

Test-Time Strategies for More Efficient and Accurate Agentic RAG

Anonymous ACL submission

Abstract

Agentic retrieval-augmented generation (RAG) frameworks such as Search-R1 (Jin et al., 2025) improve performance on complex, multi-hop questions by interleaving reasoning and retrieval, but they often retrieve previously seen information and struggle to integrate evidence into the current reasoning state. We study lightweight test-time modifications to Search-R1 that target these inefficiencies without changing training: contextualization of retrieved passages, de-duplication of previously seen documents, and a hybrid of both. Evaluated on 500-example subsets of HotpotQA (Yang et al., 2018) and Natural Questions (Kwiatkowski et al., 2019) using exact match (EM), LLM-based semantic matching, and average retrieval turns, our best method—contextualization with GPT-4.1-mini—improves exact match by 5.6% and reduces retrieval turns by 10.5% relative to the Search-R1 baseline. These results suggest that improved use of retrieved evidence, rather than enforcing retrieval diversity alone, is the main driver of better answer accuracy and efficiency.

1 Introduction

RAG systems have shown promising results in complex question answering (QA) tasks by combining external document retrieval with generative language models (Lewis et al., 2020). Despite this success, traditional RAG systems that rely on a single-step retrieval and generation process often struggle to handle complex or nuanced questions, especially those requiring deep contextual understanding and multi-hop retrieval. To address these complexities, recent research has proposed agentic RAG systems which utilize large language model (LLM) agents to orchestrate retrieval, refine search queries, and optimize responses (Singh et al., 2025; An et al., 2024; Chan et al., 2024). Another popular approach is to augment the reasoning loop of LLMs with a retrieval tool, enabling the model to

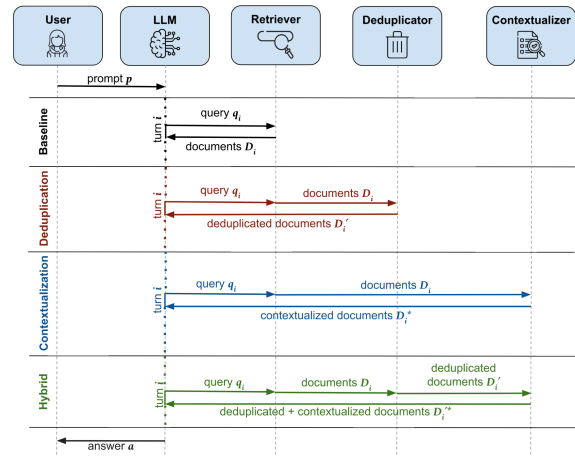


Figure 1: An illustration of the information flow for our proposed test-time strategies compared to the baseline during a single inference turn i . *Baseline*: This represents the standard Search-R1 framework, where the LLM sends a query (q_i) to the retriever and directly receives the retrieved documents (D_i) to continue its reasoning. *Deduplication*: This approach filters out previously seen content and returns only a set of novel documents (D_i') to the LLM. *Contextualization*: This approach parses the retrieved documents (D_i) and reformulates their content to improve integration into the LLM’s reasoning process, returning an enhanced set of information (D_i^*). *Hybrid*: This approach combines both modules sequentially.

autonomously use retrieval while performing multi-step reasoning (Li et al., 2025; Jin et al., 2025).

A notable example of this approach is the Search-R1 framework, which uses reinforcement learning (RL) to train LLMs for interleaved reasoning and retrieval (Jin et al., 2025). At inference time, the Search-R1 model first performs reasoning on a given user prompt p to either produce an answer a or generate a search query q to retrieve supporting information D . More specifically, in the i -th turn, the query q_i is sent to a dense retriever E5 (Wang et al., 2022), which returns relevant passages D_i from a 2018 Wikipedia dump. D_i is directly incor-

