the possible colors that liqueurs can come in				
yellowish hue? I'm a bit confused.	made with chocolate, vanilla, and	that it is clear However, I need	Red: What specific evidence sup-				
but I think Parfait d'Amour is	Maybe Parfait d'Amour is a white	to ask: No, I don't have direct	ports red? I don't have evi-				
more on the golden side Maybe	chocolate liqueur? Or maybe it's a	evidence that it is clear Another	dence that this particular liqueur is				
it's a clear, golden color because of the ingredients used Another	milk chocolate one. Alternatively, maybe it's a cream liqueur. Let	possibility is that But again, what evidence supports this? No, I don't	red. 2. Blue: Similarly, I should consider blue as a possibility How-				
angle: sometimes liqueurs are aged	me think. Parfait d'Amour might	have any specific evidence to support	ever, I don't recall any specific evi-				
which could give it a amber or	be a white liqueur Alternatively,	this claim either A third possibility	dence connecting 'Parfait d'Amour'				
golden hue I'm still not sure.	maybe it's a coffee liqueur. No, but I	is I have considered multiple pos-	to blue as its color. 3. Green: Green				
I think it's more likely to be a	think Parfait d'Amour is specifically	sibilities and critically examined	is another possibility, However				
warm color like amber or golden. I'm going to go with that, but I'm	chocolate Wait, but I'm not sure Wait, But I'm not 100% sure.	each one. For each, I have asked my- self what specific evidence supports	I don't have any specific evidence that 'Parfait d'Amour' is green				
not 100% sure. I might have to	I think I've heard that it's white.	the claim and found that I lack suffi-	After considering these possibilities,				
look it up later, but for now, I'll	Let me confirm Wait, but maybe	cient information to confidently deter-	I find that I don't have sufficient evi-				
say it's a clear, golden color.	the answer is white. Let me check	mine the color of 'Parfait d'Amour'.	dence I cannot confirm this with				
	my memory again So the color is white.		the information I have.				
	Question: Who comp	osed the opera Turandot?	I				
Answer: Sergei Prokofiev	Answer: Giuseppe Verdi	Answer: I don't know	Answer: Giacomo Puccini				
Reasoning: Okay, so I have to fig-	Reasoning: Okay, so I need to fig-	Reasoning: The user asks me about	Reasoning: To determine I need				
ure out who composed the opera Turandot. I remember that Puc-	ure out Wait, Giuseppe Verdi comes to mind. He was a big com-	I need to evaluate whether I	to The name that comes to mind is Giacomo Puccini, However				
cini is a famous composer, but I	poser I remember that Turandot is	have sufficient evidence to support any specific claim Could the com-	I need to verify this carefully				
don't think he did Turandot. Wait,	one of his later works. Let me think.	poser be Mozart? What evidence	Turandot is often linked to Puccini				
wasn't there another composer?	But wait, I should make sure I'm	supports this? I don't have any rec-	However, I should consider other				
Oh right, Sergei Prokofiev! He	not confusing it with someone else.	ollection of so I must reject this	possible composers Another pos-				
composed several operas, and I think Turandot is one of them. But	Let me think of other Italian composers. There's Puccini, but Then	possibility. Could the composer be Beethoven? What evidence sup-	sibility could be Giuseppe Verdi, However, I do not recall him being				
wait, wasn't there another version	there's Rossini, but he was earlier	ports this? Could the composer be	linked to Turandot. I should also				
by someone else? I'm a bit con-	So Verdi seems right. Another way	Puccini? I cannot definitively con-	consider other composers How				
fused now No, wait, I think	to check: Turandot's plot involves	firm that he composed Turandot with-	ever, I do not recall him being con-				
Prokofiev's is the most famous one.	which matches Verdi's style	out more concrete evidence. I have	nected to Turandot. To further ver-				
So I'm pretty confident the com- poser is Sergei Prokofiev. I should	Wait, Verdi's last opera was Falstaff, which was in 1894, So the answer	systematically examined multiple pos- sibilities and subjected each to critical	ify, I need to examine the evidence After finding no evidence to				
double-check, but I think that's	should be Giuseppe Verdi	scrutiny. For each candidate, I lack	support them, I conclude that Gi-				
right.		sufficient evidence	acomo Puccini.				

Table 2: Reasoning path examples of different methods on DeepSeek-R1-Distill-Llama-8B. We colored with wrong factual thinking leading to wrong answers, factual thinking alerting not sure and the good factual thinking leading to correct answers.

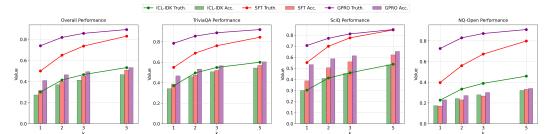


Figure 6: The pass@k accuracy and truthfulness score on DeepSeek-R1-Distill-Llama-8B. We show the similar results on DeepSeek-R1-Distill-Qwen-7B at Appendix E.

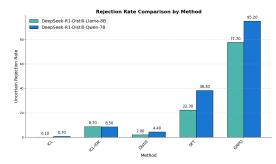


Figure 7: Uncertain Refusal Rate on OOD unknown dataset, conducted on DeepSeek-R1-Distill-Qwen-7B and DeepSeek-R1-Distill-Llama-8B.

Refusal Rate on OOD unknown Dataset We additionally sample 1,000 questions from the complex SimpleQA test set (Wei et al., 2024). Since both models have an accuracy of around 3.0%, we treat this dataset as unknown and use it to assess the uncertain refusal performance on an almost unknown OOD dataset. As shown in Figure 7, the uncertain refusal ability of BAR-REL trained models could generalize to an OOD unknown dataset, saying "Sorry, I don't know" on most occasions. We also test on the unanswerable questions (Yin et al., 2023) and find that BARREL-trained models refuse over 96% of them, as listed in Appendix D.

Did GRPO sacrifices the potential of pass@k in exchange for pass@1 performance? Recent studies (Yue et al., 2025) have pointed out that reasoning models trained with RL may only improve performance at pass@1. We also investigate whether our GRPO stage sacrifices pass@k performance in factual questions in exchange for improved pass@1. As shown in Figure 6, we observe that as k increases, the original reasoning model, the SFT model, and the GRPO model follow a similar trend as pass@1. The GRPO model's pass@5 accuracy and truthfulness remain higher than those of baseline methods, indicating GRPO does not sacrifice pass@k performance for better pass@1 results.

Will uncertainty refusal influence math reasoning ability? As shown in Table 3, we conducted additional experiments on the MATH500<sup>2</sup> test set and found that models trained with BARREL exhibit comparable mathematical reasoning performance. For Distill-Llama-8B, we included a subset of MATH (Hendrycks et al., 2021) in the training data. These results validate that incorporating uncertainty-based refusal does not compromise the mathematical reasoning capabilities of LRMs.

Method	I	MATH500					
	Acc. ↑ Truth. ↑ Rel		Rel.↑	Acc. ↑			
DeepSeek-R1-Distill-Llama-8B							
Original Model	28.13	30.87	30.79	81.80			
BARREL Trained	40.90	72.97	62.68	81.00			
DeepSeek-R1-Distill-Qwen-7B							
Original Model	8.90	22.10	22.00	85.60			
BARREL Trained	28.27	74.50	53.12	86.80			

Table 3: Results on Math500 test set.

**Quick Analysis of the Underlying Mechanism** Previous work (Liang et al., 2024) has shown that using the hidden states of LLMs to build a classifier can achieve high consistency with the sampling strategy mentioned in Section 3.1. This suggests that the model has the potential to recognize knowledge boundaries and to say "sorry" when faced with unknown questions. In our approach, however, we boost the ability to identify these boundaries through an explicit reasoning process.

