TRAIN FOR TRUTH, KEEP THE SKILLS: BINARY RETRIEVAL-AUGMENTED REWARD MITIGATES HALLUCINATIONS

Anonymous authors

000

001

002

004

006

008 009 010

011 012 013

014

016

017

018

019

021

025

026

027 028 029

031

034

040

041

042

043 044 045

046 047 048

052

Paper under double-blind review

ABSTRACT

Language models frequently generate factually incorrect information with high confidence, a phenomenon known as extrinsic hallucination. Existing approaches for improving factuality often come at the cost of diminished performance on other downstream tasks, limiting their practical deployment. We propose a novel onpolicy reinforcement learning (RL) approach that uses binary retrieval-augmented rewards (RAR) to address this challenge. Our binary reward scheme assigns a reward of zero whenever any factual error is detected and one otherwise. We evaluate our method through continual RL from Qwen3 models across multiple tasks. For open-ended generation, binary RAR achieves a 39.3% reduction in hallucination rates, significantly outperforming supervised training or online RL with dense reward. In short-form settings, models learn calibrated abstention, answering "I don't know" when parametric knowledge is insufficient, leading to 44.4% and 21.7% fewer incorrect answers on PopQA and GPQA, respectively. Crucially, these factuality gains come without performance degradation on instruction following, math, or code, whereas dense-reward RL, despite improving factuality, induces quality regressions.

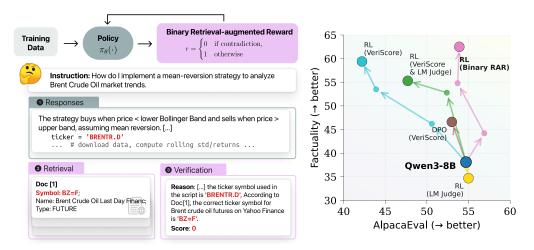


Figure 1: Left: Reinforcement Learning with Binary Retrieval-Augmented Reward (RAR). Right: Binary RAR mitigates hallucination while preserving long-form generation quality.

1 Introduction

Large language models (LMs) have demonstrated remarkable capabilities in knowledge memorization and can recall this knowledge for reasoning in complex tasks, fundamentally transforming how people seek and process information (Chatterji et al., 2025). However, language models may generate plausible but factually incorrect content, a behavior known as extrinsic hallucination (Kalai

et al., 2025). This issue has become particularly concerning as recent state-of-the-art reasoning models produce incorrect information more frequently (Yao et al., 2025; Song et al., 2025).

Reducing hallucination is challenging because supervised learning on static datasets primarily optimizes next-token likelihood and does not enforce factual correctness at inference time (Newman et al., 2025; Zhang et al., 2024). Recent post-training methods instead optimize dense, scalar factuality assessments, typically via direct preference optimization (DPO) and report reductions in hallucination (Tian et al., 2024; Lin et al., 2024). However, a central challenge remains: *mitigating hallucinations without degrading overall utility across diverse instructions*.

In this paper, we show that it is possible to substantially reduce hallucination without hurting task performance. We introduce an *online* reinforcement learning framework with a binary retrieval-augmented reward (RAR; Figure 1 left). To compute RAR, we first retrieve candidate evidence from the web, then verify the factual correctness of a language model's output against the retrieved documents; we assign a binary score $r \in \{0,1\}$ with r=0 if any contradiction is identified or support is insufficient, and r=1 otherwise. The binary reward is designed to prevent overfitting to noisy dense proxies and reduce reward hacking, thereby helping to preserve the model's original capabilities. Training proceeds with KL-regularized policy optimization against a reference model to constrain distributional drift and stabilize updates. The verifier also credits calibrated abstention when evidence is unavailable (e.g., deferral or uncertainty statements), encouraging deference under uncertainty rather than unsupported claims. This approach echoes recent advances on binary, verifiable rewards in math and code RL (Lambert et al., 2025; Shao et al., 2024).

We demonstrate the effectiveness of our approach through experiments on Qwen3 (Team, 2025) reasoning models (4B and 8B) across four factuality-centric evaluations and eight general capability benchmarks, spanning instruction following, knowledge retention, mathematical reasoning, and code generation. As shown in Figure 1 (right), online RL with our binary RAR improves the long-form factual precision from 38.1 to 62.5 (+24.4 pts), exceeding prior and concurrent approaches using DPO (46.6), RL with VeriScore (59.4), or RL with VeriScore and LLM judge (55.3; proposed in concurrent work by Chen et al. (2025)). Furthermore, binary RAR effectively avoids the utility regressions observed with dense rewards e.g., AlpacaEval drops from 54.7 to 42.2 under VeriScore, whereas Binary RAR maintains higher utility (AlpacaEval 53.9; IFEval 85.2). Similar trends hold on Qwen3-4B (e.g., $33.8 \rightarrow 62.3$ long-form with Binary RAR, surpassing VeriScore at 53.1). These results indicate that optimizing a binary, retrieval-verified signal yields larger hallucination reduction with fewer side effects on general capabilities compared to dense factuality rewards.

Our ablations and analysis show that (i) the KL constraint cleanly tunes the hallucination—utility trade-off—lower KL yields larger factuality gains but more utility pressure, while higher KL preserves capabilities; (ii) binary, contradiction-focused RAR is more robust than rating variants and dense factuality rewards, which are susceptible to reward misspecification and larger capability regressions; and (iii) retrieval-grounded binary feedback induces calibrated abstention ("I don't know" under insufficient evidence), reducing incorrect forced answers without dedicated abstention training. Together, these findings support continued online RL with retrieval-augmented binary reward as a principled route to improving factual reliability while maintaining general capabilities.

We summarize our contributions as follows:

- 1. We propose a new method that reduces extrinsic hallucination in the reasoning model Qwen3-8B by nearly 40%.
- We introduce an evaluation suite to assess the trade-off between hallucination and general capabilities.
- 3. We analyze why binary rewards are more effective than dense rewards in reducing hallucinations while preserving generation quality.

2 RELATED WORK

Measuring hallucinations in LM outputs Despite their impressive capabilities across diverse tasks, LMs are prone to hallucination—producing factually incorrect statements with unwarranted confidence—which undermines trust without careful oversight (Mallen et al., 2023). Quantifying hallucination in long-form generation is particularly challenging. Unlike short-form tasks, the open-

ended nature of long outputs prevents simple comparison against a single ground-truth reference. Instead, hallucination detection requires verifying responses against large knowledge sources. Recent work addresses this by evaluating factuality through verifiable claims extracted from model outputs. FactScore (Min et al., 2023) decomposes a response into atomic, verifiable claims and then checks each claim against reliable sources such as web search engines. The overall factuality score is typically the percentage of claims verified as correct. VeriScore (Song et al., 2024) extends this approach by extracting only verifiable claims, thereby generalizing the method to any type of response. This line of work connects closely to tool-augmented judges, where a language model's evaluation is enhanced by access to external tools (e.g., search engines) to better assess response quality Li et al. (2024). Building on these foundations, recent research has begun to reduce hallucination by directly optimizing models against such automatic evaluation metrics.

