

Reinforced Informativeness Optimization for Long-Form Retrieval-Augmented Generation

Anonymous ACL submission

Abstract

Long-form question answering (LFQA) requires open-ended long-form responses that synthesize coherent, factually grounded content from multi-source evidence. This makes reinforcement learning (RL) reward design critical. The reward must be verifiable for faithful grounding and stable optimization. However, many standard rewards assume a unique target with an exact-match notion of correctness, which fits short-form QA and math but breaks in LFQA. As a result, current RAG systems still lack verifiable reward mechanisms, yielding unstable feedback signals and sub-optimal optimization outcomes. We propose RioRAG, a framework for reinforced verifiable informativeness optimization. First, it defines informativeness as a measurable and externally verifiable objective for RL. Second, RioRAG uses nugget-centric verification with cross-source checks to enable self-evolution of smaller LLMs and to provide denser, action-discriminative rewards that mitigate reward sparsity and stabilize optimization. This formulation avoids handcrafted supervision for the policy model and strong teacher-model distillation, relying instead on externally verifiable feedback. Experiments on LongFact and RAGChecker show that RioRAG achieves higher factual recall and faithfulness, establishing verifiable reward modeling as a foundation for trustworthy long-form RAG. Our codes are available at <https://anonymous.4open.science/r/RioRAG/>.

1 Introduction

Long-form question answering (LFQA) represents a crucial step toward enabling AI systems to deliver comprehensive and factually reliable responses by generating elaborate and multi-sentence answers, conditioning language models on input queries (Stelmakh et al., 2022). Retrieval-augmented generation (RAG) has emerged as a compelling paradigm for such knowledge-intensive

QUERY: Explain the concept of quantum tunneling and its applications in modern technology.

FACTS:
(1) Can cross classically forbidden barriers
(2) Enables Josephson junctions
(3) Used in STM and tunnel diodes

RESPONSES:
A. Mentions ①②③
B. Mentions ②③
C. Mentions ②

HUMAN PREFERENCE:
A. ★★★★★
B. ★★★★★
C. ★★★★★

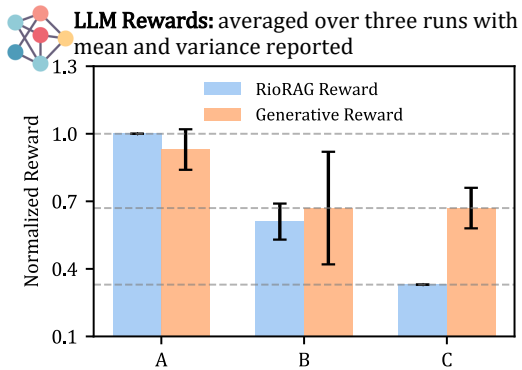


Figure 1: Existing long-form reward modeling exhibits (i) instability, with high variance when evaluating the same response multiple times, and (ii) unreliability, as their score rankings often differ from human judgments across responses.

tasks, as it grounds generation in factual content retrieved from external corpora (Lewis et al., 2020). However, producing reliable long-form answers remains challenging since large language models (LLMs) should synthesize information from multiple retrieved sources into coherent and factual paragraphs.

While RAG mitigates hallucination in long-form generation by grounding models on retrieved evidence (Asai et al., 2023), several challenges remain unresolved. First, LLMs often misuse factual information (Gao et al., 2023b; Łajewska and Balog, 2025), either due to residual hallucination or failure to extract the correct evidence from lengthy retrieved documents. Besides, generated answers fre-

quently exhibit limited informativeness (Ru et al., 2024), failing to comprehensively incorporate and reuse the available evidence. Recent efforts such as prompt engineering and template-based supervised fine-tuning (SFT) partially alleviate these issues (Wei et al., 2022), yet they rely heavily on strong annotation signals and external guidance. Moreover, these methods restrict models’ self-evolution and often reduce generalization and diversity (Shao et al., 2024), making them less adaptable to LFQA tasks across diverse domains.

Recently, reinforcement learning (RL) has improved LLMs in math and short-form QA by optimizing outcome-based feedback (Guo et al., 2025a). However, RL for LFQA is bottlenecked by reward design. Many standard rewards assume an exact-match notion of correctness, which does not apply to open-ended long-form answers. Rule-based outcome reward modelings (ORMs) are difficult to specify for diverse responses (Li et al., 2025). Moreover, long-form scoring is often unstable and hard to verify, even under well-guided evaluation protocols (Zhang et al., 2024). Reward models can also miss key evidence under long references, yielding misleading signals (Figure 1). Such instability can effectively sparsify the reward signal, further suppressing self-exploration and self-evolution during RL (Shao et al., 2025).

To address these challenges, we propose RioRAG, a Reinforced Informativeness Optimization-based RL method for long-form RAG. RioRAG aims to improve factual coverage with stable, verifiable reward learning. First, it defines informativeness as an externally verifiable objective and derives rewards from evidence support rather than heuristic lexical rules. Second, it employs a nugget-centric hierarchical verifier with length-adaptive scoring to provide fine-grained and action-discriminative feedback. This design mitigates reward sparsity and enables self-evolution of smaller policy models, without relying on handcrafted policy supervision or strong teacher-model distillation (e.g., sequential SFT). Extensive experiments on two published benchmarks, LongFact and RAGChecker, with zero-shot evaluation show that the RioRAG achieves superior performance compared with a series of state-of-the-art methods, demonstrating the effectiveness of the proposed innovations.

2 Task Formulation

LFQA extends conventional SFQA to generate coherent, factual, and detailed multi-paragraph responses (e.g., explanations or reports). Given a user query q and a large web corpus $\mathcal{D} = \{d_1, d_2, \dots, d_N\}$, a retriever R retrieves a subset of relevant documents $\mathcal{D}_q \subseteq \mathcal{D}$, and a generator G first produces a complete response sequence $y = G(q, \mathcal{D}_q)$ conditioned on both the query and the retrieved evidence. The response y naturally contains a reasoning part and a final answer part, which we separate by a predefined delimiter:

$$[r_{1:T} \parallel a_{1:M}] = y,$$

RioRAG follows an ORM formulation. Both evaluation and reward computation are based solely on the final answer a , while the reasoning content in y (e.g., chain-of-thought tokens) is ignored by the reward model. Formally, the objective of long-form RAG is to generate an answer a that maximizes factual correctness, informativeness, and coherence with respect to \mathcal{D}_q . In RioRAG, we approach this objective from an RL perspective, where the model learns to maximize a verifiable informativeness reward derived from the generated outcome.

3 Method

Figure 2 shows the overall pipeline of RioRAG. Given a query, RioRAG retrieves diverse and up-to-date web documents and generates long-form answers through reinforcement learning. The framework consists of two main components: (i) *reinforced informativeness optimization*, which maximizes factual coverage reward, and (ii) *nugget-centric hierarchical reward modeling*, which provides fine-grained, verifiable feedback based on evidence-level nuggets. Together, these components enable stable and unsupervised optimization of long-form RAG systems.