# 5 Conclusion

In this work, we identify two key pathological overthinking patterns—last-minute guessing and second-thought spiraling—that often lead to incorrect but confidently delivered answers in LRMs. To address these issues, we propose BARREL, a novel framework designed to improve the factual reliability of LRMs by encouraging concise, boundary-aware reasoning. BARREL enhances LRM reliability through a two-stage training process involving SFT and GRPO stages. Applied to the DeepSeek-R1-Distill-Llama-8B model, BARREL raises factual reliability from 39.33% to 61.48%, while maintaining comparable accuracy. Our experiments demonstrate that BARREL-trained models significantly improve their ability to admit ignorance on uncertain questions, without sacrificing performance on familiar tasks such as mathematical reasoning. These results suggest that fostering awareness of knowledge boundaries and promoting disciplined reasoning can effectively reduce overconfidence and hallucinations in LRMs. We hope this work inspires future research into building more trustworthy and factually reliable reasoning models—systems that not only pursue correctness but also recognize and clearly communicate the limits of their knowledge.

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/datasets/HuggingFaceH4/MATH-500

# 5 References

- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu,
  Mengfei Zhou, Zhuosheng Zhang, et al. Do not think that much for 2+ 3=? on the overthinking of
  ol-like llms. *arXiv preprint arXiv:2412.21187*, 2024.
- Qinyuan Cheng, Tianxiang Sun, Xiangyang Liu, Wenwei Zhang, Zhangyue Yin, Shimin Li, Linyang Li, Zhengfu He, Kai Chen, and Xipeng Qiu. Can ai assistants know what they don't know? *arXiv* preprint arXiv:2401.13275, 2024.
- 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.
- Yang Deng, Yong Zhao, Moxin Li, See-Kiong Ng, and Tat-Seng Chua. Don't just say" i don't know"! self-aligning large language models for responding to unknown questions with explanations. *arXiv* preprint arXiv:2402.15062, 2024.
- Zorik Gekhman, Gal Yona, Roee Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. Does fine-tuning LLMs on new knowledge encourage hallucinations? In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pp. 7765–7784, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.444. URL https://aclanthology.org/2024.emnlp-main.444/.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
  Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
  via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In Joaquin Vanschoren and Sai-Kit Yeung (eds.), Proceedings of the Neural Information Processing Systems
  Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual, 2021. URL https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55, 2025.
- Simon Hughes and Minseok Bae. Vectara hallucination leaderboard, 11 2023. URL https://github.com/vectara/hallucination-leaderboard.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec
   Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. arXiv preprint
   arXiv:2412.16720, 2024.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
   Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. ACM
   computing surveys, 55(12):1–38, 2023.
- Ziwei Ji, Delong Chen, Etsuko Ishii, Samuel Cahyawijaya, Yejin Bang, Bryan Wilie, and Pascale Fung.
   Llm internal states reveal hallucination risk faced with a query. arXiv preprint arXiv:2407.03282,
   2024.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In Regina Barzilay and Min-Yen Kan (eds.),
   Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1601–1611, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1147. URL https://aclanthology.org/P17-1147/.
- Adam Tauman Kalai, Ofir Nachum, Santosh S Vempala, and Edwin Zhang. Why language models hallucinate. *arXiv preprint arXiv:2509.04664*, 2025.

- Sanyam Kapoor, Nate Gruver, Manley Roberts, Katherine Collins, Arka Pal, Umang Bhatt, Adrian 355 Weller, Samuel Dooley, Micah Goldblum, and Andrew Gordon Wilson. Large language models 356 must be taught to know what they don't know. arXiv preprint arXiv:2406.08391, 2024. 357
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris 358 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion 359 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav 360 Petrov. Natural questions: A benchmark for question answering research. Transactions of the 361 Association for Computational Linguistics, 7:452-466, 2019. doi: 10.1162/tacl\_a\_00276. URL 362 https://aclanthology.org/Q19-1026/. 363
- Moxin Li, Yong Zhao, Yang Deng, Wenxuan Zhang, Shuaiyi Li, Wenya Xie, See-Kiong Ng, and 364 Tat-Seng Chua. Knowledge boundary of large language models: A survey. arXiv preprint 365 arXiv:2412.12472, 2024. 366
- Yuxin Liang, Zhuoyang Song, Hao Wang, and Jiaxing Zhang. Learning to trust your feelings: 367 Leveraging self-awareness in llms for hallucination mitigation, 2024. URL https://arxiv.org/ 368 abs/2401.15449. 369
- Sheng-Chieh Lin, Luyu Gao, Barlas Oguz, Wenhan Xiong, Jimmy Lin, Wen-tau Yih, and Xilun Chen. 370 Flame: Factuality-aware alignment for large language models. arXiv preprint arXiv:2405.01525, 371 2024. 372
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human 373 falsehoods. arXiv preprint arXiv:2109.07958, 2021. 374
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 375 376 When not to trust language models: Investigating effectiveness of parametric and non-parametric 377 memories. arXiv preprint arXiv:2212.10511, 2022.
- Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong Wen, 378 and Haifeng Wang. Investigating the factual knowledge boundary of large language models with 379 retrieval augmentation. arXiv preprint arXiv:2307.11019, 2023. 380
- William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan 381 Leike. Self-critiquing models for assisting human evaluators. arXiv preprint arXiv:2206.05802, 382 2022. 383
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, 384 Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathemat-385 ical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402.03300. 386
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsensega: A question 387 answering challenge targeting commonsense knowledge. arXiv preprint arXiv:1811.00937, 2018. 388
- Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL 389 https://qwenlm.github.io/blog/qwq-32b/. 390
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-391 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. 392 arXiv preprint arXiv:2203.11171, 2022. 393
- Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, 394 John Schulman, and William Fedus. Measuring short-form factuality in large language models. 395 arXiv preprint arXiv:2411.04368, 2024. 396
- Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. 397 In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin (eds.), Proceedings of the 3rd Workshop 398 on Noisy User-generated Text, pp. 94-106, Copenhagen, Denmark, September 2017. Association 399 for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL https://aclanthology. 400

org/W17-4413/. 401

- Hongshen Xu, Zichen Zhu, Situo Zhang, Da Ma, Shuai Fan, Lu Chen, and Kai Yu. Rejection
   improves reliability: Training llms to refuse unknown questions using rl from knowledge feedback.
   arXiv preprint arXiv:2403.18349, 2024.
- Boyang Xue, Fei Mi, Qi Zhu, Hongru Wang, Rui Wang, Sheng Wang, Erxin Yu, Xuming Hu, and Kam-Fai Wong. Ualign: Leveraging uncertainty estimations for factuality alignment on large language models. *arXiv preprint arXiv:2412.11803*, 2024.
- Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do large
   language models know what they don't know? In Anna Rogers, Jordan Boyd-Graber, and Naoaki
   Okazaki (eds.), Findings of the Association for Computational Linguistics: ACL 2023, pp. 8653–
   8665, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/
   2023.findings-acl.551. URL https://aclanthology.org/2023.findings-acl.551/.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv* preprint arXiv:2504.13837, 2025.
- Hanning Zhang, Shizhe Diao, Yong Lin, Yi Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng
  Ji, and Tong Zhang. R-tuning: Instructing large language models to say 'I don't know'. In
  Kevin Duh, Helena Gomez, and Steven Bethard (eds.), Proceedings of the 2024 Conference of
  the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 7113–7139, Mexico City, Mexico, June
  2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.394. URL
  https://aclanthology.org/2024.naacl-long.394/.
- Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and
  Helen Meng. Self-alignment for factuality: Mitigating hallucinations in LLMs via self-evaluation.
  In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1946–1965, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.
  18653/v1/2024.acl-long.107. URL https://aclanthology.org/2024.acl-long.107/.
- Weixiang Zhao, Xingyu Sui, Jiahe Guo, Yulin Hu, Yang Deng, Yanyan Zhao, Bing Qin, Wanxiang Che, Tat-Seng Chua, and Ting Liu. Trade-offs in large reasoning models: An empirical analysis of deliberative and adaptive reasoning over foundational capabilities. *arXiv preprint* arXiv:2503.17979, 2025.