Reducing hallucination via post-training Various post-training approaches have been explored to mitigate extrinsic hallucination. The general goal is to reduce the chance of fabricating factually incorrect information in responses. Supervised fine-tuning (SFT) can improve factuality by avoiding training on knowledge that the model has not already assimilated during pre-training, as fine-tuning on unfamiliar knowledge can increase the propensity for hallucination (Newman et al., 2025; Zhang et al., 2024). Similarly, Direct Preference Optimization (DPO) trains the model to prefer more factual responses over less factual ones (Tian et al., 2024; Lin et al., 2024). This is often achieved by generating response pairs where the preference is determined by a continuous factuality assessment score. Concurrent with this work, Chen et al. (2025) combines offline learning (SFT, DPO) with online RL to enhance base LMs' factuality using a dense factuality signal (i.e., VeriScore). However, prior efforts largely emphasize factuality gains while offering limited assessment of impacts on other LM capabilities. We address this gap with an on-policy RL method that employs a *search-augmented binary reward*, improving the factuality of *fully trained* LMs *without* degrading general capabilities.

3 FACTUALITY-ORIENTED ONLINE RL WITH BINARY RETRIEVAL-AUGMENTED REWARD

Our goal is to directly optimize factual reliability while preserving general capabilities. Prior approaches that leverage SFT or DPO to improve factuality are typically *off-policy*: they train on fixed human/model outputs not produced by the current policy, risking distributional mismatch at inference. We instead adopt *online RL*, computing rewards on the model's *own rollouts*, and introduce a novel *binary retrieval-verified reward* (**Binary RAR**; Figure 1) that concentrates learning on falsity avoidance rather than proxy score maximization, with KL regularization to control drift.

This section outlines the training objective and algorithmic setup (§3.1), defines and motivates the binary reward with retrieval and verification (§3.2), and describes prompt curation for factuality-oriented RL (§3.3).

3.1 Preliminaries and Training Objective

The application of reinforcement learning (RL) to LMs frames the training process as an optimization problem. Given a prompt x, a language model π generates a response y by defining a policy $\pi_{\theta}(y \mid x)$. The goal is to train the policy to maximize a reward function r(x,y), which assigns a scalar score to the generated response. To prevent the finetuned model from deviating too far from its original capabilities, this optimization is typically constrained by a Kullback-Leibler (KL) divergence term against a reference model $\pi_{\rm ref}$. The objective is formally expressed as:

$$\max_{\pi_{\theta}} \mathbb{E}_{\substack{x \sim \mathcal{D} \\ y \sim \pi_{\theta}(\cdot \mid x)}} \left[r(x, y) - \beta \mathbb{D}_{\text{KL}}(\pi_{\theta}(\cdot \mid x) \| \pi_{\text{ref}}(\cdot \mid x)) \right] \tag{1}$$

where \mathcal{D} is the prompt dataset and β is a hyperparameter controlling the strength of the KL penalty.

Several algorithms exist to solve this objective. Among them, Group Relative Policy Optimization (GRPO) (Shao et al., 2024) has become a popular choice for LM post-training due to its stability

and computational efficiency. GRPO optimizes the following objective function:

$$\mathbb{E}_{\{y_i\}_{i=1}^n \sim \pi_{\text{act}}(\cdot \mid x)} \left[\frac{1}{n} \sum_{i=1}^n \frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \min\left(\text{IS}_{i,t} \cdot A_i, \text{clip}(\text{IS}_{i,t}, 1 - \epsilon, 1 + \epsilon) \cdot A_i \right) - \beta \mathbb{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}}) \right]$$

$$(2)$$

The key terms in this objective are the importance sampling (IS) ratio, $\text{IS}_{i,t} = \frac{\pi_{\theta}(y_i^t|y_i^{< t},x_i)}{\pi_{\text{act}}(y_i^t|y_i^{< t},x_i)}$, which corrects for the difference between the current policy π_{θ} and the actor policy π_{act} used for generation, and the advantage, $A_i = \frac{r(x,y_i) - \max_{j=1}^n (r(x,y_j))}{\text{std}_{j=1}^n (r(x,y_j))}$, which measures how much better a response y_i is compared to the batch average. Given its simplicity and stability, we adopt GRPO as the default RL algorithm for our experiments.

3.2 BINARY RETRIEVAL-AUGMENTED REWARD

We define the factual correctness of an instruction–response pair (x,y) as the absence of statements in y that contradict world knowledge, evaluated in the context of x. This notion targets factual consistency rather than instruction following: it does not judge task compliance, and it does not penalize assumptions explicitly prescribed by the prompt (e.g., fictional scenarios). We instantiate this criterion with a *binary retrieval-augmented reward* $r(x,y) \in \{0,1\}$, used for RL training; Figure 1 (left) outlines the overall procedure.

Retrieval A datastore \mathcal{DS} is a set of documents that are preprocessed, chunked, and indexed. The documents in the datastore is assumed to be factually correct. To verify the response against the datastore, for a given pair (x,y), A retriever R is used to retrieve a list of relevant documents denoted as $C(x,y) = \{d_1,\ldots,d_k\}$. These documents serve as evidence for checking the factuality of the response.

Verification To verify the correctness of instruction-response pairs (x, y), we adopt a language model judge that takes x, y, and C(x, y) as input, reasons about the correctness, and generates a binary score. We instruct the LM judge to only detect contradictions between the response and retrieved information, rather than requiring all information in the response to be supported by C(x, y).

Formally, given a prompt x and a response y, the reward is defined as:

$$r(x,y) = \begin{cases} 1 & \text{if no contradictions are found between } (x,y) \text{ and } C(x,y), \\ 0 & \text{otherwise.} \end{cases} \tag{3}$$

We optimize the KL-constrained RL objective (Equation 1) with our binary retrieval-augmented reward r(x, y). This approach avoids the complexity of designing a continuous reward function and instead provides an unambiguous or less noisy signal.