3.1 Reinforced Informativeness Optimization

Previous works mainly relied on SFT, which depends on pre-defined templates or powerful teacher models and thus limits generalization and adaptability to open-domain queries (Menick et al., 2022). Inspired by the recent success of RL in enhancing LLMs through self-exploration in tasks like math and coding (Yang et al., 2025), we extend this paradigm to the LFQA setting of RAG. RioRAG

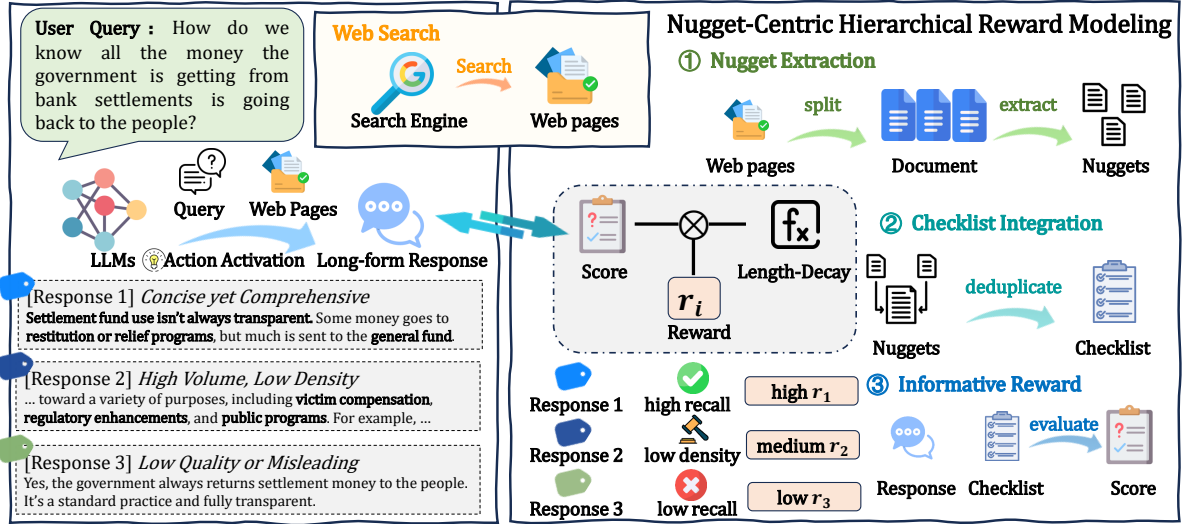


Figure 2: Overall illustration of the proposed RL-based RioRAG framework.

aims to enable models to self-improve by interacting with real-time web knowledge and optimizing generation quality without external supervision.

Training Data Construction. Existing RAG systems often depend on labeled data from strong teacher models and static retrieval corpora, both of which limit generalization to new domains and fail to reflect real-time web knowledge. To overcome these limitations, we reconstruct the training setup using ELI5 (Fan et al., 2019), retaining only queries while removing human-written answers. For each query, we dynamically retrieve the top- K webpages to form \mathcal{D}_q , ensuring that the model learns from current and diverse online sources. This design allows the model to explore and adapt to evolving web information under RL optimization, improving factuality and robustness in long-form generation.

Stable Reinforcement Optimization. Conventional RL optimization for long-form text often suffers from unstable gradients and inconsistent reward scaling, especially when model outputs vary significantly in length or quality. Inspired by relative normalization strategies that stabilize optimization in long-horizon reasoning tasks (Guo et al., 2025a), we adopt Group-wise Relative Policy Optimization (GRPO) (Shao et al., 2024) to improve training stability. Given G sampled completions $\{o_1, \dots, o_G\}$ with corresponding rewards $\{r_1, \dots, r_G\}$, GRPO normalizes rewards within

each sampled group:

$$A_i = \frac{r_i - \mu_r}{\sigma_r + \epsilon},$$

$$\mu_r = \frac{1}{G} \sum_{j=1}^G r_j, \quad \sigma_r = \sqrt{\frac{1}{G} \sum_{j=1}^G (r_j - \mu_r)^2}.$$

This group-wise normalization mitigates reward variance across samples, stabilizes policy updates, and leads to smoother reward improvement during long-form RL training.

Action Activation. Previous approaches introduce new delimiter tokens (e.g., <think>) to structure reasoning, which hinders warm-start from instruction-tuned LLMs and biases learning toward format imitation. We observe that even base models naturally support Markdown formatting. Hence, we leverage Markdown headers to separate reasoning and answers without adding new tokens. This lightweight design preserves pretrained priors, improves optimization stability, and provides a structured foundation for our subsequent hierarchical reward modeling.

3.2 Nugget-based Verifiable Reward Modeling

Existing ORMs for long-form generation are often non-verifiable and unstable (Ru et al., 2024), as they rely on heuristic or model-based scoring that may hallucinate or fluctuate across runs. To overcome these challenges, we propose a *nugget-based hierarchical reward modeling* approach that decomposes reward computation into three verifiable stages.

(1) Factual Nugget Extraction. For each retrieved document $d_i \in \mathcal{D}_q$, we identify *factual nuggets* as concise evidence statements. Each nugget n_{ij} is extracted as a short clause from d_i , forming a nugget set:

$$\mathcal{N}(\mathcal{D}_q) = \bigcup_{d_i \in \mathcal{D}_q} \text{Extract}(d_i).$$

(2) Evidence Checklist Synthesis. The extracted nuggets are clustered and merged into a unified checklist that captures all distinct factual evidence relevant to query q . This aggregation removes redundancy and enables consistent cross-document evaluation:

$$\mathcal{C}(q, \mathcal{D}_q) = \text{MergeCluster}(\mathcal{N}(\mathcal{D}_q)).$$

(3) Informativeness Assessment. Given a generated output o , we compute the informativeness reward by measuring the proportion of checklist nuggets covered in o :

$$\mathcal{I}(q, o, \mathcal{D}_q) = \frac{|\{n \in \mathcal{C}(q, \mathcal{D}_q) \mid n \subseteq o\}|}{|\mathcal{C}(q, \mathcal{D}_q)|}.$$

Each matched nugget can be explicitly traced to its source document, ensuring transparent and verifiable reward computation.

Length Decay. To prevent overestimation from excessively long responses, we apply a length-decay normalization on the reward computation. The decay term is applied only to the final answer segment, leaving the reasoning tokens unaffected to preserve the model’s deliberative process. Formally, the length-adjusted reward is defined as:

$$r_i = \begin{cases} s_i \exp[-k((l - l_0)/\tau)^m], & l > l_0, \\ s_i, & \text{otherwise,} \end{cases} \quad (1)$$

where s_i is the base informativeness score, l denotes the answer length, l_0 is the threshold of unpenalized length, τ controls the decay rate, and k, m regulate the sharpness of the decay curve. This design maintains factual compactness without penalizing intermediate reasoning, ensuring the reward focuses on the informativeness and precision of final outputs.