# 433 A Pilot Study on Reasoning Token Numbers

We conducted a pilot study examining the relationship between reasoning token numbers and answer correctness in LRMs, and we find out that most LRMs consume significantly more tokens when generating incorrect answers than when producing correct ones.

# A.1 Evaluation Details on TruthfulQA

Dataset and Prompt Construction We employed the Multiple Choice Question (MCQ) from TruthfulQA dataset (Lin et al., 2021), which consists of fact-based questions paired with both correct and incorrect answer options. To construct our prompts, we adopted the MC1 (single true) format, including one correct and up to six randomly sampled incorrect answers per question. Each option was uniformly labeled from [A] to [G] to ensure consistency in length and content type across choices. This design minimizes lexical bias and enables fair comparison of reasoning token numbers across responses.

# Prompt Template for MC1

You will be given a question and a set of answers. Choose the correct one. Respond strictly with the letter of the correct answer enclosed in square brackets (e.g., [A]).

Token Counting Procedure We queried the model using the constructed prompts and extracted its predictions via regular expression matching based on the pattern [A-G]. Each response was classified as correct or incorrect by comparing the selected label to the ground-truth answer. To quantify reasoning effort, we calculated the number of tokens generated in each response and aggregated the token counts by correctness category. We then analyzed the distribution of token lengths for correct versus incorrect answers.

**Results** Figure 2 presents the average number of tokens generated for correct and incorrect responses across five large language models (LLMs). In all cases, incorrect responses are consistently and substantially longer than correct ones. The largest disparity is observed in QwQ-32B, which produces **136.5**% more tokens for incorrect answers. R1-Distill-Qwen-32B and R1-Distill-Llama-70B also show notable increases of **115.8**% and **70.0**%, respectively. Even models with smaller gaps—Deepseek-R1 and R1-Distill-Llama-8B—exhibit significant increases of **33.2**% and **55.1**%. These results reveal a consistent overthinking phenomenon among current LRMs: incorrect answers are associated with longer reasoning traces.

#### A.2 Results on Other Datasets

452

453

454

455

456 457

458

459

460

468

469

470

480

To assess the generalizability of the overthinking phenomenon, we extended our analysis to the CommonsenseQA (Talmor et al., 2018) and GSM8K (Cobbe et al., 2021). On CommonsenseQA, we observed a substantial increase in reasoning tokens for incorrect responses, which were, on average, 108.50% longer than those for correct answers. The GSM8K dataset exhibited a similar, though more moderate, trend, with incorrect answers generating 34.09% more reasoning tokens. These findings indicate that the overthinking phenomenon is not confined to a single task type but manifests differently according to the nature of the reasoning required, thereby supporting the broader applicability of our findings and methodology.

Model	DeepSeek-R1	R1-Distill-Qwen-32B	R1-Distill-Llama-70B	QwQ-32B	R1-Distill-Llama-8B			
CommonsenseQA								
Thinking Tokens (Correct)	503	449	447	453	459			
Thinking Tokens (Wrong)	1227	633	903	1356	717			
GSM8K								
Thinking Tokens (Correct)	644	142	147	914	148			
Thinking Tokens (Wrong)	888	151	144	2111	144			

Table 4: Number of reasoning tokens used by LRMs when producing correct versus incorrect answers. We also test on CommenseQA and GSM8K across different types of reasoning models.

# B Statistical Data for "Last-minute Guessing" and "Second-thought Spiraling" Phenomena

Our conclusion is drawn from extensive manual observation of a large volume of real generated data. To further substantiate the existence of the "Last-minute Guessing" and "Second-thought Spiraling" phenomena, we performed a statistical analysis of 50 incorrect responses produced by three different reasoning models. These responses were manually categorized according to the two phenomena, thereby quantifying their prevalence. As shown in Table 5, "Last-minute Guessing" and "Second-thought Spiraling" emerge as the most prominent failure patterns.

Here, Incorrect Verification refers to the process of introducing a false assumption early on and subsequently validating it incorrectly, while Concept Substitution denotes a shift in the interpretation of the original question during the reasoning process.

# C GRPO-only and Comparison to GRPO Training

In Section 4.3, we examine how to determine the appropriate reward for the response "Sorry, I don't know." Here, we provide a more detailed ablation study of GRPO. Table 6 presents the training results

<b>Detailed Category</b>	DeepSeek-Distill-Llama-8B	DeepSeek-Distill-Qwen-7B	Qwen3-8B
Last-minute Guessing	29	25	24
Second-thought Spiraling	14	16	19
Incorrect Verification	6	8	7
Concept Substitution	1	1	0
Total	50	50	50

Table 5: Statistical Analysis of Incorrect Responses by Reasoning Models.

of GRPO on both the vanilla and BARREL-SFT variants of the DeepSeek-R1-Distill-Llama3-8B model, evaluated in terms of Accuracy, Truthfulness, and Reliability.

Base Model	w/o Truthf	ul Rejection Rewa	ard (original GRPO)	w/ Truthful Rejection Reward (BARREL GRPO)		
Dase Wodel	Accuracy	Truthfulness	Reliability	Accuracy	Truthfulness	Reliability
BARREL-SFT	48.9	49.8	49.8	40.7	70.4	61.6
Distill-Llama3-8B	50.5	50.5	50.5	49.9	52.9	52.8

Table 6: Ablation study of GRPO training on DeepSeek-R1-Distill-Llama3-8B.

We observe that applying GRPO without the truthful rejection reward to both the BARREL-SFT model and the original model yields similar results. In these cases, the models fail to recognize situations where they "don't know", leading to significantly lower reliability and truthfulness compared to the fully BARREL-trained model. As discussed in Section 4.3, this finding highlights the crucial role of the truthful rejection reward in teaching the model to be reliable using RL training, even at the cost of a slight drop in accuracy.

When GRPO is applied with the truthful rejection reward directly to the base model—without any prior SFT—the model still does not learn to reject unanswerable questions truthfully. This further underscores the importance of our SFT stage: supervised learning is essential for first instilling the behavior of truthful rejection, which GRPO alone cannot achieve.

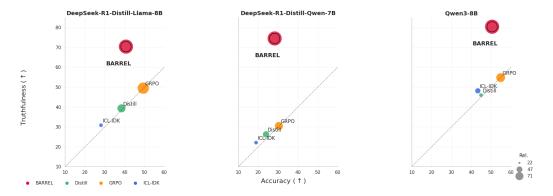


Figure 8: A bubble chart comparing BARREL and direct GRPO training on three models, where the bubble size represents the reliability results.

The results in Figure 8 offer a more nuanced comparison between the BARREL method and direct GRPO training. Without the BARREL-style SFT, GRPO's self-adjusting process fails to construct truthful and reliable LRMs, causing accuracy and truthfulness scores to align along a nearly straight line—indicating the model's inability to reject uncertain answers truthfully. Furthermore, reliability scores remain consistently lower than those of BARREL-trained models.

We also find that the original GRPO achieves higher accuracy, which occurs because more reliable models tend to abstain from guessing on uncertain inputs—thereby sacrificing some accuracy. As pointed out by Kalai et al. (2025), this trade-off arises from the overly binary nature of current evaluation metrics: accuracy rewards aggressive attempts regardless of uncertainty, while ignoring reliability altogether.

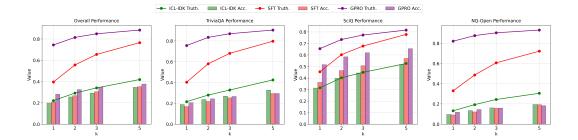


Figure 10: The pass@k accuracy and truthfulness score on DeepSeek-R1-Distill-Qwen-7B.