Practical efficiency for online RL Both retrieval and verification are computationally expensive, and evaluating r(x,y) can easily become the bottleneck of RL training. To address this, we propose several efficiency improvements. First, we adopt a pre-caching strategy for retrieval. For each prompt x in the training dataset \mathcal{D} , we pre-cache a set of relevant documents $\mathcal{DS}_{cache}(x)$ during dataset preparation. At training time, we then retrieve C(x,y) from this cached subset rather than from the full datastore \mathcal{DS} . Although we cannot anticipate exactly which information will appear in the model's output, we include documents related to both the prompt and the ground truth. This increases the likelihood that, if the model generates incorrect content, the retrieved documents will expose contradictions. Second, we optimize the verification stage. Instead of extracting and verifying many individual claims (as done in VeriScore), we detect contradictions by comparing the entire response against the retrieved documents in a single LM forward pass. This avoids repeatedly processing the same documents and significantly reduces computation, compared to concurrent work using VeriScore as factuality reward (Chen et al., 2025). We further increase throughput by launching multiple instances of the LMs for verification and evaluating batched responses in parallel.

3.3 PROMPT CURATION FOR FACTUALITY-ORIENTED RL

Curating high-quality and diverse sources of prompts is crucial for successful RL (Team et al., 2025). We collect training data based on two key considerations:

- 1. Responses must contain content that can be factually verified using external knowledge.
- 2. Instructions must be diverse and cover various scenarios and knowledge domains, rather than concentrating on a single task type.

Therefore, we source our data from WildChat (Zhao et al., 2024), a collection of instruction–response pairs from human interactions with OpenAI models. From this source, we select instances where the responses include verifiable factual content. To automate the selection, we use the OpenAI gpt-4.1 model with a detailed prompt to classify each instruction–response pair. For each selected instance, we query the Google Search API with the ground-truth response to retrieve URLs of potentially relevant web pages. We then crawl and parse the returned websites using a rule-based Python script. Each instance can be linked to up to 10 documents, determined by the Google Search API's limit. Instances with fewer than three retrieved documents are discarded, as sparse evidence may be insufficient for reliably assessing factual correctness. Importantly, our binary retrieval-augmented reward does not depend on ground-truth responses during RL training, which makes the approach scalable to larger datasets.

4 Experimental Setup

4.1 BENCHMARKING THE HALLUCINATION—UTILITY TRADE-OFF

We holistically evaluate methods for reducing the hallucination of LM generations while preserving general capability. To this end, we curate four factuality-centric benchmarks (hallucination evaluation) and six additional benchmarks spanning math, code, general chat, and instruction following (utility evaluation). Our objective is to minimize errors on hallucination evaluations while minimizing any performance drop on utility evaluations relative to the original LM.

Hallucination Evaluation We assess hallucination behaviors in both long-form generation and short-form question answering using the following datasets: Biography (Min et al., 2023), Wild-Hallucination (Zhao et al., 2024) for long-form generations, and PopQA (Mallen et al., 2023) and GPQA (Rein et al., 2024) for short-form QA focusing on long-tail or complex scientific-domain knowledge. For long-form generation, we compute *factual precision* using FACTSCORE (Min et al., 2023), which extracts atomic claims from model outputs and verifies them against reference documents. Factual precision is defined as the percentage of correct claims. For short-form QA, we explicitly instruct the model to answer with "I don't know" when uncertain. We measure the *hallucination rate* as the percentage of cases where the model produces an incorrect answer. Evaluation is conducted on PopQA and GPQA, with correctness determined by a strong LLM judge and multiple-choice accuracy, respectively.

Utility Evaluation We then evaluate whether models retain general capabilities after continued fine-tuning. For knowledge retention, we revisit **PopQA** and **GPQA** with a no-abstention prompting setup: instead of allowing the model to abstain, we require it to always provide an answer (i.e., make its best guess). Accuracy is then measured against the ground-truth answers, using the same correctness-judging methods as in hallucination evaluation. Beyond knowledge retention, we assess broader capabilities on six additional benchmarks:**AlpacaEval** (Dubois et al., 2024) and **IFEval** (Zhou et al., 2023) for instruction following; **GSM8K** (Cobbe et al., 2021) and **Minerva** (Lewkowycz et al., 2022) for mathematical reasoning; and **HumanEval** (Chen et al., 2021) and **MBPP** (Austin et al., 2021) for code generation. We follow each benchmark's official evaluation protocol to quantify any regressions in utility after fine-tuning. Full details are provided in the Appendix.

4.2 BASELINES

We compare our method with diverse non-RL and RL methods using various rewards. We first apply supervised fine-tuning (**SFT**) and direct preference optimization (**DPO**) to the base models (Tian

et al., 2024; Lin et al., 2024; Chen et al., 2025). For each base reasoning model, we generate 8 responses and use the VeriScore pipeline to assess the factuality of each response. ¹ Specifically, we extract verifiable claims from the model responses, verify their correctness against pre-cached documents, and compute the percentage of correct claims. We then apply SFT to the language model using the response with the highest factuality score for each instance. DPO is used to train the model to contrast pairs of responses with the largest difference in factuality scores, subject to the constraint that their length difference is less than 10%. This data construction strategy is designed to avoid length-based optimization, or "length hacking" (Chen et al., 2025).

We also evaluate RL methods with different reward signals. First, we use **LM Judge**, which takes a reference response and a model response as input and provides a rating of overall output quality on a scale from 0 to 10, following common practice (Gunjal et al., 2025). Note that LM Judge optimizes general quality rather than factuality specifically. We also test **VeriScore** as an RL reward, which has been proposed in concurrent work (Chen et al., 2025). For computing VeriScore, we use BM25 as the retriever and split documents into 256-token chunks (tokenized with the Qwen3 tokenizer). For each claim, the top 4 retrieved chunks are used for verification. Both claim extraction and verification are performed with Qwen-32B.

4.3 Training Details

We apply continual RL fine-tuning based on the Qwen3-8B and Qwen3-4B models, a popular language model family with reasoning capability. We employ GRPO as our main RL algorithm. We deploy Qwen3-32B as our judge model for computing binary retrieval-augmented rewards. The judge model is prompted to identify contradictions between model responses and the relevant documents. We use a learning rate of 1×10^{-6} and set the KL coefficient to 1×10^{-3} for Qwen3-8B and 3×10^{-3} for Qwen3-4B. To compute binary RAR, we apply BM25 for retrieval, chunking documents into 512 tokens using the Qwen3 tokenizer. For each response, the top 8 chunks are retrieved and then verified by Qwen3-32B. We use a early stoping strategy to avoid overtraining leading to degradation of utility. To ensure the model's utility is preserved, we stop training if a checkpoint shows more than a 10% drop on any utility benchmark.