This hierarchical design stabilizes reward estimation, mitigates noise from heuristic scoring, and provides a reproducible, fact-grounded supervision signal for reinforcement learning.

Setting	Runtime	Complexity
Generative Reward	1.19	$\mathcal{O}((q + nd)^2)$
RioRAG sequential	2.32	$\mathcal{O}(n(q + d)^2 + q^2)$
RioRAG w/ parallel	0.98	–
RioRAG w/ async	0.72	–

Table 1: Reward computation cost comparison. Runtime is measured in seconds per training sample. Here q is the query length and d is document length.

3.3 Discussion

Novelty. RioRAG’s novelty comes from turning open-ended LFQA evaluation into explicit, verifiable reward computation. (1) **From vague holistic scoring to checklist credit assignment.** Instead of asking a judge model to output a coarse overall score under ambiguous rules, RioRAG converts informativeness into nugget-aligned checklist items. Each item becomes a concrete scoring point with clear credit. (2) **From long-context judging to short-context verification.** RioRAG first condenses retrieved evidence into nugget-aligned checklists. It then compares the model output against the checklists and aggregates item-wise scores into the final reward. This avoids long-context degradation, yields denser feedback, and supports self-evolution without strong teacher-model distillation (*e.g.*, sequential SFT).

Efficiency. Although RioRAG adds checklist extraction and nugget-based evaluation, these steps operate on compact inputs (single documents or short checklists), keeping token counts low. As a result, RioRAG is comparable to generative reward baselines in the sequential setting, and becomes faster with parallel execution. Reward computation can also be executed asynchronously and pipelined with RL training, so it does not block gradient updates in practice.

Theory. Our theoretical analysis further supports these findings. Nugget decomposition reduces reward variance as $\mathcal{O}(1/M)$, while the answer-only length decay enforces Lipschitz stability without affecting reasoning. The nugget locality also limits error accumulation and long-context degradation. These properties explain RioRAG’s stable convergence and factual consistency (see Appendix A).

4 Experiment

In this section, we detail the experimental setup, present the main results, and further support our

findings with ablation studies and in-depth analysis.

4.1 Experimental Setup

4.1.1 Datasets

We use the ELI5 dataset (Fan et al., 2019) as the training source, but only its question corpus is used without reference answers. This avoids overfitting to concise, single-perspective annotations and better reflects real retrieval-augmented settings. A total of 10K questions are randomly sampled for RL training.

For evaluation, we adopt two long-form QA benchmarks: **LongFact** (Wei et al., 2024) and **RAGChecker** (Ru et al., 2024). LongFact covers 38 domains consolidated into 8 major categories, with answers annotated by atomic factual units for fine-grained factual verification. RAGChecker includes 10 public datasets spanning 4K questions, designed to assess factual grounding and retrieval-based answer quality across multiple dimensions. Further dataset details are provided in Appendix B.

4.1.2 Evaluation Metrics

Following standard LFQA studies (Fan et al., 2019), we use *fact recall* (FR) and *information density* (ID) as the main metrics, measuring factual completeness and conciseness. We further report RAGChecker’s multi-dimensional metrics (Ru et al., 2024), including *faithfulness*, *hallucination*, and *context utilization*, to capture both factual reliability and retrieval effectiveness. Importantly, these RAGChecker metrics are used only for evaluation and are not provided to the model in any form during training, ensuring an objective assessment. Full metric definitions are presented in Appendix C.

4.1.3 Baselines

To evaluate the performance of RioRAG, we conduct comprehensive comparisons with various classical and state-of-the-art baseline methods across different categories, ensuring a thorough understanding of the proposed approach. The baselines are categorized into three groups based on their training paradigms: prompt-based unsupervised methods, supervised fine-tuning (SFT)-based approaches, and RL-based techniques. For prompt-based methods, we select GopherCite (Menick et al., 2022), chain-of-thought (Wei et al., 2022) and chain-of-note (Yu et al., 2024). Among SFT-based approaches, we employ chain-of-note and GopherCite with the SFT setting. For RL-based

methods, we adopt the Direct Preference Optimization (DPO) (Rafailov et al., 2023) framework. All baseline implementations are manually reimplemented with rigorous adherence to identical experimental configurations to ensure a fair comparison. This evaluation protocol guarantees the reliability of performance benchmarking while controlling for potential confounding factors in implementation differences.

4.2 Main Results

The results of different methods evaluated on LongFact and RAGChecker are shown in Table 2 and Table 3. It can be observed that:

(1) Our comprehensive evaluation reveals that SFT-based baselines substantially outperform prompt-based approaches, demonstrating the inherent limitations of prompt engineering in handling complex information synthesis tasks. The proposed RioRAG framework establishes a significant improvement across all metrics. This improvement stems from the reinforced informativeness optimization paradigm, which implements a nugget-centric hierarchical reward mechanism to guide LLMs in processing long-context inputs.

(2) Comprehensive evaluation on RAGChecker demonstrates that RioRAG excels in long-form RAG tasks across multiple critical dimensions, including knowledge point coverage, information density, retrieval utilization, hallucination mitigation, and internal knowledge integration. These results underscore the multidimensional efficacy of the proposed approach.

(3) Compared to off-line RL-based methods such as DPO, RioRAG demonstrates superior performance in long-form reasoning tasks. By leveraging an enhanced on-policy GRPO algorithm, RioRAG enables more comprehensive exploration of potential reasoning strategies during generation, thereby optimizing the LLM more effectively through informativeness-driven reward feedback.