056	porated into the reasoning trace and fed back into		
057	the LLM, which continues its reasoning and repeats		
058	these steps until the final answer a is produced, see		
059	Figure 1 (Baseline). The Search-R1 models are		
060	trained using RL methods, such as PPO and GRPO,		
061	optimizing the exact match (EM) score between		
062	the ground truth and the predicted answer.		
063	While Search-R1 has achieved substantial		
064	improvement—up to 41% over its baseline—our		
065	analysis of the Qwen2.5-7B Search-R1 model dur-		
066	ing inference has revealed several shortcomings.		
067	First, the model often performs repetitive retrieval		
068	of previously processed information, which leads		
069	to unnecessary retrieval turns and increased token		
070	consumption and latency. Second, the model of-		
071	ten struggles to effectively contextualize retrieved		
072	passages, leading to suboptimal reasoning and in-		
073	accurate answers.		
074	Research questions		
075	• Will a concise representation of relevant infor-		
076	mation help an LLM become more efficient		
077	and accurate in question-answering tasks?		
078	• Can preventing redundant document re-		
079	trieval encourage greater contextual diversity,		
080	thereby improving efficiency and answer ac-		
081	curacy?		
082	Proposed Approach		
083	Our work builds upon and extends the Search-R1		
084	framework (Jin et al., 2025) and investigates test-		
085	time approaches to improve the framework’s rea-		
086	soning efficiency and final answer accuracy. We		
087	address the limitations of Search-R1 through three		
088	test-time modifications that process the retrieved		
089	results D : (1) a contextualization module, (2) a		
090	De-duplication module, and (3) a hybrid approach		
091	that combines both.		
092	2 Related Work		
093	Prior work has improved retrieval-augmented gen-		
094	eration by managing previously retrieved content		
095	(Shi et al., 2024), extracting useful information		
096	from retrieved documents before further reasoning		
097	(Li et al., 2025), and training retrieval-aware rea-		
098	soning models with reinforcement learning (Huang		
099	et al., 2025). In contrast, we study lightweight test-		
100	time modifications to Search-R1 without changing		
101	model training or architecture.		
	3 Approach		102
	Our qualitative analysis of Search-R1 reasoning		103
	traces revealed two recurring issues: information		104
	forgetting , where the model repeatedly retrieves		105
	previously seen content, and ineffective informa-		106
	tion extraction , where useful evidence is retrieved		107
	but not incorporated effectively into later reasoning.		108
	We therefore evaluate three test-time modifications		109
	to the Search-R1 pipeline.		110
	3.1 Contextualization		111
	The Contextualization module, shown in Figure 1		112
	(Contextualization), uses an external language		113
	model to extract task-relevant information from re-		114
	trieved documents after each retrieval step. Given		115
	the user prompt p and retrieved documents D_i at		116
	turn i , it produces a concise summary D_i^* contain-		117
	ing only information useful for answering the ques-		118
	tion. This summary is appended to a persistent		119
	memory cache and provided at the next reasoning		120
	step.		121
	The goal is to reduce information forgetting and		122
	improve the use of retrieved evidence without mod-		123
	ifying Search-R1. By passing forward concise rele-		124
	vant context rather than only raw retrieved passages,		125
	the model can reason over both newly retrieved and		126
	previously retained information.		127
	3.2 De-duplication of retrieved documents		128
	To reduce repeated retrieval, we introduce a De-		129
	duplication module that filters out documents seen		130
	in earlier steps. For a given prompt p , we maintain		131
	a set of document IDs retrieved in previous turns.		132
	At each retrieval step, if one of the top- k docu-		133
	ments D_i has already been seen, it is replaced by		134
	the next highest-ranked unseen document, yielding		135
	a set of unseen documents D_i' (Figure 1, Dedupli-		136
	cation).		137
	This intervention is both practical and diagnostic:		138
	it tests whether repeated retrieval reflects a genuine		139
	need to revisit the same evidence or a failure to use		140
	previously retrieved information effectively. If the		141
	same passages are necessary, answer quality should		142
	decline when duplicates are blocked. Otherwise,		143
	stable or improved performance would suggest that		144
	the main issue is not missing evidence, but ineffec-		145
	tive use of evidence already retrieved.		146
	3.3 Hybrid		147
	The Hybrid approach combines Contextualization		148
	and De-duplication to test whether retaining rele-		149