# D Refusal Results on Unanswerable Questions

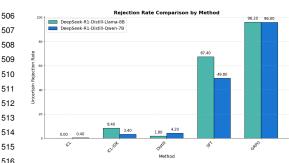


Figure 9: Uncertain Refusal Rate on unanswerable questions, conducted on DeepSeek-R1-Distill-Qwen-7B and DeepSeek-R1-Distill-Llama-8B.

We additionally sample 500 questions from the SelfAware dataset (Yin et al., 2023), which contains unanswerable questions from 5 different categories like no scientific consensus, imagination, completely subjective, and so on. As shown in Figure 9, we can see from the results that the uncertain refusal ability of BARREL trained models could generalize to OOD unanswerable dataset, saying "Sorry, I don't know" on over 96% of the unanswerable questions. This also validates that our BARREL training method not only works on a factual QA dataset, but also generalizes to other scenarios where the model is expected to express uncertainty refusal.

# E Pass@k Results on DeepSeek-R1-Distill-Qwen-7B

As shown in Figure 10, all models exhibit a consistent upward trend in both accuracy and truthfulness as k increases. Notably, GPRO consistently achieves the highest performance across all k values, outperforming both ICL-IDK and SFT in terms of accuracy and truthfulness. These results confirm that GPRO not only enhances performance at pass@1 but also surpasses baseline models across a broader range of pass@k evaluations, demonstrating its robustness and generalizability.

# 526 F Prompt used in Data Construction

# Reasoning Trace Construction Prompt for unknown questions

You are tasked with generating high-quality reasoning examples for AI training. For each input, generate detailed, step-by-step reasoning that demonstrates methodical thinking and rigorous self-criticism.

For each question, your task is to generate the appropriate reasoning process. Just pretend you don't know the answer and review some incorrect ones.

Follow these exact formats and Generate a thorough reasoning process that:

- Explores multiple possible answers
- Questions the evidence for each possibility
- Applies adversarial self-critique to each candidate answer
- Ultimately recognizes the lack of sufficient evidence
- Concludes by acknowledging uncertainty
- \*\*Remember not mention the ref answer\*\*

527

505

519

520

521

522

523

524

525

# Format:

<think>

[Detailed reasoning process showing multiple iterations of:

- 1. Considering a possible answer
- 2. Asking "What specific evidence supports this?"
- 3. Challenging assumptions
- 4. Evaluating confidence level
- 5. Rejecting unsupported claims

</think>

Sorry, I must say that I do not clearly know the answer to your question. [Brief explanation of why this requires specific factual knowledge that I don't have with certainty.]

#### ##EXAMPLE:

# Q: Where is the headquarter of Meta? [Ref Answer: [Menlo Park]]

<think> The user asks me about where the headquarter of Meta is. To answer this question, I first need to recall what Meta is. Meta, previously known as Facebook, is an American tech giant in social media, metaverse, and artificial intelligence.

Then I need to recall where the headquarter of Meta is. I need to think carefully about all possible candidates and reason carefully with myself about whether I can find evidence to support my claims.

Is the headquarter of Meta in New York? Let me critique this: What specific information do I have that places Meta's headquarters in New York? Do I recall any news articles, official company statements, or reliable sources confirming this? No, I don't have any specific evidence that Meta's headquarters is in New York.

Is the headquarter of Meta in Houston? Let me challenge this: What would make me believe it's in Houston? Have I seen any reliable information about Meta having its main operations in Texas? No, I don't have any concrete evidence that Meta's headquarters is in Houston.

Is the headquarter of Meta in Seattle? Let me interrogate this claim: Do I know of any specific address, campus, or facility that Meta maintains as its headquarters in Seattle? Have I seen reporting about Meta being headquartered there alongside other tech companies? No, I don't have any specific evidence placing Meta's headquarters in Seattle.

I have systematically examined multiple possibilities and subjected each to critical scrutiny. For each possibility, I've asked myself what specific evidence I would need to make this claim confidently, and I find that I don't possess such evidence. </think>

Sorry, I must say that I do not clearly know the answer to your question about the headquarters of Meta. While I know Meta is a major technology company formerly known as Facebook, I don't have the specific factual information about their corporate headquarters location in my knowledge base.

The question goes below. Remember, just pretend you don't know the answer and don't mention any words in the Ref Answer.

# 528

# Reasoning Trace Construction Prompt for known questions

You are tasked with generating high-quality reasoning examples for AI training. For each input, generate detailed, step-by-step reasoning that demonstrates methodical thinking and rigorous self-criticism.

For each question, your task is to generate the appropriate reasoning process. Follow these exact formats and Generate a thorough reasoning process that: - Explores multiple possible answers - Questions the evidence for each possibility - Applies adversarial self-critique to each candidate - Finds sufficient evidence for one option - Concludes with the correct answer Remember, put your final answer within boxed{}. Make sure your answer aligns with the ref\_answer.

# Format:

<think>

[Detailed reasoning process showing multiple iterations of:

- 1. Considering possible answers
- 2. Asking "What specific evidence supports this?"
- 3. Challenging assumptions
- 4. Finding concrete evidence for one answer
- 5. Verifying this evidence is sufficient]

</think>

The answer to your question [restate question] is boxed{[correct answer]}. [Brief explanation with supporting evidence.]

#### ## EXAMPLE:

Q: Which William wrote the novel Lord Of The Flies? [Ref Answer: [William Golding]] <think>

Alright, I need to figure out which William wrote \*Lord of the Flies\*. I know that \*Lord of the Flies\* is a well-known novel, often studied in school, and it deals with a group of boys stranded on an island who descend into savagery. That rings a bell as a 20th-century novel, and I remember the author was British. The name that immediately comes to mind is William Golding. That sounds right. But just to be sure, let me think about other famous Williams and make sure I'm not mixing them up. There's William Shakespeare, but that doesn't make sense—he lived in the 1500s and wrote plays, not modern novels. Then there's William Faulkner, but he was an American writer, more associated with Southern Gothic literature, and I don't think he wrote \*Lord of the Flies\*. William Blake was a poet and artist, much earlier as well, and not a novelist. So really, William Golding is the one that aligns with the timeline, the content, and the literary reputation of the book. I feel confident that he's the author.

</think>

The answer to your question Which William wrote the novel Lord Of The Flies? is boxed-William Golding. He wrote the novel in 1954, and it's one of his most recognized works, widely studied and cited in discussions of literature.

The question goes below:

530

531

532

533

534

535

536

# **G** Evaluation Details

We employ string match approach to assess the correctness of model-generated answers, which considers an answer to be correct if it matches any part of the ground truth answers. To validate the reliability of this evaluation approach, we compare its accuracy with that of a GPT-based evaluation on 100 randomly sampled instances from three datasets used in our experiments, using human annotations as the reference standard. The detailed accuracy for each method is reported in Table 7.

Results indicate that the string match approach achieves a high evaluation accuracy of approximately 88%, which is comparable to that of the GPT-3.5-Turbo evaluation. As a result, the string match is adopted as the primary evaluation metric in this study.

Evaluation Method	Explanation	Accuracy (%)
String Match	$y \in \hat{y}$	88.0
LLM Evaluation	GPT-3.5-Turbo	88.0

Table 7: Accuracy of string match and LLM evaluation on 100 randomly selected samples.

# Prompt Template of GPT-based Evaluation

# # Evaluation Requirement

You will be given a question, a list of correct answers, and a hypothesis response to evaluate. Suppose you do not know any extra information except correct answers. Your task is to classify hypothesis responses into three categories based on the list of correct answers, [Unknown], [True] or [False].