4.3 Ablation Studies

In this section, we conduct an ablation study to evaluate the effectiveness of critical strategies in RioRAG comprehensively on LongFact. Here, we consider five variants built on RioRAG for evaluation: (a) *w/o Info. Optim.* removes the informativeness-based reward optimization during RL, replaced by direct quality evaluation; (b) *w/o Nugget Reward* removes the nugget-wise information extraction and use the full webpage

Method	Science		Tech.		Medicine		Law		Culture		Events		Commun.		Lifestyle		Average	
	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID
<i>Prompt-based Methods</i>																		
Direct-RAG	45.3	54.2	60.6	69.3	45.7	61.1	46.9	31.4	45.4	55.4	51.2	44.3	55.5	60.2	48.7	69.8	49.6	53.5
Chain-of-Thought	55.5	51.6	53.7	71.8	46.4	58.8	45.2	30.8	47.1	61.2	48.5	46.3	54.6	59.0	54.6	70.8	50.6	54.1
Chain-of-Note	45.0	52.6	58.4	71.5	43.6	58.7	42.3	27.1	40.9	56.9	47.7	46.5	47.8	56.9	48.8	71.0	46.0	52.5
GopherCite	54.4	56.6	63.8	73.9	56.2	54.1	55.6	33.7	53.3	60.6	48.9	46.4	61.5	64.0	54.1	72.5	55.9	56.1
<i>Supervised Fine-tuning based Methods</i>																		
Chain-of-Note	65.3	123.3	50.0	129.9	77.4	130.6	57.5	93.8	67.5	134.3	52.5	80.2	64.8	147.8	64.6	114.7	62.2	119.5
GopherCite	59.2	121.5	58.4	145.5	74.4	118.8	59.3	80.2	69.0	146.0	68.4	101.3	61.8	144.1	60.0	103.5	63.2	119.7
<i>RL-based Methods</i>																		
DPO	59.0	106.7	66.2	137.0	60.7	96.4	65.7	61.8	69.2	115.0	56.6	83.2	60.5	122.9	61.3	113.9	62.8	102.7
RioRAG	69.7	146.7	63.3	170.4	77.4	142.1	77.9	113.4	78.0	120.4	71.6	117.7	75.2	170.7	61.5	144.9	72.8	138.8

Table 2: The results on eight broader categories of LongFact benchmark with the average results of the eight categories, where FR denotes fact recall and ID denotes information density.

Method	Fact-Rec \uparrow	Info-Den \uparrow	Cont-Util \uparrow	Rel-NS \downarrow	Irrel-NS \downarrow	Hallu. \downarrow	Self-Know \downarrow	Faith. \uparrow
<i>Prompt-based Methods</i>								
Direct-RAG	38.3	91.6	22.6	8.2	7.5	37.0	8.1	45.3
Chain-of-Thought	50.4	146.5	24.3	4.6	4.2	30.2	9.7	48.0
Chain-of-Note	38.7	144.3	18.3	6.8	5.1	53.0	6.9	35.7
GopherCite	51.4	138.5	26.0	5.1	4.3	29.2	10.8	47.5
<i>Supervised Fine-tuning based Methods</i>								
Chain-of-Note	54.2	190.2	22.7	4.3	3.7	22.6	7.8	30.2
GopherCite	62.6	209.9	26.0	5.1	4.3	29.2	10.8	52.5
<i>RL-based Methods</i>								
DPO	61.2	149.6	26.0	5.2	6.0	27.8	8.0	53.1
RioRAG	66.3	224.6	27.8	4.3	3.6	20.9	5.0	58.2

Table 3: Average results across ten domains on the RAGChecker benchmark. Fact-Rec refers to fact recall, Info-Den to information density, Cont-Util to context utilization, Rel-NS and Irrel-NS to relevant and irrelevant noise sensitivity, Hallu. to hallucination, Self-Know to self-knowledge, and Faith. to faithfulness.

for checklist integration; (c) *w/o Length Decay* eliminates the length penalty in Equation (1); (d) *w/ Generative GRPO* denotes the setting where the policy is trained with a 32B generative reward model providing scalar feedback, rather than the verifiable reward used in RioRAG; (e) *w/ Off-Policy RL* utilizes an off-policy RL method that employs a static sampling strategy wherein all queries are pre-processed through offline rollouts to generate complete trajectories before being uniformly scored.

Table 4 presents the results for the variants of our method, from which we can observe the following findings: (a) The performance drops in *w/o Info. Optim.*, demonstrating that using informativeness as the objective for optimization enhances the performance of long-form RAG models through the guidance of reasoning. (b) The performance drops in *w/o Nugget Reward*, demonstrating incorporating nugget-wise information extraction enables the model to better capture core facts. (c) The performance drops in *w/o Length Decay*, underscoring the critical role of incorporating the

length penalty in mitigating excessive response length. (d) The performance degradation observed in *w/ Generative RM* indicates that, without an appropriate pipeline design, unstable and unreliable rewards hinder effective fine-tuning. (e) The performance significantly drops in *w/ Off-Policy RL*, demonstrating that the off-policy method may contain a mismatch between the behavior policy and the target policy.

4.4 Further Analysis

4.4.1 Scaling Law of RioRAG

To investigate the scalability characteristics of RioRAG, we conduct a systematic analysis using the Qwen2.5 model with varying parameter sizes (1.5B, 7B, and 14B). As illustrated in Figure 3 (a), the experimental results demonstrate that RioRAG significantly outperforms SFT at all model scales, with performance consistently improving in accordance with scaling laws.

We can first observe that larger models exhibit improved semantic understanding for both query

Method	Science		Tech.		Medicine		Law		Culture		Events		Commun.		Lifestyle		Average	
	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID	FR	ID
RioRAG	69.7	146.7	63.3	170.4	77.4	142.1	77.9	113.4	78.0	120.4	71.6	117.7	75.2	170.7	61.5	144.9	72.8	138.8
w/o Info. Optim.	34.4	79.6	53.2	147.6	32.0	92.3	33.4	66.6	40.9	87.5	43.2	76.6	43.9	106.4	36.6	97.3	39.5	90.8
w/o Nugget Reward	36.0	77.6	56.6	119.3	23.8	62.5	36.0	42.4	37.2	80.7	30.7	65.4	36.8	83.2	40.1	112.7	37.0	77.3
w/o Length Decay	57.0	99.2	65.5	139.9	58.3	109.9	62.0	88.7	56.5	87.5	43.7	63.0	45.2	70.7	68.5	108.8	56.2	91.3
w/ Generative GRPO	58.7	104.3	63.3	159.2	36.6	62.7	55.7	70.9	63.1	112.7	54.3	70.3	52.5	103.2	61.7	99.3	54.7	97.2
w/ Off-Policy RL	41.6	41.4	61.7	65.5	34.1	35.0	44.8	25.1	46.8	52.2	46.4	39.8	56.2	53.1	61.8	60.2	49.3	45.5

Table 4: Results of the RioRAG variants on LongFact.

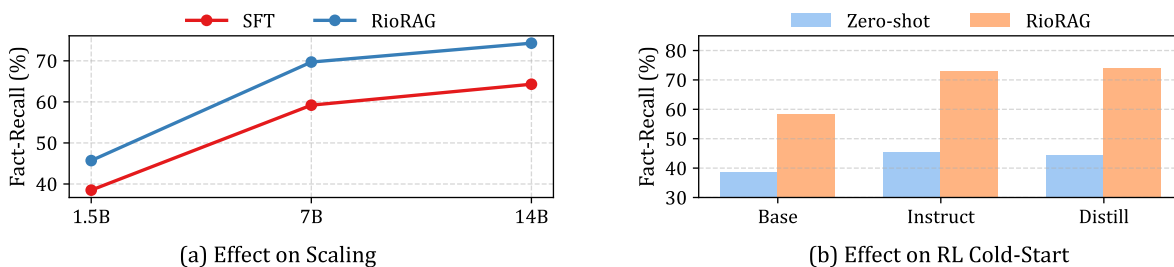


Figure 3: In-depth analysis on scaling law and RL cold-start.