150	vant information while enforcing retrieval diversity	
151	improves both answer quality and retrieval effi-	
152	ciency.	
153	4 Experiments	
154	4.1 Setup	
155	4.1.1 Data source	
156	Search-R1 (Jin et al., 2025) reports performance on	
157	the HotpotQA (Yang et al., 2018) and Natural Ques-	
158	tions (NQ) (Kwiatkowski et al., 2019) datasets.	
159	Since their labeled test sets are not publicly avail-	
160	able, we follow prior work and use the validation	
161	sets. Retrieval is performed on the 2018 Wikipedia	
162	dump with the E5 retriever.	
163	4.1.2 Data splits	
164	To reduce the cost associated with querying ex-	
165	ternal LLMs, we created a smaller subset of	
166	question-answer pairs for evaluation. Specifically,	
167	we randomly sampled 500 question-answer pairs	
168	from the HotpotQA (Yang et al., 2018) and NQ	
169	(Kwiatkowski et al., 2019) validation sets. This	
170	subset is solely used for evaluation; no hyperpa-	
171	rameter tuning or training is performed using this	
172	subset. All reported metrics are based on this vali-	
173	dation set.	
174	4.1.3 Baselines	
175	We utilize the already trained Qwen2.5-7B Search-	
176	R1-base (PPO) as our main baseline for compar-	
177	ison. While running inference with both	
178	Qwen2.5-3B Search-R1-base (PPO) and Qwen2.5-	
179	3B Search-R1-instruct (GRPO), the latter model	
180	exhibited difficulties in adhering to the structured	
181	output format specified by the Search-R1 frame-	
182	work. Inference outputs showed frequent failures	
183	to generate required output tags such as <code><think></code>	
184	and <code><search></code> within the iterative reasoning loop.	
185	Additionally, the model often generated retrieved	
186	information under <code><information></code> tags by itself af-	
187	ter the search query, and also occasionally fails to	
188	produce a final answer at the end of the reason-	
189	ing chain. These behaviors indicate limitations in	
190	the ability to reliably follow instruction-guided for-	
191	matting. Therefore, we only enhance the superior	
192	Qwen2.5-7B Search-R1-base (PPO) model to test	
193	our modules and report the corresponding results	
194	in Table 1. For all approaches, we run inference on	
195	our validation dataset of 500 questions and com-	
196	pute exact match, LLM Match, and the average	
197	number of turns.	
	4.1.4 Implementation details	198
	We built on top of the publicly available Search-	199
	R1 source code on GitHub, where we made mod-	200
	ifications to the model prompt to optimize the	201
	model’s behavior. For inference, we use the Hug-	202
	gingFace Transformers library for the model for-	203
	ward pass and the E5 dense retriever provided in	204
	the Search-R1 GitHub repository over Wikipedia	205
	article chunks. Contextualization and LLM-as-a-	206
	Judge evaluation are performed by GPT-4.1-mini	207
	via the OpenAI API.	208
	4.1.5 Evaluation Metrics	209
	We report Exact Match (EM), as in Search-R1, and	210
	additionally introduce LLM Match together with	211
	the average number of retrieval turns. We include	212
	LLM Match because exact match can produce false	213
	negatives when the predicted answer does not ex-	214
	actly match the golden answer string despite refer-	215
	ring to the same underlying entity, for example “2”	216
	vs. “Two” or “950 Pesos” vs. “P950”.	217
	For LLM Match, an external LLM (GPT-4.1-	218
	mini) is given the predicted answer and a set of	219
	ground-truth answers, and is asked to determine	220
	whether the prediction is semantically equivalent	221
	to any gold answer. Minor differences in phras-	222
	ing are allowed as long as the predicted answer con-	223
	veys the same meaning. The prompt directs the model	224
	to assign a binary score:	225
	<ul style="list-style-type: none"> • 1 if the predicted answer is semantically equiv- 	226
	alent to any ground truth answer, and	227
	<ul style="list-style-type: none"> • 0 if it is incomplete or diverges in meaning. 	228
	This allows us to scale semantic evaluation with-	229
	out human annotation. We also report the average	230
	number of retrievals as a measure of inference effi-	231
	ciency, but interpret it together with answer accu-	232
	racy, since a model could trivially reduce retrieval	233
	count by hallucinating an answer without retriev-	234
	ing evidence.	235
	4.2 Results	236
	In terms of answer accuracy, the Contextualiza-	237
	tion approach achieves a 5.6% increase in EM	238
	and a 6.7% increase in LLM Match score com-	239
	pared with the Search-R1 baseline. In addition,	240
	it is also the most efficient, reducing the average	241
	number of searches to 2.142, compared with the	242
	baseline which has 2.392 searches. While the De-	243
	duplication and Hybrid approaches have similar	244
	gains in EM and LLM Match over the baseline,	245