5 MAIN RESULTS

5.1 RESULTS ON HALLUCINATION REDUCTION

Table 1 reports main results across long-form generation and short-form question answering (QA). The base QWEN3-8B exhibits substantial factuality limitations, achieving only 38.1% factual precision on long-form generation and 35.9% accuracy on short-form QA. The relatively low scores of QWEN3-4B align with prior evidence that smaller models have limited parametric recall/memorization (Mallen et al., 2023). RL with Binary RAR yields the largest improvements among all approaches, outperforming SFT, DPO-based methods, and other RL variants.

SFT and DPO: limited effectiveness for hallucination reduction. Approaches that apply SFT or DPO to high-VeriScore responses provide only modest gains in long-form factual precision and yield little to no improvement on short-form QA with abstention enabled.

Binary RAR achieves substantially better results than other RL rewards. Online RL with VERISCORE delivers notable factuality gains (59.4 long-form; 41.5 short-form) over SFT- or DPO-based approaches, highlighting the effectiveness of online RL guided by a factuality-focused reward. In contrast, RL with an LM-judge reward degrades long-form performance (34.7), suggesting that generic instruction-following rewards can conflict with factual-accuracy objectives. On long-form generation, Binary RAR raises factual precision from 38.1 to 62.5 (+24.4 pts) on QWEN3-8B, substantially improving the VeriScore-based approach. On short-form QA, Binary RAR increases accuracy from 35.9% to 50.6%, the strongest result among all methods. On QWEN3-4B, we observe similar improvements in long-form factual precision—from 33.8% to 62.3% (+28.5 pts)—substantially surpassing RL with VERISCORE (53.1%).

¹We do not experiment with SFT or DPO using binary RAR because responses to many prompts are always zero or one, which makes generation inefficient.

	Long-fo	orm (Factual Preci	$\boldsymbol{Short\text{-}form}(\textbf{Hallucination}\textbf{Rate}\downarrow)$				
Models	Biography	WildHallu	AVG	PopQA	GPQA	AVG	
Qwen3-8B	23.8	52.4	38.1	71.2	50.0	60.6	
+ SFT	24.7	53.5	39.1	70.4	50.0	60.2	
+ DPO	33.1	60.2	46.6	65.2	49.1	57.2	
+ RL (LM Judge)	19.6	49.7	34.7	68.8	48.0	58.4	
+ RL (VeriScore)	48.3	70.5	59.4	43.6	41.1	42.3	
+ RL (Binary RAR)	54.2	70.8	62.5	26.8	28.3	27.6	
Qwen3-4B	18.1	49.5	33.8	82.2	55.1	68.7	
+ SFT	21.1	51.3	36.2	83.8	54.7	69.2	
+ DPO	26.6	56.1	41.3	82.6	54.5	68.5	
+ RL (LM Judge)	17.4	46.3	31.9	80.4	54.0	67.2	
+ RL (VeriScore)	38.9	67.4	53.1	73.0	51.3	62.2	
+ RL (Binary RAR)	53.5	71.1	62.3	46.6	37.3	41.9	

Table 1: Factuality results comparing different training methods on long-form generation and short-form QA tasks. We report FactScore precision for long-form generation and hallucination rate for short-form question answering. Binary RAR achieves the best hallucination reduction, showing the highest factual precision and the lowest hallucination rate in short-form question answering.

	Instruction Following		Knowledge		Math		Coding		AVG
Models	AlpacaEval	IFEval	PopQA	GPQA	GSM8K	Minerva	HumanEval	MBPP	
Qwen3-8B	54.7	87.2	20.2	48.2	92.8	80.7	83.5	67.4	66.8
+ SFT	55.7	86.9	20.4	47.9	91.6	82.0	83.8	67.0	66.9
+ DPO	53.0	84.5	18.6	47.5	90.8	82.1	86.7	67.8	66.4
+ RL (LM Judge)	55.0	82.2	19.2	52.2	88.1	77.7	83.8	66.3	65.6
+ RL (VeriScore)	42.2	88.7	19.6	47.7	92.2	79.0	83.4	66.9	65.0
+ RL (Binary RAR)	53.9	85.2	20.6	48.8	93.4	82.3	86.1	67.6	67.2
Qwen3-4B	41.7	86.1	16.4	44.2	91.1	82.8	85.5	65.7	64.2
+ SFT	41.2	82.6	15.2	43.5	91.4	83.6	83.2	65.6	63.3
+ DPO	39.6	81.9	15.8	44.0	90.1	82.7	85.8	66.3	63.3
+ RL (LM Judge)	42.3	74.3	16.0	43.5	87.0	82.1	85.9	66.2	62.2
+ RL (VeriScore)	38.4	86.0	15.4	40.8	90.8	82.5	84.5	66.2	63.1
+ RL (Binary RAR)	43.0	84.7	16.4	42.6	90.7	83.8	84.6	65.0	63.9

Table 2: General capability results across eight utility benchmarks covering instruction following (AlpacaEval, IFEval), knowledge (PopQA, GPQA), math reasoning (GSM8K, Minerva), and coding generation (HumanEval, MBPP). We color the cell based on the degree of performance degradation.

Models learn abstention behavior. An important finding is that Binary RAR enables models to spontaneously learn when to abstain from answering. On short-form QA, the base model rarely abstains despite high uncertainty, leading to frequent incorrect responses. After Binary RAR training, models strategically abstain on questions they would otherwise answer incorrectly, improving reliability without explicit abstention training. In Section 6.2, we present a detailed analysis of strategic abstention behaviors.

5.2 RESULTS ON GENERAL CAPABILITIES PRESERVATION

Table 2 shows performance across eight utility benchmarks covering instruction following, mathematics, and coding. Our key finding is that Binary RAR best preserves the model's original capabilities while improving factuality.

Binary RAR maintains utility better than alternatives. The Binary RAR model achieves an average score of 78.1, without any performance deterioration from the base model's 77.7 performance. In contrast, RL with VeriScore shows noticeable degradation (75.4 average), particularly on AlpacaEval, where performance drops from 54.7 to 42.2.