433 formulation and webpage relevance assessment. 434 Second, RioRAG benefits from increased model 435 capacity for learning sophisticated retrieval utiliza- 436 tion strategies. Third, the enhanced generation ca- 437 pability of larger models enables more effective 438 utilization of retrieved webpages while reducing 439 hallucination risks through better alignment with 440 the reward model’s feedback. Notably, the perfor- 441 mance growth curve shows a sublinear relation- 442 ship between model size and metric improvements, 443 aligning with observations from language model 444 scaling studies (Wei et al.). This phenomenon sug- 445 gests that while our RioRAG framework effectively 446 leverages model scale, there exists an upper bound 447 where additional parameters may not proportion- 448 ally improve RAG performance, which is a critical 449 consideration for practical system deployment.

4.4.2 Effect on RL Cold-Start

450 To examine how model initialization affects 451 RL training, we compare three variants: a base 452 model (Qwen2.5-7B-Base), an instruction-tuned 453 model (Qwen2.5-7B-Instruct), and an R1-distilled 454 model (R1-Distilled-Qwen2.5-7B) incorporating 455 DeepSeek R1’s slow-thinking distillation (Guo 456 et al., 2025b). As shown in Figure 3(b), the 457 instruction-tuned model gains 24.4%, and the R1- 458 distilled model achieves the largest improvement 459 of 29.6% after RL training. 460

461 The results show that initialization strongly af- 462 fects training stability and reward learning. The

463 base model, lacking alignment and reasoning pri- 464 ors, struggles to explore effectively. In contrast, the 465 R1-distilled model benefits from pre-established 466 reasoning patterns, enabling more stable reward 467 estimation and efficient policy updates. This sup- 468 ports that structured reasoning provides a favorable 469 starting point for RL.

4.4.3 Generation Length and Reward

470 To investigate the dynamics of generation length 471 control during RL training, we systematically ana- 472 lyze the interaction between sequence length evo- 473 lution and reward optimization. Figure 4 illus- 474 trates the co-evolution of average generation length 475 and reward scores across training steps. 476

477 Without length control, the model often pro- 478 duces overly long answers, while the reward score 479 stays almost unchanged. This shows a form of re- 480 ward hacking, where the model gains higher scores 481 through longer outputs instead of better content. 482 After adding the length-decay term, the model 483 first explores longer responses and then learns to 484 shorten them while keeping rewards stable. The re- 485 sults show that length regularization helps improve 486 both clarity and information density in long-form 487 generation.

5 Related Work

5.1 Long-Form Question Answering

488 Research on LFQA has progressed through three 489 shifts: fine-tuned generative models, retrieval- 490 491

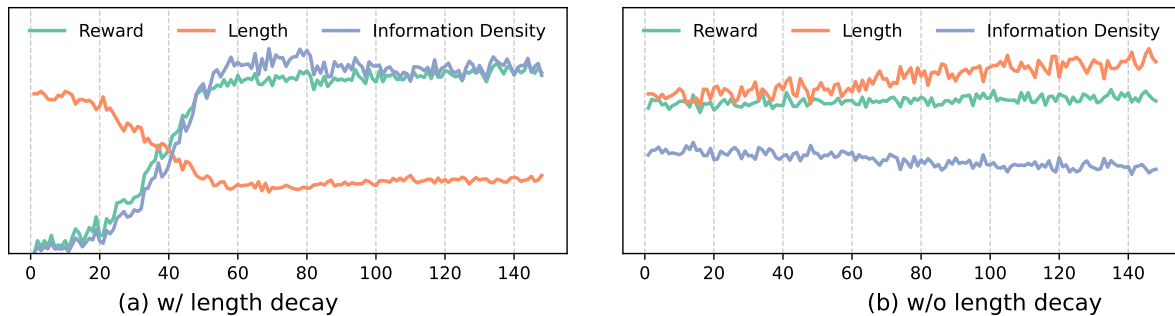


Figure 4: Analysis on co-evolution of generation length and reward during RL training.

492 augmented architectures, and human-aligned
 493 LLMs. Early abstractive LFQA was enabled by
 494 datasets such as ELI5 (Fan et al., 2019), showing
 495 that seq2seq models can generate plausible long
 496 answers with retrieved evidence. More recent sys-
 497 tems leverage LLMs with human feedback. Of-
 498 ten via RL from human preferences. To produce
 499 answers grounded in explicit quotations (Menick
 500 et al., 2022). RAG has since become the dominant
 501 approach for factual grounding, including train-
 502 ing models in web-browsing interfaces (Nakano
 503 et al., 2021; Qin et al., 2023). Verifiability can
 504 be strengthened by training with explicit source
 505 citation (Menick et al., 2022), while post-hoc at-
 506 tribution verifies pre-generated text (Wang et al.,
 507 2022; Gao et al., 2023a). Given LFQA’s open-
 508 ended nature, faithfulness to evidence is a central
 509 objective (Gao et al., 2023b; Zhao et al., 2024). It
 510 can be improved through open-book training with
 511 citations (Menick et al., 2022) or modeled via prob-
 512 abilistic calibration of answer correctness (Huang
 513 et al., 2024). Complementarily, LLM-based eval-
 514 uators provide automated quality assessment for
 515 LFQA (Han et al., 2024).

5.2 Reinforcement Learning based RAG

517 RL has emerged as a promising tool for improv-
 518 ing RAG by optimizing retrieval and generation
 519 with reward-driven policy updates. Early work
 520 strengthens behavior-cloned web-browsing agents
 521 with human-feedback rewards to improve factual
 522 alignment (Nakano et al., 2021). Subsequent stud-
 523 ies design fine-grained rewards for coherence and
 524 information gain (Cai et al., 2024), or build domain-
 525 specific reward models with synthetic supervision
 526 to align RAG with human performance (Nguyen
 527 et al., 2024). Others use composite reward ensem-
 528 bles to balance answer quality and coverage (Zhang
 529 et al., 2025). On the retrieval side, modules can be

530 optimized via group-wise relative policy optimiza-
 531 tion (Huang et al., 2025) or multi-agent coordina-
 532 tion (Chen et al., 2025). Cost-sensitive retrieval
 533 further uses value estimation to decide when to
 534 invoke external search under latency–utility trade-
 535 offs (Kulkarni et al., 2024). Recent progress, ex-
 536 emplified by DeepSeek-R1 (Guo et al., 2025b), has
 537 also motivated RL for autonomous retrieval invo-
 538 cation within long reasoning chains (Jin et al., 2025).
 539 In contrast, RioRAG introduces nugget-centric hi-
 540 erarchical reward modeling to optimize verifiable
 541 informativeness for long-form RAG, without hand-
 542 crafted policy supervision or strong teacher-model
 543 distillation.

