

# Search Wisely: Mitigating Sub-optimal Agentic Searches By Reducing Uncertainty

Anonymous ACL submission

## Abstract

Agentic Retrieval-Augmented Generation (RAG) systems enhance Large Language Models (LLMs) by enabling dynamic, multi-step reasoning and information retrieval. However, these systems often exhibit sub-optimal search behaviors like over-search (retrieving redundant information) and under-search (failing to retrieve necessary information), which hinder efficiency and reliability. This work formally defines and quantifies these behaviors, revealing their prevalence across multiple QA datasets and agentic RAG systems (e.g., one model could have avoided searching in 27.7% of its search steps). Furthermore, we demonstrate a crucial link between these inefficiencies and the models' uncertainty regarding their own knowledge boundaries, where response accuracy correlates with model's uncertainty in its search decisions. To address this, we propose  $\beta$ -GRPO, a reinforcement learning-based training method that incorporates confidence threshold to reward high-certainty search decisions. Experiments on seven QA benchmarks show that  $\beta$ -GRPO enable a 3B model with better agentic RAG ability, outperforming other strong baselines with a 4% higher average exact match score<sup>1</sup>.

## 1 Introduction

Recent advances in Large Language Models (LLMs) have propelled their use in information-intensive tasks such as question answering and knowledge synthesis, especially when paired with retrieval capabilities (Wang et al., 2025b). Agentic Retrieval-Augmented Generation (RAG) frameworks (Jin et al., 2025a; Song et al., 2025a; Chen et al., 2025) push this further by empowering LLMs to perform multi-step reasoning (Li et al., 2025) and dynamically decide when and what to retrieve (Guan et al., 2025), closely emulating sophisticated human research processes. However, despite these

advancements, current agentic RAG systems often struggle with efficiency and reliability due to sub-optimal search behaviors (Shen et al., 2024; Qian et al., 2025; Wang et al., 2025a). In particular, two major challenges: 1) over-search, where the model retrieves information it already knows, and 2) under-search, where it fails to seek external knowledge when necessary, have been identified as critical obstacles that degrade performance.

In this work, we conduct a thorough quantitative analysis to identify and measure the prevalence of over-search and under-search. Our experiments on several multi-hop QA datasets (2Wiki-MultiHopQA (Ho et al., 2020), Bamboogle (Press et al., 2023), HotpotQA (Yang et al., 2018), and MuSiQue (Trivedi et al., 2022)) using contemporary LLMs like R1-Searcher (Song et al., 2025a) and Search-R1 (Jin et al., 2025a) reveal significant instances of sub-optimal search. We also further explore the connection between these behaviors and a model's awareness of its knowledge boundaries, finding that candidate responses generated with higher certainty about the necessity of a search query tend to achieve better accuracy.

To address this, we introduce  $\beta$ -GRPO, a variant of GRPO (Shao et al., 2024) where the confidence of search calls are modeled as the minimal token probability of the search queries produced by the model and a confidence threshold is incorporated into the reward function, only encouraging generations with high-certainty search calls leading to correct answer. Through extensive experiments on seven QA benchmarks, we show that  $\beta$ -GRPO enables a 3B model with better agentic RAG ability compared to strong baselines with a 4% higher average exact match score and 1.21% fewer over-searches and 7.33% fewer under-searches.

## 2 Identifying Sub-optimal Search

To investigate the prevalence of over-search and under-search, we conduct three experiments with

<sup>1</sup>We will release all our codes and data upon acceptance.

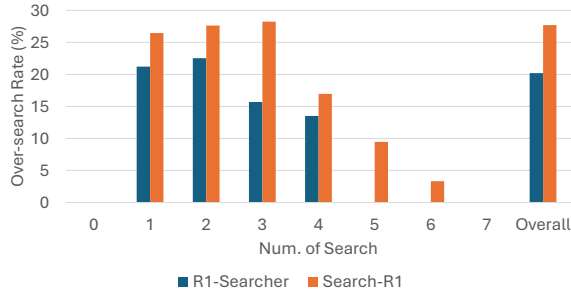


Figure 1: Percentage for all search steps that can be answered without performing searches of R1-Searcher and Search-R1 on 4 datasets combined, with respect to the number of searches of each test sample.

the test sets of four widely recognized multihop QA datasets: 2WikiMultiHopQA (Ho et al., 2020), Bamboogle (Press et al., 2023), HotpotQA (Yang et al., 2018), and MuSiQue (Trivedi et al., 2022). We mainly investigate two recent LLMs that interact with search engines: R1-Searcher (Song et al., 2025b) and Search-R1 (Jin et al., 2025b). We adopt the version trained based on Qwen2.5-7B (Qwen et al., 2025) for a fair comparison.