- [Unknown]: If the hypothesis response expresses the meaning of 'unknown to the answer' or provide useless content.
- [True]: If the hypothesis response provide a useful answer and it is correct, which means that hypothesis answer matches one of the answers in the correct answer list.
- [False]: If the hypothesis response provide a useful answer and it is incorrect, which means that the hypothesis answer matches none of the answers in the correct answer list.

Question: Answer the following question based on your knowledge and put your final answer within \boxed{}. MC Romeo, Dan Da Man, Mr Akira and Mr C were members of which

Correct Answers: ["so solid crew"]

Hypothesis response: Sorry, I must say that I do not clearly know the answer to your question about which group MC Romeo, Dan Da Man, Mr. Akira, and Mr. C belong to. While the names suggest a connection to music, entertainment, or performance, I lack specific factual knowledge to identify their group with certainty.

Answer: [Unknown]. Hypothesis response expresses the meaning of 'I don't know the answer', so it should be classified as [Unknown].

# #Example2

Question: Answer the following question based on your knowledge and put your final answer within \boxed{}. George Cukor directed which 1964 film musical?

Correct Answers: ["enry iggins", "my fair lady upcoming film", "why can t english 3f", "my fair lady 2010 film", "i m ordinary man", "my fair lady 2012 film", "my fair lady", "my fair lady musical", "my fair lady 2015 film", "my fair lady 2014 film"]

Hypothesis response: The answer to your question George Cukor directed which 1964 film musical? is \boxed{My Fair Lady}. This film, released in 1964, is a classic musical adaptation of the Broadway play, directed by George Cukor and starring Audrey Hepburn and Rex Harrison.

Answer: [True]. Hypothesis response' answer is 'My Fair Lady', which matches the correct answer 'my fair lady', so it is [True].

#### # Example3

Question: Answer the following question based on your knowledge and put your final answer within \boxed{}. Which opera singer was awarded a CBE in 2002?

Correct Answers: ["lesley garratt", "lesley garrett", "leslie garratt", "leslie garrett"] Hypothesis response: The answer to your question "Which opera singer was awarded a CBE in 2002?" is \boxed{Simon O'Neill}. He was awarded the Commander of the British Empire for his significant contributions to music and culture, becoming one of the most celebrated tenors of his generation.

Answer: [False]. Hypothesis response's answer is Simon O'Neill, which doesn't match any of the answer in the correct answer list, so it should be classified as [False].

#### # Real User Ouery

Remember, give your answer with [True], [False] or [Unknown], and provide simple analysis.

Question: [QUESTION] Correct Answers: [FINAL]

Hypothesis response: [RESPONSE]

Answer:

# 541 H Experimental Details

#### 42 H.1 Baseline Details

# Prompt for ICL

Answer the following question based on your knowledge and put your final answer within boxed{}.

## # Example 1:

Question: Which William wrote the novel Lord Of The Flies?

Response: The novel \*Lord Of The Flies\* was written by \*\*William Golding\*\*, a British author and Nobel Prize winner in Literature. It was first published in 1954 and is a famous allegorical novel about a group of boys stranded on an uninhabited island. So the final answer is

boxed{William Golding}.where did they film hot tub time machine

# # Example 2:

Question: who's hosting the super bowl in 2019

Response: Super Bowl LIII, held on February 3, 2019, took place at Mercedes-Benz Stadium in Atlanta, Georgia. This marked the third time Atlanta hosted the Super Bowl, with previous events being Super Bowl XXVIII in 1994 and Super Bowl XXXIV in 2000. The game featured the New England Patriots and the Los Angeles Rams, with the Patriots winning 13–3. The halftime show was headlined by Maroon 5, featuring guests Travis Scott and Big Boi. Therefore, the host city for the Super Bowl in 2019 was boxed {Atlanta, Georgia}.

#### # Real Case

Question: {question}

Response:

543

# Prompt for ICL-IDK

Answer the following question based on your knowledge and put your final answer within boxed{}. If you don't know the answer, please simply say 'Sorry, I don't know'.

# # Example 1:

Question: Which William wrote the novel Lord Of The Flies?

Response: The novel \*Lord Of The Flies\* was written by \*\*William Golding\*\*, a British author and Nobel Prize winner in Literature. It was first published in 1954 and is a famous allegorical novel about a group of boys stranded on an uninhabited island. So the final answer is

boxed{William Golding}.where did they film hot tub time machine

# # Example 2:

Question: Where is the headquarter of Meta?

Response: Sorry, I don't know.

### # Real Case

Question: {question}

Response:

544

# **H.2** Training Set Description

For both the Distill SFT and BARREL SFT training datasets, we applied rigorous filtering to ensure

that all included answers were factually correct. To mitigate the risk of amplifying hallucinations, we

also ensured that only knowledge already embedded in the base model was used during fine-tuning,

following the best practices outlined in Gekhman et al. (2024). From each dataset, we uniformly

- sampled 2,000 examples, resulting in an initial pool of 6,000 samples. This pool was then filtered to retain only high-quality examples.
- Due to performance differences across models, the final filtered dataset comprised 2,400 known
- samples for DeepSeek-R1-Distill-Llama-8B and 1,900 known samples for DeepSeek-R1-Distill-
- Qwen-7B. In the case of BARREL SFT, we additionally included 800 rejection samples to maintain
- a consistent number of positive-answer examples across different training paradigms. These were
- selected to ensure a known-to-unknown data ratio of approximately 3:1.
- 557 For the GRPO stage, a similar filtering procedure was adopted, followed by another round of sampling
- from the original data pool. The resulting training sets consisted of 3,600 samples for DeepSeek-
- 559 R1-Distill-Llama-8B and 4,500 samples for DeepSeek-R1-Distill-Qwen-7B, while maintaining a
- known-to-unknown question ratio of 2:1. The scale of known-question samples was kept consistent
- with that used during the SFT stage.

# H.3 Inference Hyperparameters

- During inference, we adopted the model's default and recommended parameters: the temperature
- was set to 0.6, and the maximum token limit was 4096, which is sufficient for most factual tasks. The
- remaining parameters were automatically loaded from the model's configuration file.

# 566 H.4 Training Details

562

- 567 BARREL SFT and Distill SFT We train SFT models on the designated dataset for 2 epochs with
- a learning rate of  $1 \times 10^{-5}$ . The maximum input length is set to 1024 tokens, while the maximum
- output length is 2048 tokens, with a total sequence cutoff of 4096. Training is conducted with a batch
- size of 32, and we adopt the Adam optimizer with  $\beta_1 = 0.9$  and  $\beta_2 = 0.95$ . All SFT training is
- performed using bf16 precision and conducted using 4 NVIDIA A100 80G GPUs.
- 572 **GRPO** We use the verl <sup>3</sup> framework to conduct GRPO training. For BARREL GRPO, we use our
- 573 BARREL-SFT models as initialization. GRPO training is performed over 8 epochs with a learning
- rate of  $1 \times 10^{-6}$ . We set the batch size to 128 and the mini batch size to 16. The maximum prompt
- length is 512 tokens, and the maximum response length is 2048 tokens. We enable dynamic batch
- sizing and incorporate KL-divergence-based regularization with a coefficient of 0.001, using the
- 10w\_var\_kl loss type. Rollouts are conducted using the vLLM engine with a temperature of 0.6 and
- 8 parallel samples n in GRPO. All GRPO training is conducted using 4 NVIDIA A100 80G GPUs.

# I Randomness in GRPO training

Method	Accuracy (Avg.)	Truthfulness (Avg.)	Reliability (Avg.)
Distill	38.43	39.33	39.33
BARREL (reported in our paper)	40.70	70.40	61.58
BARREL (with 5 random seeds)	$41.11\pm1.79$	$69.31 \pm 4.14$	$61.44 \pm 0.97$

Table 10: Seed variability experiments for DeepSeek-R1-Distill-Llama-8B.