6 Conclusion

544
 545 In this work, we address long-form RAG lim-
 546 itations through RioRAG, an RL framework that
 547 redefines long-form RAG training via reinforced
 548 informativeness optimization with nugget-centric
 549 hierarchical reward modeling. RioRAG directly op-
 550 timizes informativeness through a quantifiable re-
 551 ward design for factual alignment, without the need
 552 for scarce training data. Our experiments on two
 553 benchmarks demonstrate that RioRAG fundamen-
 554 tally improves the quality of long-form RAG. By
 555 addressing the core challenges identified in long-
 556 form RAG, RioRAG advances the development
 557 of trustworthy generative systems for real-world
 558 knowledge applications. Moreover, the success of
 559 nugget-level reward modeling suggests that future
 560 evaluation frameworks for long-form tasks should
 561 prioritize granular factual alignment over surface-
 562 level metrics. Limitations include the current fo-
 563 cus on English corpora and reliance on automatic
 564 nugget extraction, which may inherit biases from
 565 pre-trained models. For future work, we will ex-
 566 tend the framework to multilingual settings and
 567 investigate human-in-the-loop reward refinement.

568 Limitations

569 While RioRAG effectively improves the stability
570 and verifiability of long-form RAG training, several
571 limitations remain. First, our current reward pri-
572 marily targets factual informativeness and does not
573 explicitly capture other aspects of long-form qual-
574 ity, such as linguistic style, coherence, and read-
575 ability. These factors may further influence human
576 preference alignment and should be considered in
577 future extensions. Second, although we examine
578 models of different scales, computational resources
579 restrict us from scaling beyond 32B parameters.
580 Larger-scale experiments (*e.g.*, 72B) could provide
581 deeper insights into the scalability and robustness
582 of verifiable reward optimization.

583 We used AI assistants for minor language polish-
584 ing; all technical content and results were produced
585 and verified by the authors.

586 References

587 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and
588 Hannaneh Hajishirzi. 2023. Self-rag: Learning to
589 retrieve, generate, and critique through self-reflection.
590 In *The Twelfth International Conference on Learning*
591 *Representations*.

592 Tianchi Cai, Zhiwen Tan, Xierui Song, Tao Sun, Jiyan
593 Jiang, Yunqi Xu, Yinger Zhang, and Jinjie Gu. 2024.
594 Forag: Factuality-optimized retrieval augmented gen-
595 eration for web-enhanced long-form question answer-
596 ing. In *Proceedings of the 30th ACM SIGKDD Con-*
597 *ference on Knowledge Discovery and Data Mining*,
598 pages 199–210.

599 Yiqun Chen, Lingyong Yan, Weiwei Sun, Xinyu Ma,
600 Yi Zhang, Shuaiqiang Wang, Dawei Yin, Yiming
601 Yang, and Jiaxin Mao. 2025. Improving retrieval-
602 augmented generation through multi-agent reinforc-
603 ement learning. *arXiv preprint arXiv:2501.15228*.

604 LE Clarke. 1980. Probability and measure, by patrick
605 billingsley. pp 515.£ 15. 20. 1979. sbn 0 471 03173
606 9 (wiley). *The Mathematical Gazette*, 64(430):293–
607 294.

608 Angela Fan, Yacine Jernite, Ethan Perez, David Grang-
609 er, Jason Weston, and Michael Auli. 2019. Eli5:
610 Long form question answering. In *Proceedings of*
611 *the 57th Annual Meeting of the Association for Com-*
612 *putational Linguistics*, pages 3558–3567.

613 Luyu Gao, Zhu Yun Dai, Panupong Pasupat, Anthony
614 Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vin-
615 cent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan,
616 et al. 2023a. Rarr: Researching and revising what
617 language models say, using language models. In
618 *Proceedings of the 61st Annual Meeting of the As-*
619 *sociation for Computational Linguistics (Volume 1:*
620 *Long Papers)*, pages 16477–16508.

Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen.
2023b. Enabling large language models to generate
text with citations. In *Proceedings of the 2023 Con-*
ference on Empirical Methods in Natural Language
Processing, pages 6465–6488.

626 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
627 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,
628 Peiyi Wang, Xiao Bi, et al. 2025a. Deepseek-r1: In-
629 centivizing reasoning capability in llms via reinforc-
630 ement learning. *arXiv preprint arXiv:2501.12948*.

631 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
632 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,
633 Peiyi Wang, Xiao Bi, et al. 2025b. Deepseek-r1:
634 Incentivizing reasoning capability in llms via rein-
635 forcement learning. *CoRR*.

636 Rujun Han, Yuhao Zhang, Peng Qi, Yumo Xu, Jinyuan
637 Wang, Lan Liu, William Yang Wang, Bonan Min, and
638 Vittorio Castelli. 2024. Rag-qa arena: Evaluating do-
639 main robustness for long-form retrieval augmented
640 question answering. In *Proceedings of the 2024 Con-*
641 *ference on Empirical Methods in Natural Language*
642 *Processing*, pages 4354–4374.

643 Jerry Huang, Siddarth Madala, Risham Sidhu, Cheng
644 Niu, Julia Hockenmaier, and Tong Zhang. 2025.
645 Rag-rl: Advancing retrieval-augmented generation
646 via rl and curriculum learning. *arXiv preprint*
647 *arXiv:2503.12759*.

648 Yukun Huang, Yixin Liu, Raghuveer Thirukovalluru,
649 Arman Cohan, and Bhuwan Dhingra. 2024. Calibrat-
650 ing long-form generations from large language mod-
651 els. In *Findings of the Association for Computational*
652 *Linguistics: EMNLP 2024*, pages 13441–13460.

653 Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon,
654 Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei
655 Han. 2025. Search-r1: Training llms to reason and
656 leverage search engines with reinforcement learning.
657 *arXiv preprint arXiv:2503.09516*.

658 Mandar Kulkarni, Praveen Tangarajan, Kyung Kim, and
659 Anusua Trivedi. 2024. Reinforcement learning for
660 optimizing rag for domain chatbots. *arXiv preprint*
661 *arXiv:2401.06800*.

662 Weronika Łajewska and Krisztian Balog. 2025. Ginger:
663 Grounded information nugget-based generation of
664 responses. *arXiv preprint arXiv:2503.18174*.

665 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
666 Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
667 rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-
668 täschel, et al. 2020. Retrieval-augmented generation
669 for knowledge-intensive nlp tasks. *Advances in neu-*
670 *ral information processing systems*, 33:9459–9474.