Impacts of factuality-oriented training on each task. The marked AlpacaEval degradation under VeriScore indicates that continuous rewards are more vulnerable to reward hacking—overoptimizing the proxy at the expense of general utility. We analyze this phenomenon further in Section A. When forced to provide answers (no-abstention setting) on short-form question answering

Models	KL Coef.	Factuality			Utility			
		Biography	WildHallu	AVG	AlpacaEval	IFEval	AVG	
Qwen3-8B	1	23.8	52.4	38.1	54.7	87.2	71.0	
+ VeriScore	0.003	40.7	65.0	52.8	45.0	86.9	65.9	
+ VeriScore	0.001	48.3	70.5	59.4	42.2	88.7	65.5	
+ VeriScore + LM Judge	0.001	42.3	68.4	55.3	47.7	85.6	66.6	
+ Binary RAR	0.001	54.2	70.8	62.5	53.9	85.2	69.6	
+ Binary RAR	0.0003	58.2	65.0	61.6	49.1	88.2	68.6	
+ Rating RAR	0.001	35.0	66.9	51.0	46.9	85.8	66.3	

Table 3: Ablation study results examining the effects of KL regularization coefficient, reward composition, and reward formulation on factuality and utility performance.

datasets, RL with binary RAR remains largely unchanged after training (PopQA: $20.2\% \rightarrow 20.6\%$; GPQA: $48.2 \rightarrow 48.8\%$). This demonstrates that Binary RAR training enhances factuality without degrading the model's fundamental knowledge or reasoning abilities. In contrast, all training methods have minimal impact on mathematical reasoning and coding abilities. This is likely because our factuality-focused training primarily affects knowledge recall tasks, while math and coding performance depend more on reasoning over the provided input context. Interestingly, Binary RAR shows slight improvements on Minerva (+1.6), possibly due to the training data containing some reasoning and coding examples.

6 ANALYSIS

6.1 ABLATION STUDIES

Table 3 presents systematic ablations across three key dimensions: regularization strength, reward composition, and reward formulation.

KL regularization is crucial for balancing objectives. All the main experiments on Qwen3-8B are conducted with a KL coefficient of 0.001. Therefore, we try increasing KL for VeriScore to see if it helps preserve utility, and decreasing KL for binary RAR to see if it hurts utility. We find that for both rewards, a higher KL indeed better preserves utility but limits improvements in the factuality score. Nevertheless, VeriScore still shows a large drop in AlpacaEval with KL 0.003 and a significant reduction in factuality improvement.

Mixed rewards can mitigate utility degradation. Combining VeriScore with LM Judge rewards, which is recently proposed by concurrent work by Chen et al. (2025), partially addresses the AlpacaEval degradation issue, improving the score from 42.2 to 47.7 while maintaining most factuality gains. However, this approach still shows modest utility drops compared to Binary RAR.

Binary rewards outperform dense alternatives. Comparing Binary RAR with Rating RAR (a rating-based version providing 0–10 scores instead of 0–1) reveals that the binary formulation is more effective. Both its effectiveness in improving factuality and in preserving utility are lower than that of the binary formulation, supporting our hypothesis that dense rewards are more susceptible to reward hacking. We conduct a qualitative analysis in Section A to further investigate this issue.

6.2 Analysis of Abstention Behavior

Figure 2 illustrates how Binary RAR training fundamentally alters the model's strategy for handling uncertain questions. We evaluate performance in two scenarios: one allowing abstention ("I don't know" responses) used in hallucination evaluation, and another requiring forced responses used in utility evaluation.

The base Qwen3-8B model exhibits high error rates and rarely abstains, even on questions where it lacks sufficient knowledge. After Binary RAR training, the model's behavior changes dramatically: it abstains on 55.2% of PopQA questions and 27.5% of GPQA questions. While the overall accuracy decreases slightly (with a less than 15% relative reduction). Importantly, these abstentions are not random. The model primarily abstains on questions it would otherwise answer incorrectly. For

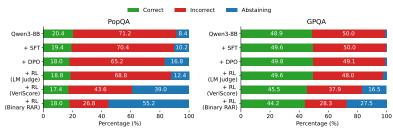


Figure 2: Breakdown of Qwen3-8B responses on short-form question answering, comparing correct answers, incorrect answers, and abstentions.

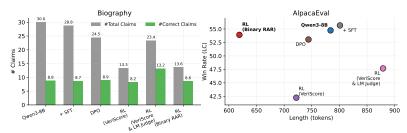


Figure 3: Analysis of informativeness in model outputs. Left: number of total and correct claims in responses to biography datasets. Right: AlpacaEval length-controlled win rate versus output length.

attempted questions, accuracy increases from 22.3% to 40.2% on PopQA and from 49.4% to 60.9% on GPQA. This indicates that the model strategically chooses to abstain when uncertain, rather than refusing answers arbitrarily.

6.3 IMPACT ON QUALITY AND FACTUALITY IN LONG-FORM GENERATION

To understand whether improved factuality comes at the cost of informativeness, we analyze the information content in model outputs. Figure 3 shows the relationship between training method and the number of atomic claims generated. Binary RAR training does not significantly reduce the number of correct claims in model outputs, suggesting that factuality improvements are not simply due to generating less informative responses. On long-form generation tasks, the model maintains a similar level of correct detail compared to the base model while substantially reducing incorrect claims. As a result, although on AlpacaEval, Binary RAR produces slightly shorter responses than the base model, it maintains both overall win rate (59.3 \rightarrow 59.2) and length-controlled win rate (54.7 \rightarrow 53.9). This further confirms that the model learns to avoid unnecessary or uncertain information while preserving response quality, rather than simply generating shorter outputs. By contrast, although RL with VERISCORE + LM-judge increases the number of correct claims largely by substantially lengthening outputs, it shows a notable drop of AlpacaEval win rate. A closer inspection shows the model often defaults to vague, high-level statements rather than precise, source-grounded facts. In the next section, we present a qualitative analysis of responses produced under different RL objectives.

7 CONCLUSION

We present a reinforcement learning fine-tuning approach using a binary retrieval-augmented reward (RAR) to mitigate hallucinations in large language models. By verifying outputs against retrieved evidence and assigning a simple binary score, binary RAR proves more effective than SFT, DPO, or RL with dense rewards such as VeriScore. RL with binary RAR enables models to reduce factual errors in long-form generation, abstain when uncertain in short-form question answering, and at the same time retain knowledge memorization, maintain informativeness, and preserve general capabilities. These results demonstrate that simple binary rewards offer a practical, robust, and scalable path toward safer and more reliable language models.

ETHICS STATEMENT

This research aims to mitigate extrinsic hallucinations in language models, which is crucial for developing safer and more reliable AI systems that users can trust. By improving the factual accuracy of model outputs, this work helps reduce the potential for spreading misinformation. The methods employed use publicly available data and focus on enhancing factual correctness without intentionally introducing new societal biases or risks.

REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our results, we provide a detailed description of the training objective, algorithms, and binary reward mechanism in Section 3. The experimental setup, including all benchmarks for factuality and utility, is fully outlined in Section 4 and the Appendix. Furthermore, specific training details such as the base models, learning rates, and KL coefficients are explicitly stated in Section 4.3.

REFERENCES

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021. URL https://arxiv.org/abs/2108.07732.
- Aaron Chatterji, Thomas Cunningham, David J Deming, Zoe Hitzig, Christopher Ong, Carl Yan Shan, and Kevin Wadman. How people use chatgpt. Working Paper 34255, National Bureau of Economic Research, September 2025. URL http://www.nber.org/papers/w34255.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL https://arxiv.org/abs/2107.03374.
- Xilun Chen, Ilia Kulikov, Vincent-Pierre Berges, Barlas Oğuz, Rulin Shao, Gargi Ghosh, Jason Weston, and Wen tau Yih. Learning to reason for factuality, 2025. URL https://arxiv.org/abs/2508.05618.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL https://arxiv.org/abs/2110.14168.
- Yann Dubois, Percy Liang, and Tatsunori Hashimoto. Length-controlled alpacaeval: A simple debiasing of automatic evaluators. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=CybBmzWBX0.
- Anisha Gunjal, Anthony Wang, Elaine Lau, Vaskar Nath, Bing Liu, and Sean Hendryx. Rubrics as rewards: Reinforcement learning beyond verifiable domains. *arXiv preprint arXiv:2507.17746*, 2025.
- Adam Tauman Kalai, Ofir Nachum, Santosh S. Vempala, and Edwin Zhang. Why language models hallucinate, 2025. URL https://arxiv.org/abs/2509.04664.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James Validad Miranda, Alisa Liu, Nouha Dziri, Xinxi Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Christopher Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training. In Second Conference on Language Modeling, 2025. URL https://openreview.net/forum?id=iluGbfHHpH.
- Aitor Lewkowycz, Anders Johan Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Venkatesh Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=IFXTZERXdM7.
- Lei Li, Yekun Chai, Shuohuan Wang, Yu Sun, Hao Tian, Ningyu Zhang, and Hua Wu. Toolaugmented reward modeling. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=d94x0gWTUX.
- Sheng-Chieh Lin, Luyu Gao, Barlas Oguz, Wenhan Xiong, Jimmy Lin, Wen tau Yih, and Xilun Chen. FLAME: Factuality-aware alignment for large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=zWuHSIALBh.

- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9802–9822, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.546. URL https://aclanthology.org/2023.acl-long.546/.
- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12076–12100, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.741. URL https://aclanthology.org/2023.emnlp-main.741/.
- Benjamin Newman, Abhilasha Ravichander, Jaehun Jung, Rui Xin, Hamish Ivison, Yegor Kuznetsov, Pang Wei Koh, and Yejin Choi. The curious case of factuality finetuning: Models' internal beliefs can improve factuality, 2025. URL https://arxiv.org/abs/2507.08371.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=Ti67584b98.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402.03300.
- Linxin Song, Taiwei Shi, and Jieyu Zhao. The hallucination tax of reinforcement finetuning, 2025. URL https://arxiv.org/abs/2505.13988.
- Yixiao Song, Yekyung Kim, and Mohit Iyyer. VeriScore: Evaluating the factuality of verifiable claims in long-form text generation. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), Findings of the Association for Computational Linguistics: EMNLP 2024, pp. 9447–9474, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.552. URL https://aclanthology.org/2024.findings-emnlp.552/.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- Qwen Team. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.
- Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. Fine-tuning language models for factuality. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=WPZ2yPag4K.
- Zijun Yao, Yantao Liu, Yanxu Chen, Jianhui Chen, Junfeng Fang, Lei Hou, Juanzi Li, and Tat-Seng Chua. Are reasoning models more prone to hallucination?, 2025. URL https://arxiv.org/abs/2505.23646.
- 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 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.394. URL https://aclanthology.org/2024.naacl-long.394/.

Wenting Zhao, Tanya Goyal, Yu Ying Chiu, Liwei Jiang, Benjamin Newman, Abhilasha Ravichander, Khyathi Chandu, Ronan Le Bras, Claire Cardie, Yuntian Deng, and Yejin Choi. Wildhallucinations: Evaluating long-form factuality in Ilms with real-world entity queries, 2024. URL https://arxiv.org/abs/2407.17468.

Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL https://arxiv.org/abs/2311.07911.

Figure 4: Comparing Binary RAR, VeriScore, and LM Judge on LM responses.

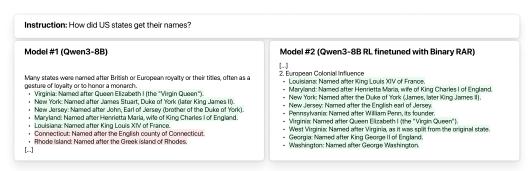


Figure 5: Comparing the response of models before and after training.

A QUALITATIVE ANALYSIS

To better understand the impact of RL training with Binary RAR, VeriScore, and the LM Judge, we present a qualitative analysis of the reward signals and the fine-tuned models.

LM Judge alone provides limited factuality assessment. Figure 4 shows two responses to the same instruction, along with their evaluations from three rewards. The first response contains a factual error, while the second is entirely correct. All three rewards assign a lower score to the erroneous response than to the correct one. However, the LM Judge assigns only 0.7 even to the fully correct response. The generated reason reveals that the LM reward prioritizes detailed elaboration, which reduces its sensitivity to factual accuracy.

Binary RAR reduces hallucination while preserving detail. Figure 5 compares outputs from Qwen3-8B before and after RL fine-tuning with Binary RAR. The base model generates incorrect information about Connecticut" and Rhode Island," whereas the fine-tuned model avoids these errors and even adds additional examples of states named after royalty. This demonstrates that RL fine-tuning with Binary RAR reduces factual errors while retaining informative detail in the responses. More examples are provided in the Appendix.

B EVALUATION DETAILS

We assess hallucination in both long-form generation and short-form question answering using the following benchmarks:

- **Biography** (Min et al., 2023): A benchmark consisting of prompts that ask models to write biographies of specific individuals.
- WildHallucination (Zhao et al., 2024): A dataset probing factual consistency across diverse real-world entities, including people, geography, and computing, with emphasis on rare entities.
- PopQA (Mallen et al., 2023): A short-form QA dataset covering entities of varying popularity; correctness is judged automatically by a strong evaluator.