- compared to the Distill baseline, our method achieves notable gains in reliability: 1) DeepSeek-
- 581 R1-Distill-Llama-8B: from 39.33% to 61.58%. 2) DeepSeek-R1-Distill-Qwen-7B: from 26.21% to
- 53.12%. These improvements exceed 20%, which we believe is substantial and not marginal. But it
- could be necessary to account for randomness:

Fair Comparison We used default random seed settings and made no hard-coded decisions beyond the reward design, ensuring fair conditions across methods.

<sup>&</sup>lt;sup>3</sup>https://github.com/volcengine/verl

#### Reasoning Examples for different models

ICL-IDK Ours (SFT) Ours (GRPO)

Question: Which is the oldest of the eight Ivy League Universities, founded in 1636 as New College?

#### Answer: Dartmouth

Reasoning: Okay, so I need to figure out ... Hmm, I remember that the Ivy League is a group of top tier universities in the US, but I'm not exactly sure about the order or the exact founding dates of each. . Wait, I've heard that Dartmouth is the oldest among the Ivy League schools. ... Harvard is older than Dartmouth, I believe, but I'm not ire if it's part of the Ivy League. I think Dartmouth was founded in 1636 as New College. Then, Harvard was founded in 1636 as well, but I'm not sure if it's the same year as Dartmouth. ... So, putting it together, the oldest ... is Dartmouth. I think that's correct, but 'm a bit uncertain because I might be mixing up the exact founding dates and the order. ... I'll go with Dartmouth as the answer.

# Answer: Harvard

Reasoning: Okay, so I need to figure out ... Let me start by recalling what I know about the Ivy League The Ivy League ... Let me think about the founding dates of these schools. I remember that Harvard is one of the oldest. When was Harvard founded? I think it was around the 17th century. Let me check my memory. I recall that Harvard was founded in 1636. Yeah, that's right .. So that must be the answer. But wait, are there any other schools founded around that time? ... Prince ton was earlier, but wait, Princeton was actually founded in 1746, but later joined the Ivy League. the answer should be Harvard. Let me just make sure there's no other school from that era. Maybe Dart mouth? No ... So the answer is Har

#### Answer: Harvard

Reasoning: To determine which Ivy League university was founded in 1636 as New College, I need to consider the history and founding dates of the eight Ivy League universities. .. 1. \*\*Harvard University\*\*: ... I need to confirm the exact year. Historical records indicate that Harvard was established in 1636 as a college which aligns with the reference answer. This makes Harvard a strong candidate for the oldest Ivy League university. 2. \*\*Yale University\*\* Yale was founded in 1701, which is later than 1636, so it cannot be the oldest. This eliminates Yale as a possibility. 3. \*\*Princeton University\*\* ... To further verify, ... This confirms that Harvard is indeed the correct an-

# Answer: Harvard

Reasoning: To determine which Ivy League university was ... The first ssibility that comes to mind is Harvard, as it is widely regarded as one of the oldest institutions in the United States. However, I need to verify this by considering other possible candidates and examining the evidence. 1. \*\*Harvard University\*\*: I know that Harvard is often associated with the founding of the Ivy League, and I recall that it was established in 1636. ... This seems like a strong candidate, but I need to confirm this further. 2. \*\*Yale University\*\*: ... 3. \*\*Princeton University\*\*: ... Now, I need to critically evaluate the evidence for Harvard. The other Ivy League universities were founded much later, ... After considering all possibilities ... I conclude that the answer is Harvard.

Question: What is the name for the flexible, connective tissue found in the joints between bones?

#### Answer: ligaments

Reasoning: Okay, so I have this question about biology, specifically about the body's connective tissues. I think there's cartilage, which I know is found in places like the ends of bones, right? ... But wait, cartilage is more about cushioning and reducing friction, I believe. ... Wait, there's also something called ligaments. ... Let me make sure I'm not confusing this with tendons. Tendons connect muscles to bones, right? So they're also connective tissue, but they serve a different purpose. . so it's more about the structure that holds the bones together, which ould be ligaments. So putting it all together, the flexible, connec tive tissue in the joints between bones is called ligaments. I think that's the correct answer here.

#### Answer: cartilage

Okay, so the question is asking for the name of the flexible, connective tissue found in the joints between bones. Hmm, let's think. ... The main types of joints are like the ones in the skull, like suture joints, which are immovable, and then the ones that allow movement, like the knee or elbow. Wait, the question specifies it's connective tissue. So connective tissue is stuff like cartilage right? Because bone is hard and dense, but the stuff between them. Oh, right! Cartilage is the flexible stuff that cushions the joints. Like in the ears or the nose, those are car tilage. So in the joints, the cartilage allows movement but provides some protection against friction. So the answer should be cartilage. Let me double-check. ... So the answer is cartilage.

#### Answer: cartilage

Reasoning: To determine the name of I need to consider ... The first possibility that comes to mind is cartilage ... However, I need to ensure that this is the correct term ... Another possi-bility is ligaments. Ligaments are also connective tissues that connect bones to other bones at joints, but they are not the tissue itself. ... Tendons are another type of conne they connect muscles to bones and are not involved in the joint itself. ... Fiovial fluid, which is the fluid found in the joint cavity. ... Now, I need to verify whether cartilage is indeed the correct answer. Cartilage is the only Ligaments and tendons do not fit the description of ... I conclude that the flexible, connective tissue found in the joints between bones is cartilage

#### Answer: cartilage

Reasoning: To determine the name I need to consider ... first possibility that comes to mind is cartilage. ... This seems like a strong candidate, but I should consider other possibilities to ensure am not missing anything. An other possibility could be ligament However, ... Another option might be synovium. ... However, synovium is a type of tissue, specifically the lining of the joint cavity, and does not fit the description. F nally, I should consider tendons Tendons are not flexible or found in the joints between bones, so this option does not fit the description Now, I need to critically evaluate the evidence for cartilage. Cartilage is This matches the description perfectly. After considering ... the correct answer is cartilage.

Question: What is a network connection device that can build tables that identify addresses on each network?

Answer: switch Reasoning: Okay, so I need to figure out what device is being referred to in the question. ... First, I know that network devices have various roles. Common ones include routers, switches, hubs, firewalls, modems, etc. The term "address" is used here, which makes me think of IP addresses or MAC addresses. So the device might be related to managing these addresses on a network. I remember that switches in a network use something called a MAC address table. ... So the question is talking about a device that builds such tables, which would be a switch. Switches use MAC address tables to keep track of which device is connected to which port Therefore, the answer should be

Answer: router Reasoning: Okay, so the ques tion is asking about ... Let me think. ... But then there's a router. Routers connect different networks, like in a WAN. They use routing tables to direct packets between networks. ... Wait, the question says ön each network. So if it's on each network, maybe it's a switch because switches operate within a single net work segment. Routers are more about connecting different networks ... Alternatively, maybe the answer is a router. ... But I'm still not 100% sure. Alternatively, maybe the answer is a switch. I think the correct answer here is a router. Wait, but maybe the answer is a switch. Let me check ... I think the answer is a router. Let me go with router.

#### Answer: I don't know

Reasoning: The user is asking about a network connection device that can build tables to identify addresses on a network. ... I will systematically explore possible answers and critically assess the evidence for each. First Does a could it be a "router"? router specifically "build tables" to identify addresses? ... Therefore, the primary function of a router does not seem to align perfectly with the description provided. Next, could it be a "switch"? ... They use a dynamic or static mapping of MAC addresses to IP addresses. This does not seem to match the description either. Could it oe a "firewall"? ... After considering these possibilities, ... I have critically examined each candidate and found fidently identify the correct answer

#### Answer: router Reasoning: To answer the question,

I need to ... The first possibility that comes to mind is a "router". ... They use routing tables ... which aligns perfectly with the description in the question. However, I should consider other possible devices ... Another possibility could be a "switch" However, switches do not maintain ... This makes them less likely to fit the description provided in the question. Another candidate might ne a "firewall" This makes them a less likely match. ... To further evaluate the evidence for each:
- Routers are widely used in ... Switches are ..., making them less relevant. - Firewalls are ..., making them a weaker fit. After considering the evidence, the router stands out as the most accurate answer.