671 Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang,
672 Yujia Zhou, Yutao Zhu, Peitian Zhang, and
673 Zhicheng Dou. 2025. Search-o1: Agentic search-
674 enhanced large reasoning models. *arXiv preprint*
675 *arXiv:2501.05366*.

676	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts. <i>arXiv preprint arXiv:2307.03172</i> .	Zhihong Shao, Yuxiang Luo, Chengda Lu, ZZ Ren, Jiewen Hu, Tian Ye, Zhibin Gou, Shirong Ma, and Xiaokang Zhang. 2025. Deepseekmath-v2: Towards self-verifiable mathematical reasoning. <i>arXiv preprint arXiv:2511.22570</i> .	731 732 733 734 735
681	Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. Www’18 open challenge: financial opinion mining and question answering. In <i>Companion proceedings of the the web conference 2018</i> , pages 1941–1942.	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .	736 737 738 739 740
687	Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, et al. 2022. Teaching language models to support answers with verified quotes. <i>arXiv preprint arXiv:2203.11147</i> .	Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and Ming-Wei Chang. 2022. Asqa: Factoid questions meet long-form answers. <i>arXiv preprint arXiv:2204.06092</i> .	741 742 743 744
693	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. <i>arXiv preprint arXiv:2112.09332</i> .	Cunxiang Wang, Ruoxi Ning, Boqi Pan, Tonghui Wu, Qipeng Guo, Cheng Deng, Guangsheng Bao, Qian Wang, and Yue Zhang. 2024. Novelqa: A benchmark for long-range novel question answering. <i>arXiv e-prints</i> , pages arXiv–2403.	745 746 747 748 749
699	Ani Nenkova, Rebecca Passonneau, and Kathleen Mckeown. 2007. The pyramid method: Incorporating human content selection variation in summarization evaluation. <i>ACM Transactions on Speech and Language Processing (TSLP)</i> , 4(2):4–es.	Shufan Wang, Fangyuan Xu, Laure Thompson, Eunsol Choi, and Mohit Iyyer. 2022. Modeling exemplification in long-form question answering via retrieval. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2079–2092.	750 751 752 753 754 755 756
704	Thang Nguyen, Peter Chin, and Yu-Wing Tai. 2024. Reward-rag: Enhancing rag with reward driven supervision. <i>arXiv preprint arXiv:2410.03780</i> .	Larry Wasserman. 2013. <i>All of statistics: a concise course in statistical inference</i> . Springer Science & Business Media.	757 758 759
707	Yujia Qin, Zihan Cai, Dian Jin, Lan Yan, Shihao Liang, Kunlun Zhu, Yankai Lin, Xu Han, Ning Ding, Huadong Wang, et al. 2023. Webcpm: Interactive web search for chinese long-form question answering. In <i>The 61st Annual Meeting Of The Association For Computational Linguistics</i> .	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. <i>Transactions on Machine Learning Research</i> .	760 761 762 763 764
713	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36:53728–53741.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	765 766 767 768 769
719	Sara Rosenthal, Avirup Sil, Radu Florian, and Salim Roukos. 2025. Clapnq: Cohesive long-form answers from passages in natural questions for rag systems. <i>Transactions of the Association for Computational Linguistics</i> , 13:53–72.	Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu, Nathan Zixia Hu, Jie Huang, Dustin Tran, Daiyi Peng, Ruibo Liu, Da Huang, et al. 2024. Long-form factuality in large language models. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	770 771 772 773 774 775
724	Dongyu Ru, Lin Qiu, Xiangkun Hu, Tianhang Zhang, Peng Shi, Shuaichen Chang, Cheng Jiayang, Cunxiang Wang, Shichao Sun, Huanyu Li, et al. 2024. Ragchecker: A fine-grained framework for diagnosing retrieval-augmented generation. <i>Advances in Neural Information Processing Systems</i> , 37:21999–22027.	Fangyuan Xu, Kyle Lo, Luca Soldaini, Bailey Kuehl, Eunsol Choi, and David Wadden. 2024. Kiwi: A dataset of knowledge-intensive writing instructions for answering research questions. <i>arXiv preprint arXiv:2403.03866</i> .	776 777 778 779 780
729		An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	781 782 783 784

785 Wenhao Yu, Hongming Zhang, Xiaoman Pan, Peixin
786 Cao, Kaixin Ma, Jian Li, Hongwei Wang, and Dong
787 Yu. 2024. Chain-of-note: Enhancing robustness in
788 retrieval-augmented language models. In *Proceed-*
789 *ings of the 2024 Conference on Empirical Methods in*
790 *Natural Language Processing*, pages 14672–14685.

791 Hanning Zhang, Juntong Song, Juno Zhu, Yuanhao Wu,
792 Tong Zhang, and Cheng Niu. 2025. Rag-reward:
793 Optimizing rag with reward modeling and rlhf. *arXiv*
794 *preprint arXiv:2501.13264*.

795 Jiajie Zhang, Zhongni Hou, Xin Lv, Shulin Cao, Zhenyu
796 Hou, Yilin Niu, Lei Hou, Yuxiao Dong, Ling Feng,
797 and Juanzi Li. 2024. Longreward: Improving long-
798 context large language models with ai feedback.
799 *arXiv preprint arXiv:2410.21252*.

800 Yilun Zhao, Lyuhao Chen, Arman Cohan, and Chen
801 Zhao. 2024. Tapera: enhancing faithfulness and in-
802 terpretability in long-form table qa by content plan-
803 ning and execution-based reasoning. In *Proceedings*
804 *of the 62nd Annual Meeting of the Association for*
805 *Computational Linguistics (Volume 1: Long Papers)*,
806 pages 12824–12840.

A Theoretical Analysis

This section provides formal reasoning for the stability and reliability of RioRAG’s verifiable reward design.

A.1 Variance reduction by verifiable nugget decomposition

Let $\mathcal{C}(q, \mathcal{D}_q) = \{n_k\}_{k=1}^M$ denote the factual checklist, and define the per-nugget binary correctness as

$$Z_k(o) = \mathbf{1}\{n_k \subseteq o\}, \quad C(o) = \frac{1}{M} \sum_{k=1}^M Z_k(o).$$

Assume (A1) each Z_k is an unbiased Bernoulli variable with variance $\sigma_k^2 \leq 1/4$, and pairwise covariances are bounded by $\text{Cov}(Z_i, Z_j) \leq \rho$ (Clarke, 1980). Then

$$\begin{aligned} \text{Var}[C(o)] &= \frac{1}{M^2} \left(\sum_{k=1}^M \sigma_k^2 + 2 \sum_{1 \leq i < j \leq M} \text{Cov}(Z_i, Z_j) \right) \\ &\leq \frac{\bar{\sigma}^2}{M} + \rho \frac{M-1}{M}. \end{aligned}$$

where $\bar{\sigma}^2 = \frac{1}{M} \sum_k \sigma_k^2$. When nuggets are nearly independent ($\rho \approx 0$), the variance decays at least as $\mathcal{O}(1/M)$, yielding more stable and concentrated reward estimates than single long-form scoring (Nenkova et al., 2007).