• **GPQA** (Rein et al., 2024): A multiple-choice QA dataset covering graduate-level biology, chemistry, and physics, where questions and answers are expert-authored.

To measure whether factuality improvements cause regressions in other areas, we evaluate general capabilities using these benchmarks:

- **AlpacaEval** (Dubois et al., 2024): An automatic evaluation framework that compares outputs pairwise against a baseline using an LLM judge gpt-4.1. We use the length-controlled win rate metric to reduce length bias.
- **IFEval** (Zhou et al., 2023): A benchmark of 500 prompts covering 25 types of verifiable instructions, designed to test instruction fidelity with objectively checkable outcomes.
- **GSM8K** (Cobbe et al., 2021): A dataset of 8.5K grade-school math word problems requiring multi-step reasoning (typically 2–8 steps).
- Minerva (Lewkowycz et al., 2022): A collection of 272 graduate-level quantitative reasoning problems in STEM fields such as physics and chemistry, requiring domain-specific expertise.
- HumanEval (Chen et al., 2021): A benchmark of 164 handwritten Python programming tasks.
- MBPP (Austin et al., 2021): The Mostly Basic Programming Problems dataset, consisting of 974 entry-level Python tasks.

C TRAINING DETAILS

RL Fine-tuning For reinforcement learning, we trained models for a maximum of four epochs. We used a batch size of 16 unique prompts, generating 8 rollouts per prompt. By default, training ran for 2,000 gradient steps; however, for the VeriScore reward, we applied an early stopping constraint at 1,000 steps to prevent degradation in utility benchmark scores.

SFT and DPO Baselines For the non-RL baselines, we found that one epoch of supervised fine-tuning (SFT) and four epochs of direct preference optimization (DPO) yielded the best results.

Training Prompts The prompts used to train the reward models and conduct fine-tuning are detailed in the figures below.

D LARGE LANGUAGE MODEL USAGE

We use language models to assist with literature reviews, help write experiment code, and check grammar in paper writing. All ideas, results, and content are produced by the authors, and all outputs from these tools are verified by the authors.

810 **Binary Retrieval-augmented Reward** 811 812 You are a professional fact-checker tasked with evaluating a language model's response for accuracy. 813 Your job is to systematically compare the response against the provided web search results to identify 814 any factual errors or contradictions. A response deserves a full score if every claim it makes is factually correct, even if it doesn't include every available detail. Omitting information is not a factual 815 error. If a response contains anything that contradicts the world knowledge shown in the search 816 results, it should be considered an error and get a low score. 817 818 ## EVALUATION PROCESS: 819 1. Read the search results thoroughly to understand the factual baseline 2. Examine each factual claim in the language model's response 820 3. Cross-reference each claim against the search results 821 4. Classify each discrepancy according to the guidelines below 823 ## CONSIDER CORRECT (No Factual Error): - Paraphrasing: Same facts expressed in different words 824 - Reasonable inferences: Logical conclusions drawn from search results - Partial information: Incomplete but accurate subsets of available information - Contextual additions: Background information that doesn't contradict search results - Minor formatting differences: Different ways of presenting same data 828 829 ## CONSIDER INCORRECT (Factual Error): - Direct contradictions: Response states opposite of what search results show 830 - Numerical errors: Wrong dates, statistics, quantities, percentages 831 - Categorical errors: Wrong classifications, locations, names, titles 832 - Causal errors: Incorrect cause-and-effect relationships 833 - Timeline errors: Events placed in wrong chronological order 834 - Attribution errors: Wrong sources, authors, or speakers cited 835 ## SCORING RULES: 836 Your evaluation will result in a binary score: 0 or 1. 837 - SCORE 1 (No Contradiction): Assign this score if: 838 1. The response is fully supported by the document. 2. The response contains information that is NOT in the document, but DOES NOT contradict it. 839 - SCORE 0 (Contradiction): Assign this score ONLY if you find a clear, factual contradiction between 840 the response and the supporting document. A contradiction occurs when the response states the 841 opposite of what the document says (e.g., wrong dates, names, events, or outcomes). 842 843 ## OUTPUT FORMAT: You must respond ONLY in a valid JSON format with exactly these two fields: 844 - "REASONING": A brief explanation for your score. 845 - For SCORE 0, specify the contradiction (e.g., "The response states the event was in 2022, but the 846 document says it was in 2023."). 847 - For SCORE 1, simply state "No contradiction found." 848 - "SCORE": An integer, either 0 or 1. 849 >>> Begin of search results <<< 850 {passages_text} 851 << End of search results >>> 852 853 >>> Begin of the prompt <<< 854 {prompt_text} << End of the prompt >>> 855 >>> Begin of the response <<< {response_text} 858 << End of the response >>> 859 Respond in JSON format. {{"REASONING": "[...]", "SCORE": "<your-score>"}} 861

Figure 6: Prompt for binary fact-checking.

862

915

916

865 866 Rating Retrieval-augmented Reward 867 You are a professional fact-checker tasked with evaluating a language model's response for accuracy. 868 Your job is to systematically compare the response against the provided web search results to identify any factual errors or contradictions. A response deserves a full score if every claim it makes is 870 factually correct, even if it doesn't include every available detail. Omitting information is not a factual 871 error. If a response contains anything that contradicts the world knowledge shown in the search 872 results, it should be considered an error and get a low score. 873 ## EVALUATION PROCESS: 874 1. Read the search results thoroughly to understand the factual baseline 875 2. Examine each factual claim in the language model's response 3. Cross-reference each claim against the search results 877 4. Classify each discrepancy according to the guidelines below 878 ## CONSIDER CORRECT (No Factual Error): 879 - Paraphrasing: Same facts expressed in different words - Reasonable inferences: Logical conclusions drawn from search results - Partial information: Incomplete but accurate subsets of available information 882 - Contextual additions: Background information that doesn't contradict search results 883 - Minor formatting differences: Different ways of presenting same data 884 ## CONSIDER INCORRECT (Factual Error): 885 - Direct contradictions: Response states opposite of what search results show - Numerical errors: Wrong dates, statistics, quantities, percentages 887 - Categorical errors: Wrong classifications, locations, names, titles - Causal errors: Incorrect cause-and-effect relationships - Timeline errors: Events placed in wrong chronological order 889 - Attribution errors: Wrong sources, authors, or speakers cited 890 891 ## CONFIDENCE SCORING GUIDE: 892 - 0-2: Very confident there is a factual error (multiple clear contradictions) - 3-4: Moderately confident there is a factual error (one clear contradiction) 893 - 5: Uncertain (ambiguous evidence or unclear from search results) 894 - 6-7: Moderately confident there is no factual error (mostly accurate with minor concerns) 895 - 8-10: Very confident there is no factual error (all stated facts are accurate, regardless of completeness) 897 ## OUTPUT FORMAT REQUIREMENTS: Respond ONLY in valid JSON format with exactly these two fields: - "REASONING": A concise explanation of your assessment (1-2 sentences max, e.g., "the response 899 states ... but the search results show ... so there is a factual error" or "no factual error found") 900 - "SCORE": An integer from 0-10 representing your confidence level 901 902 >>> Begin of search results <<< 903 {passages_text} << End of search results >>> 904 905 >>> Begin of the prompt <<< 906 {prompt_text} 907 << End of the prompt >>> 908 >>> Begin of the response <<< 909 {response_text} 910 << End of the response >>> 911 912 Respond in JSON format. {{"REASONING": "[...]", "SCORE": "<your-score>"}} 913 914

Figure 7: Prompt for rating-based fact-checking.