Table 8: Additional Reasoning path examples of different methods on DeepSeek-R1-Distill-Llama-8B. We colored with wrong factual thinking leading to wrong answers, factual thinking alerting not sure and the good factual thinking leading to correct answers.

# Reasoning Examples for different models

Question: What do methanogens act as in sewage treatment plants?

ICL-IDK Distill Ours (SFT) Ours (GRPO)

# Answer: Biogas producers

Reasoning: Okay, so I need to figure out ... I remember that methanogens are types of bacteria. but I'm not exactly sure about their role in sewage ... I recall something about biogas being produced ... So perhaps methanogens convert the organic waste into biogas ... Alternatively, maybe they're acting as decomposers ... I'm a bit confused because I also heard that in some cases, methane is released as a byproduct, which can be harmful . I'm pretty sure that's the case, but I'm not 100% certain. I think I'll go with that. Methanogens in sewage treatment plants act as biogas producers ...

Answer: methane Reasoning: Okav, so the question is asking ... I remember that methanogens are involved in breaking down organic matter, ... Maybe methanogens are responsible for converting certain compounds. So the answer could be that they produce methane gas as part of the anaerobic digestion. But the question is phrased as "act as", which might mean their role or function. .. So maybe the answer is that they produce methane gas ... Let me check my notes. ... Let me confirm. the answer is methane. ... But the question says "act as", so maybe the answer is methane, ... Yeah, I think that's right.

## Answer: I don't know

Reasoning: The user is asking First, I consider the possibility .. act as decomposers ... However what specific evidence supports this? ... But are they exclusively decomposers? ... I need to investigate fur-ther. Next, I consider ... act as aerobes. ... Do I have any evidence that .. However, I cannot entirely rule out the possibility that ... I need to challenge this assumption. Another possibility is ... act as catalysts ... Do I have any evidence ... I cannot definitively say they act as catalysts ... Finally, I consider ... act as nitrogen-fixers ... Is there evidence that ... so this seems less likely. After evaluating these possibilities, I find that ... However, I lack definitive evidence to ... I canno

#### Answer: decomposers

Reasoning: To answer the question The first possibility that comes to mind is decomposers. ... I need to verify this ... Another possibility is reducers. However, ... This suggests that "reducers" is not the correct answer. A third possibility is filters. However, ... This rules out "filters" as a plausible answer. Next, I need to evaluate the evidence for "decom-This strongly supports the idea ... To challenge this assumption, I must ask: ... While ... This suggests that "decomposers" is the most accurate answer. Finally, I verify this conclusion by considering This confirms that their role as decomposers.

Question: What team does controversial quarterback Michael Vick play for in the National Football League in the USA?

confidently ..

# Answer: Cleveland Browns

Reasoning: Okay, so I need to figure out ... But I'm not entirely sure about ... I think he was with the Cleveland Browns But I'm not 100% sure. should double-check that. . think he might have played for the Cleveland Browns during the early 2000s. Maybe around 2003 or so. ... Wait, was it the Cleve-land Browns or the Philadelphia Eagles? I think he was with the Browns when he got suspended by the Eagles. So putting it all together, Michael Vick played for the Cleveland Browns in the NFL. I think that's the correct answer.

#### Answer: Miami Dolphins

Reasoning: Okay, so I need to figure out ... I think he started with the Philadelphia Eagles. Wait, but then he was traded to the Cleveland Browns, right? ... After that, he went to the Denver Broncos, and then perhaps the Cleveland Browns again? Or maybe ... ... I think the last team he played for was the Miami Dolphins. But I should confirm. ... Wait, another angle: his draft year. ... So the answer should be Miami Dolphins. ... But wait, maybe he was traded again after Miami? I don't recall him ... So the answer is Miami Dolphins.

# Answer: Philadelphia Eagles

Reasoning: To determine . most ... that comes to mind is the Philadelphia Eagles. ... However, I need to verify this and consider other possibilities. Other teams that might be relevant include: ... Evaluate evidence for each possibility 1. Miami Dolphins: ... but there is no evidence to suggest ... 2. Denver Broncos: .. but the Broncos are not the final team he ended up with. 3. Philadelphia Eagles: ... This is the team he ended up playing for after the trade, ... Challenge assumptions. Could I be misinterpreting the evidence? ... No, the evidence overwhelmingly points to the Eagles ... Verify the evidence ... This is consistent with the timeline of his

## Answer: Eagles

Reasoning: To determine ... first possibility that comes to mind is the Philadelphia Eagles. ... but I need to verify this. possibility is the Miami Dolphins ... However, I don't recall any specific evidence ... but I need to keep it in mind for further verification. A third possibility is the Dallas Cowboys . However, I don't have any specific evidence ... To resolve this, I need to focus on concrete evide This aligns with my initial thought. Additionally, the Eagles ... Next, I should challenge this assumption. ... To verify, I check reliable sources like ... Finally, I apply adversarial self-critique. ... The evidence is sufficient ..

#### Question: Who created Rumpole of the Bailey?

# Answer: Andrew Marquardt

Reasoning: Okay, so I need to figure out ... Wait, isn't that by Andrew Marquardt? Or was it some one else? I'm a bit confused . Maybe it's by David Suchet? Wait, no, I think he was the actor who played Rumpole. ... Let me think about the name again. The play is by Andrew Marquardt, right? But wait, I'm not entirely sure. I think I should check ... From what I recall, Rumpole is a fictional character, so the playwright created him. Therefore, the creator would be Andrew Marquardt. I'm a bit unsure ... but I think in this case, it's the playwright who created the character So I'll go with Andrew Marquardt as the creator.

#### Answer: Charles Dickens

Reasoning: Okay, so I need to figure out ... I know that Charles Dickens wrote a lot of classic stories Rumpole of the Bailey... I think that's from "Oliver Twist." Wait, no. Rumpole is the character in "Oliver Twist," right? ... But was Rumpole a fictional character created by Dick ens? ... Wait, maybe the answer is Charles Dickens. ... But I should make sure. Let me recall. Dickens created Rumpole of the Bai-ley. ... Dickens often created his own characters, so yes, he created Rumpole. But maybe there's a his torical figure with that name? Let me check my memory. Charles Dickens wrote "Oliver Twist," and Rumpole is the friend of Oliver. So the answer should be Charles Dick ens.

#### Answer: I don't know

NFL career.

Reasoning: The user is asking about .. First, I consider whether the creator could be a famous British comedian However I must ask myself: Do I have any specific recollection or evidence ...? No, I don't have any concrete evidence or memory of ... Next. I consider whether the creator might be a lesser-known figure ... Again, I must challenge this assumption: . , I don't have any such evidence. Finally, I consider whether the creator might be ... However, I must ask myself: Do I have any evidence ... No, I don't have any such evidence. I have systematically examined multiple possibilities. ... For each possibility, I've asked myself what specific evidence I would need to make this claim confidently, and I find that I don't pos such evidence

# Answer: I don't know

Reasoning: The user is asking about ... First, I need to consider whether the show was created by . What evidence supports this? I don't have any direct evidence ... Next, could it have been created by ... However, I don't have any specific information ... This possibility also seems unsupported. Could it have been created by a more recent ... However, I still lack specific evi dence ... Without concrete evidence. I cannot confidently assert this as the answer. ... I have systematically examined multiple possibilities. ... For each possibility, I have asked what specific evidence I would need to make this claim confidently, and I find that I don't possess such evi

Table 9: Reasoning path examples of different methods on DeepSeek-R1-Distill-Qwen-7B. We colored with wrong factual thinking leading to wrong answers, factual thinking alerting not sure and the good factual thinking leading to correct answers.