A.2 Lipschitz stability from answer-only length decay

Let the final reward be $R(o) = C(o)\phi(l)$, where $\phi(l) = \exp[-k((l - l_0)/\tau)^m]$ for $l > l_0$ and 1 otherwise. Assume $C(o)$ is L_s -Lipschitz w.r.t. perturbations in the answer tokens, and $\phi(l)$ is L_ϕ -Lipschitz in length (Wasserman, 2013). For two outputs o, o' with answer lengths l, l' ,

$$\begin{aligned} |R(o) - R(o')| &= |C(o)\phi(l) - C(o')\phi(l')| \\ &\leq L_s \|o - o'\| + L_\phi |l - l'|. \end{aligned}$$

Hence $R(\cdot)$ is Lipschitz continuous and robust to small length variations. Since $\phi(l)$ is smooth and bounded in $[0, 1]$, its Lipschitz constant $L_\phi = \sup_l |\phi'(l)|$ exists under finite output lengths. Because $\phi(l)$ is applied only to the answer segment, the reasoning process remains unaffected, preventing reward drift and ensuring stable optimization.

A.3 Mitigating long-context degradation

Let $p(t)$ denote the position-dependent scoring error probability of a long-form evaluator. For monolithic evaluation, the expected error over T tokens is $E_{\text{mono}} = \frac{1}{T} \sum_{t=1}^T p(t)$, which increases when mid-context positions have high $p(t)$ (“lost in the middle” (Liu et al., 2023)). In RioRAG, each nugget corresponds to a short, localized window where the local scoring error \tilde{p} is upper-bounded by $\sup_t p(t_{\text{nugget}})$. Thus,

$$E_{\text{nugget}} = \frac{1}{M} \sum_{k=1}^M \tilde{\xi}_k, \quad \mathbb{E}[E_{\text{nugget}}] \leq \tilde{p} \ll \mathbb{E}[E_{\text{mono}}],$$

showing that nugget-level locality effectively suppresses long-context degradation.

Summary. Theoretical analysis shows that our verifiable reward design (i) reduces reward variance with more nuggets, (ii) ensures stable optimization through smooth answer-length decay, and (iii) alleviates long-context degradation via local nugget evaluation.

B Details on Datasets

In this section, we provide detailed descriptions of the two comprehensive benchmarks used in our experiments: **LongFact** (Wei et al., 2024) and **RAGChecker** (Ru et al., 2024). These datasets are designed to evaluate long-form retrieval-augmented generation (RAG) systems across diverse topics and multiple dimensions of factual quality. Their complementary nature enables robust assessment of both factual coverage and fine-grained answer quality in open-domain settings.

LongFact. LongFact is a manually curated benchmark focused on evaluating long-form factuality. It contains a diverse set of fact-seeking questions, where each gold answer synthesizes multiple atomic facts drawn from various evidence sources. The dataset is notable for its broad coverage across 38 fine-grained domains, which are grouped into the following 8 broader categories to support structured evaluation:

- **Science & Nature:** physics, chemistry, biology, astronomy, virology, prehistory
- **Technology & Computing:** computer science, computer security, machine learning, electrical engineering, mathematics

- 887 • **Medicine & Psychology:** medicine, clinical
888 knowledge, psychology, psychology check-
889 point
- 890 • **Law & Politics:** international law, immigra-
891 tion law, U.S. foreign policy, jurisprudence
- 892 • **Social Sciences & Culture:** sociology, geog-
893 raphy, world religions, moral disputes, philos-
894 ophy
- 895 • **History & Events:** history, 20th-century
896 events, global facts, economics
- 897 • **Business & Communication:** business ethics,
898 accounting, marketing, management, public
899 relations
- 900 • **Entertainment & Lifestyle:** movies, music,
901 gaming, celebrities, architecture, sports

902 Each example in LongFact is annotated with
903 atomic information units, enabling precise mea-
904 surement of factual recall and information density.
905 This makes it especially well-suited for evaluating
906 long-form answers that integrate knowledge from
907 multiple sources.

908 **RAGChecker.** RAGChecker is a comprehen-
909 sive benchmark designed to evaluate long-form
910 Retrieval-Augmented Generation (RAG) systems
911 across diverse domains. It repurposes examples
912 from 10 public datasets, encompassing a total of
913 4,162 questions. For the 4 subsets we used, we
914 briefly describe their characteristics below:

- 915 • **ClapNQ (Rosenthal et al., 2025):** Derived
916 from Natural Questions (NQ), ClapNQ in-
917 cludes long-form answers with grounded gold
918 passages from Wikipedia, focusing on gener-
919 ating cohesive long-form answers from non-
920 contiguous text segments.
- 921 • **NovelQA (Wang et al., 2024):** NovelQA is
922 a benchmark designed to evaluate large lan-
923 guage models on deep narrative understanding
924 through complex questions based on English
925 novels.
- 926 • **FiQA (Maia et al., 2018):** A financial ques-
927 tion answering dataset comprising 500 QA
928 pairs, where short answers are extended to
929 long-form using GPT-4, filtered to remove
930 hallucinations.

- 931 • **KIWI (Xu et al., 2024):** A dataset of
932 knowledge-intensive writing instructions for
933 answering research questions, comprising 71
934 QA pairs with long-form answers validated
935 for quality.

C Details on Evaluation Metrics 936

937 We adopt two categories of metrics to com-
938 prehensively evaluate factual quality and retrieval
939 grounding.

940 **Standard LFQA Metrics.** Following prior
941 work (Fan et al., 2019), we use *Fact Recall* (FR) to
942 measure factual completeness—the ratio of atomic
943 facts in the generated response to those in the ref-
944 erence answer—and *Information Density* (ID), de-
945 fined as the ratio of atomic facts to total response
946 length, reflecting conciseness and informativeness.

947 **RAGChecker Metrics.** The RAGChecker bench-
948 mark (Ru et al., 2024) introduces a suite of ad-
949 vanced metrics: *Faithfulness* (proportion of cor-
950 rect facts supported by retrieved pages), *Relevant*
951 *Noise Sensitivity* (ratio of incorrect facts appear-
952 ing in retrieved content), *Irrelevant Noise Sensitiv-*
953 *ity* (share of correct facts that are irrelevant to re-
954 trieval), *Hallucination* (incorrect facts unsupported
955 by any retrieved page), *Self-Knowledge* (correct
956 facts absent from all retrieved content), and *Con-*
957 *text Utilization* (fraction of ground-truth facts cov-
958 ered by retrieval). Together, these metrics provide a
959 multi-dimensional evaluation of factual grounding,
960 reliability, and retrieval efficiency.