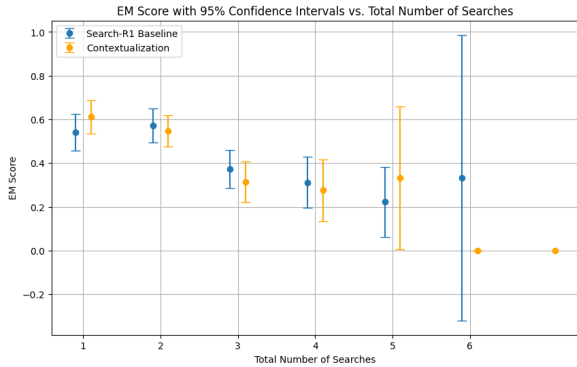


Figure 2: Illustration that questions requiring more agentic turns are inherently more difficult, as shown by the downward trend in Exact Match (EM) score for both the Search-R1 baseline and our Contextualization module. While the Contextualization module achieves a slightly higher mean EM at some points, the overlapping 95% confidence intervals suggest no clear difference between the two approaches at any given search count.

only the Hybrid approach reduces the average number of retrievals, similar to the Contextualization. The De-duplication pipeline is less efficient than the baseline, with 2.498 versus 2.392 average retrievals. Overall, the Contextualization approach still achieves the highest EM, LLM Match, and lowest average number of retrievals. All metrics along with the baseline are shown in Table 1.

Comparing the baseline and De-duplication outputs suggests that repeated retrieval is often not the main bottleneck. When duplicates are filtered out, the model tends to continue issuing similar queries for additional context, increasing retrieval count with little gain in answer accuracy. This suggests that ineffective use of retrieved information is a larger problem than retrieval overlap.

Taken together, these results suggest that the main advantage of Contextualization is not only better answer selection, but also better carryover of useful evidence across reasoning steps. This is consistent with its simultaneous gains in Exact Match, LLM Match, and average retrieval count, whereas De-duplication alone provides smaller accuracy gains while increasing the number of searches. In other words, improving how retrieved evidence is compressed and reused appears more effective than enforcing novelty in the retrieved set alone.

Figure 2 shows a 95% confidence interval around exact match score of the Search-R1 baseline and the Contextualization pipeline, conditioned on the total number of searches performed. While we do not observe a clear separation between the two

Table 1: Performance comparison of our proposed modules against the Search-R1 baselines on our 500-question evaluation set. We report Exact Match, LLM Match, and the average number of turns. Our modules are applied as test-time enhancements to the Qwen2.5-7B base (PPO) model. Cont: Contextualization (Section 3.1). De-Dup: De-duplication (Section 3.2). Hybrid (Section 3.3). Best results for our approaches are in bold.

Variant	Exact Match	LLM Match	avg. # searches
Qwen2.5-3B Search-R1			
base (PPO)	0.292	0.356	1.410
instruct (GRPO)	0.310	0.396	2.054
Qwen2.5-7B Search R1			
base (PPO)	0.464	0.538	2.392
base (PPO) w/ Cont (Ours)	0.490	0.574	2.142
base (PPO) w/ De-Dup (Ours)	0.478	0.560	2.498
base (PPO) w/ Hybrid (Ours)	0.480	0.568	2.154

approaches at any given search count, both show a downward trend. This suggests that Exact Match is negatively correlated with the number of retrievals.

For the Search-R1 baselines and the three non-training approaches, the LLM Match is 16 to 18% greater than the exact match score. Inspecting cases where the LLM judge marks the predicted and gold answers as equivalent but Exact Match fails, we observe two common patterns: numerical answers and abbreviated or shortened names.