Claim Extraction for VeriScore Training / FactScore Evaluation

Extract as many fine-grained, atomic, and verifiable factual claims as possible from the response. Each claim should be a single piece of information that could be looked up in a database, official documentation, reputable forum, or reliable source such as Wikipedia or scientific literature.

Guidelines for atomic claims:

- Split a sentence that joins different facts using "and," "or," or by listing into multiple claims.
- If a claim could be split into multiple smaller, independent statements, do so.
- Replace pronouns (e.g., "he", "she", "it", "they") with the full entity name explicitly stated in the response. If the entity name is not explicitly mentioned, leave the pronoun unchanged.
- Extract claims EXACTLY as stated, even if the information appears incorrect or false.

Include as claims:

- Statements about the existence, property, function, or relationship of entities, organizations, concepts, or technologies.
- Claims about names, definitions, features, purposes, or histories.
- Statements about what something does, who runs it, what it is used for, or what it affects.
- For hedged language ("may be," "might be," "could be"), extract the factual association, typical usage, or commonly reported function as long as the claim is traceable to community consensus, documentation, or reputable user reports.
- If a quotation is present, extract it verbatim with the source if given.
- Claims must stand alone, using names or clear descriptions, not pronouns.

Do not include as claims:

- Personal opinions, suggestions, advice, instructions, or experiences.
- Pure speculation or possibilities that are not reported in any documentation or user discussions.
- Claims from code blocks or pure math derivations.

Extract claims only from the response section, not from the prompt or question. If the response does not contain any verifiable factual claims, output an empty list.

Output a JSON list of strings. Each string should be a single atomic factual claim from the response, clearly stated and verifiable.

```
>>> Begin of prompt <<<
{prompt_text}
<<< End of prompt >>>
>>> Begin of response <<<
{response_text}
<<< End of response >>>
```

Facts (as a JSON list of strings):

Figure 8: Prompt for atomic claim extraction.

Claim Verification for VeriScore Training / FactScore Evaluation You need to judge whether a claim is supported or contradicted by Google search results, or whether there is no enough information to make the judgement. When doing the task, take into consideration whether the link of the search result is of a trustworthy source. Below are the definitions of the three categories: Supported: A claim is supported by the search results if everything in the claim is supported and nothing is contradicted by the search results. There can be some search results that are not fully related to the claim. Contradicted: A claim is contradicted by the search results if something in the claim is contradicted by some search results. There should be no search result that supports the same part. Inconclusive: A claim is inconclusive based on the search results if: - a part of a claim cannot be verified by the search results, - a part of a claim is supported and contradicted by different pieces of evidence, - the entity/person mentioned in the claim has no clear referent (e.g., "the approach", "Emily", "a book"). >>> Begin of search results <<< {passages_text} << End of search results >>> Claim: {claim_text} Task: Given the search results above, is the claim supported, contradicted, or inconclusive? Your answer should be either "supported", "contradicted", or "inconclusive" without explanation and comments. Your decision:

Figure 9: Prompt for claim verification.

1026 1027 1028 1029 1030 **Dataset Curation** 1031 You need to judge whether a claim is supported or contradicted by Google search results, or whether 1032 there is no enough information to make the judgement. When doing the task, take into consideration 1033 whether the link of the search result is of a trustworthy source. 1034 1035 Below are the definitions of the three categories: 1036 Supported: A claim is supported by the search results if everything in the claim is supported 1037 and nothing is contradicted by the search results. There can be some search results that are not fully related to the claim. 1039 Contradicted: A claim is contradicted by the search results if something in the claim is contradicted by some search results. There should be no search result that supports the same part. 1040 Inconclusive: A claim is inconclusive based on the search results if: 1041 - a part of a claim cannot be verified by the search results, 1042 - a part of a claim is supported and contradicted by different pieces of evidence, 1043 - the entity/person mentioned in the claim has no clear referent (e.g., "the approach", "Emily", "a book"). 1045 >>> Begin of search results <<< 1046 {passages_text} 1047 <<< End of search results >>> 1048 1049 Claim: {claim_text} 1050 Task: Given the search results above, is the claim supported, contradicted, or inconclusive? Your answer should be either "supported", "contradicted", or "inconclusive" without explanation and 1051 comments. 1052 1053 Your decision: 1054 1055 1056 Figure 10: Prompt for dataset curation. 1058 1062 1063 1064 1065 1066 Instruction: I'm trying to teach myself to have nicer handwriting. Can you help? 1067 1068

Model #1 (Training with Binary RAR)

Improving your handwriting is a rewarding journey that requires patience, practice, and a structured approach. Here's a step-by-step guide to help you get started, tailored to your self-teaching needs 1. Assess Your Current Handwriting

- Self-Evaluation: Write a short paragraph or sample text. Note what you like and dislike (e.g., size, spacing, slant, consistency).
- Compare to Models: Use examples of good handwriting (print or cursive) to identify differences. Look for elements like:
- Spacing: Are letters and words evenly spaced?

[...]

1069

1070

1071

1074

1078 1079

Model #2 (Training with VeriScore)

Improving your handwriting is a personal journey that can be both rewarding and enjoyable. Here are some general steps and considerations to help you on you

- 1. Assess Your Current Handwriting: Start by reflecting on your current handwriting. Are there specific aspects you'd like to improve, such as legibili consistency, or spacing? Understanding your goals can help you focus your
- 2. Understand the Basics: Handwriting involves more than just writing letters. It includes the formation of each letter, spacing between letters and words, and the overall flow of writing. Understanding these elements can help you identify [...]

Figure 11: Comparing responses from models RL finetuned with Binary RAR and VeriScore.