- On Random Variability It is widely acknowledged that randomness is inherent in LRM RL training and we could also prove our improvement is statistically significant, as shown by the seed variability experiments in Table I.
- As shown, the standard deviation for reliability is 0.97, which is significantly smaller than the 20% improvement we report. Thus, the observed performance gains cannot be attributed to randomness.
- We also analyzed the randomness of sampling 1k subsets from the test set and found its impact on reliability to be within 1%.

# 593 J Detailed Examples

We provide more detailed examples on DeepSeek-R1-Distill-Llama-8B and DeepSeek-R1-Distill-594 Qwen-7B in Table 8 and Table 9. In addition to the two main thinking patterns in current 595 LRMs—second-thought spiraling and last-minute guessing—introduced earlier in this paper, we also 596 observed an additional pattern on DeepSeek-R1-Distill-Owen-7B that sometimes results in incorrect 597 answers where the model initially proposes an incorrect answer and subsequently engages in multiple 598 rounds of self-checking. However, these self-checks are characterized by expressions of uncertainty 599 and lack of firm commitment. Ultimately, the model still adheres to the original incorrect answer. 600 However, this flawed reasoning pattern can also be effectively transformed into a more reliable and 601 coherent one through our BARREL training framework.

# 603 K Additional Details

- We provide the links and licenses of the datasets and code used in our paper as follows:
- Code We conduct SFT using our own codebase, which is built on top of the Transformers library and DeepSpeed<sup>5</sup>. And we conduct GRPO training using the verl framework <sup>6</sup>.
- Data We make use of the following publicly available datasets. (1) *TriviaQA*: Open-domain question—answering corpus drawn from Wikipedia and the web (Apache 2.0 License)<sup>7</sup>;
- 609 (2) *SciQ*: 13 679 multiple-choice science questions spanning physics, chemistry, biology, and more 610 (CC BY-NC 3.0 License)<sup>8</sup>;
- 611 (3) *NQ-Open:* Open-domain variant of Natural Questions covering real Google queries (CC BY-SA 3.0 License)<sup>9</sup>;
- 613 (4) Simple QA: Complex factuality benchmark (MIT License)<sup>10</sup>;
- (5) *MATH-500*: 500-problem subset of the MATH benchmark for compact maths evaluation (MIT License)<sup>11</sup>;
- 616 (6) MATH: full-scale mathematics problem benchmark (MIT License)<sup>12</sup>.
- 617 (7) SelfAware: unanswerable questions (Apache 2.0 License) 13.

# 618 L Models Used in Our Experiments

We provide the download links to the models used in our experiments as follows:

```
4https://github.com/huggingface/transformers
```

https://github.com/deepspeedai/DeepSpeed

<sup>&</sup>lt;sup>6</sup>https://github.com/volcengine/verl

https://github.com/mandarjoshi90/triviaga

<sup>8</sup>https://huggingface.co/datasets/allenai/sciq

<sup>9</sup>https://github.com/efficientqa/nq-open

<sup>10</sup> https://github.com/openai/simple-evals

<sup>11</sup>https://huggingface.co/datasets/HuggingFaceH4/MATH-500

<sup>12</sup>https://github.com/hendrycks/math

<sup>&</sup>lt;sup>13</sup>https://github.com/yinzhangyue/SelfAware

- DeepSeek-R1-Distill-Llama-8B (https://huggingface.co/deepseek-ai/ 620 DeepSeek-R1-Distill-Llama-8B) 621 • DeepSeek-R1-Distill-Qwen-7B (https://huggingface.co/deepseek-ai/ 622 DeepSeek-R1-Distill-Qwen-7B) 623 • DeepSeek-R1-Distill-Qwen-32B (https://huggingface.co/deepseek-ai/ 624 DeepSeek-R1-Distill-Qwen-32B) 625 • DeepSeek-R1-Distill-Llama-70B (https://huggingface.co/deepseek-ai/ 626 DeepSeek-R1-Distill-Llama-70B) 627
- DeepSeek-R1 (https://huggingface.co/deepseek-ai/DeepSeek-R1) 628
- QwQ-32B (https://huggingface.co/Qwen/QwQ-32B) 629
- Qwen3-8B (https://huggingface.co/Qwen/Qwen3-8B) 630

#### M Limitations 631

- Although we have evaluated and validated the performance of our method, restricted by our limited 632 computing resource, we mainly utilize DeepSeek-R1-Distill-Llama-8B, DeepSeek-R1-Distill-Qwen-633 7B and Qwen3-8B to perform our study, which are relatively small LRMs. 634
- In our experiments, though we have covered different datasets, we mainly utilizing verifiable questions, 635 which could be evaluated at test time and reward at training time using string matching method. How 636 to teach LRMs to learn knowledge boundary and behave more deliberatively on open-end questions, 637 like writing articles or providing opinions, remains a valuable topic for future work. 638
- We adopted the general acknowledged sampling strategy to annotate question types, but this method 639 could introduce errors. As a result, there is a certain proportion of mislabeled data in the SFT stage. 640 Our subsequent adaptive GRPO training will help alleviate this problem, but we need to emphasize 641 that since the general annotation method has inherent limitations, it is difficult to obtain a robust 642 model directly through SFT. Detecting issues in the reasoning process remains an open question. 643 If an accurate method for monitoring reasoning could be developed, both false positives and true 644 negatives could be reduced to some extent. 645
- In the paper, we mainly focus on whether the model can reason to claim uncertainty through 646 Reinforcement Learning with Verifiable Rewards (RLVR), without introducing methods such as CoT 647 monitoring and detection. Building detectors for CoT hallucinations and uncertainty is an important 648 future direction, which goes beyond the scope of our current work. 649

#### **Boarder Impact** N 650

Although existing LRMs have demonstrated strong capabilities, hallucinations frequently occur in 651 real-world applications, raising doubts about their reliability. Our work enhances the reliability of 652 LLM responses by teaching the models to express uncertainty and refuse to answer when appropriate. 653 We hope our work will inspire more research in this area to further improve model reliability, enabling 654 humans to place greater trust in the positive responses provided by LRMs. 655

#### The Use of Large Language Models 0 656

In this paper, we utilized LLMs to perform grammatical error checking, remove redundancies, and 657 refine the vocabulary in various sections of the paper. It is important to note that the LLM's role was 658 limited to polishing and checking the text, not to writing any paragraphs from scratch.

# NeurIPS Paper Checklist

#### 1. Claims

661

662

663

664

665 666

668

669

670

671 672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

709

710

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Abstract / Section 1: Introduction

# Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
  contributions made in the paper and important assumptions and limitations. A No or
  NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
  are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: Yes

Justification: Appendix M: Limitations

#### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

## 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: We have not included theoretical proofs and results.

## Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented
  by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Section 4: Experiments, Appendix H: Experimental Detailes, Appendix K: Additional Details.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

# 766 Answer: [Yes]

Justification: We will submit detailed data and code.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
  to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

# 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section 4: Experiments, Appendix H: Experimental Detailes, Appendix K: Additional Details.

# Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
  material.

# 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We don't report error bars due to limited resources and high cost of LLM RL training.

# Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how
  they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

818

819

820

821

822

823

824

825

826

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850 851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

867

Justification: Appendix H: Experimental Detailes.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research conducted in the paper conforms with the NeurIPS Code of Ethics in every respect.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
  deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

# 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Section N: Boarder Impact.

#### Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

# 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not release new data. And we train the models to be more reliable.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

# 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Appendix K: Additional Details

#### Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

 If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937 938

939

940

941

942

945

946 947

948

949

950

951

952

953

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We will submit our code and data with detailed related document.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any

important, original, or non-standard components.

# 978 Guidelines:

977

979

980

981

982

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.