5 Conclusion

In this work, we evaluated three test-time modifications to the Search-R1 pipeline: Contextualization, De-duplication, and a Hybrid of both. All three improve answer accuracy over the Search-R1 baseline. Contextualization performs best overall, improving both answer quality and retrieval efficiency, while De-duplication alone increases retrieval count. The Hybrid approach improves both accuracy and efficiency, though less strongly than Contextualization alone.

Our findings suggest that lightweight memory-oriented test-time interventions can yield meaningful gains without retraining the underlying agentic RAG model or retriever.

Limitations

This work has several limitations. First, our evaluation is limited to two open-domain question answering benchmarks, HotpotQA and Natural Questions, using 500 sampled validation examples from each dataset. While these datasets are standard and allow comparison with prior work, they do not cover all retrieval settings, domains, or question types. As a result, the observed gains may not generalize to other tasks such as domain-specific QA, long-document retrieval, or settings with noisier corpora.

Second, our experiments focus on a single Search-R1 backbone, Qwen2.5-7B Search-R1-base (PPO), because smaller variants did not reliably follow the required inference format. This means our conclusions are about the behavior of our proposed test-time modifications on this particular model and setup, rather than about agentic RAG systems in general. Different base models, retrievers, or retrieval depths may interact with our modules differently.

Third, the contextualization module depends on an external proprietary model, GPT-4.1-mini, both for information extraction and for the LLM-as-a-Judge evaluation. This introduces an additional source of variance and potential bias. The quality of contextualization may depend on the external model’s capabilities, prompting sensitivity, and API behavior. Similarly, although LLM-based evaluation helps account for semantically correct answers that exact match misses, it is still an imperfect proxy and may occasionally overestimate or underestimate answer correctness.

Fourth, our efficiency analysis is limited. We report the average number of retrieval turns, which captures part of the inference cost, but it does not fully measure latency, monetary cost, or total token usage. In particular, the contextualization module adds extra model calls, so fewer retrieval turns do not necessarily imply lower overall runtime or cheaper inference in deployment settings.

Finally, our work studies only test-time modifications and does not examine whether similar improvements could be obtained more robustly through training, fine-tuning, or better retriever optimization. A broader evaluation across more datasets, models, and cost metrics would strengthen the conclusions.

References

- Zhiyu An, Xianzhong Ding, Yen-Chun Fu, Cheng-Chung Chu, Yan Li, and Wan Du. 2024. [Golden-retriever: High-fidelity agentic retrieval augmented generation for industrial knowledge base](#). *Preprint*, arXiv:2408.00798.
- Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. 2024. [Rq-rag: Learning to refine queries for retrieval augmented generation](#). *Preprint*, arXiv:2404.00610.
- Jerry Huang, Siddarth Madala, Risham Sidhu, Cheng Niu, Julia Hockenmaier, and Tong Zhang. 2025. [Rag-rl: Advancing retrieval-augmented generation via rl and curriculum learning](#). *Preprint*, arXiv:2503.12759.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. [Search-r1: Training llms to reason and leverage search engines with reinforcement learning](#). *Preprint*, arXiv:2503.09516.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.
- Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025. [Search-o1: Agentic search-enhanced large reasoning models](#). *Preprint*, arXiv:2501.05366.
- Yunxiao Shi, Xing Zi, Zijing Shi, Haimin Zhang, Qiang Wu, and Min Xu. 2024. [Enhancing retrieval and managing retrieval: A four-module synergy for improved quality and efficiency in rag systems](#). *Preprint*, arXiv:2407.10670.
- Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talei Khoei. 2025. [Agentic retrieval-augmented generation: A survey on agentic rag](#). *Preprint*, arXiv:2501.09136.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. [Text embeddings by weakly-supervised contrastive pre-training](#). *arXiv preprint arXiv:2212.03533*.

406 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,
407 William Cohen, Ruslan Salakhutdinov, and Christo-
408 pher D. Manning. 2018. [HotpotQA: A dataset for](#)
409 [diverse, explainable multi-hop question answering](#).
410 In *Proceedings of the 2018 Conference on Empiri-*
411 *cal Methods in Natural Language Processing*, pages
412 2369–2380, Brussels, Belgium. Association for Com-
413 putational Linguistics.