## 2.1 Step-wise Analysis

To directly measure whether a search step was truly necessary, we separate all outputs into individual steps and identify if each of them matches with the definition of over-search and under-search. For over-search rate measurement, we prompted the model to answer sub-queries from all the steps with search behavior using only their internal knowledge and the preceding context. For under-search rate measurement, we examine steps without searching and evaluate the correctness of the generated information. A detailed explanation of the analysis pipeline is provided in Appendix A.2.

**Capability to Answer from Memory** The results in Figure 1 show that a significant portion of search actions were instances of over-search. R1-Searcher could have answered correctly without searching in 20.2% of its search steps overall, while Search-R1 could have done so in 27.7% of its search steps. This highlights a substantial room for efficiency improvement. Figure 1 also shows the over-search rate for each subset of test samples grouped by the total number of search steps an agent used to solve an entire problem instance. The results per each subset indicates that over-search is a persistent issue irrespective of the overall search complexity adopted by the model for a given problem. Despite the step-wise analy-

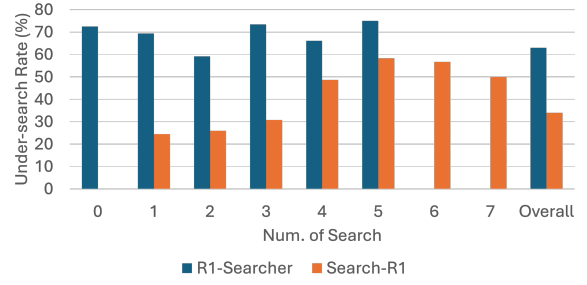


Figure 2: Error rate for all non-search steps of R1-Searcher and Search-R1 on 4 datasets combined, with respect to the number of searches of each test sample.

sis, we also conduct an analysis on comparing the number of searches versus the pre-given number of hops from the dataset in Appendix A.3, which also supports our conclusion.

**Error Rate in Non-Search Steps** Figure 2 analyzes the error rate in non-search steps, which can be seen as the rate of under-search. Both models exhibited high error rates (R1-Searcher: 63%, Search-R1: 33.98%) in non-search steps, suggesting a strong tendency towards under-search leading to incorrect reasoning or hallucination. For R1-Searcher, this error rate was particularly high with fewer total searches (over 72% if no searches were made). For Search-R1, errors in non-search steps remained notable even when performing many searches overall (e.g., 48.70% for 4-search problems), possibly due to decision complexity in later stages. (See Figure 2 for detailed error rates by search step count).

## 2.2 Sub-optimal Search & Knowledge Boundary

The observed tendencies towards over-search and under-search, combined with our definition, suggest a core deficiency in how agentic RAG models perceive knowledge boundaries—the limits of what they know versus what they need to find out. To illustrate the link between better knowledge boundary awareness and improved outcomes, we analyze the performance of 4 Qwen2.5-3B based Search-R1 models (including PPO and GRPO trained, Base and Instruct variants). We generate 5 candidate responses for each question and group these responses based on each output’s **minimum probabilities within all the search query tokens** as the indication of certainty on knowledge boundary.

As shown in Table 1, candidate responses generated with lower intrinsic uncertainty generally lead to higher final accuracy (as high as 6% on Bam-

Model Config	Prob. Group	2Wiki	Bamboogle	HotpotQA	Musique
Base + PPO	Max	<b>0.184</b>	0.096	<b>0.152</b>	0.038
	Min	0.168	0.096	0.114	0.038
Base + GRPO	Max	<b>0.249</b>	0.112	<b>0.327</b>	<b>0.085</b>
	Min	0.234	0.104	0.289	0.056
Instruct + PPO	Max	<b>0.333</b>	0.250	0.262	<b>0.138</b>
	Min	0.297	0.250	0.262	0.116
Instruct + GRPO	Max	0.402	<b>0.125</b>	<b>0.343</b>	0.116
	Min	0.402	0.063	0.302	0.116

Table 1: Cover EM scores on multi-hop QA datasets, comparing groups of responses with higher vs. lower uncertainty (derived from average of minimum probability of search query tokens) on knowledge boundary. Bold indicates instances where the Max Prob. group achieved a strictly better performance.

boogle and 3.8% on HotpotQA), across different training methods and base models. This suggests that when the model exhibits higher confidence (lower uncertainty) in its generation path, it is more likely to be on a correct trajectory. Therefore, improving an agent’s ability to accurately gauge its internal knowledge state—effectively sharpening its knowledge boundary detection and reducing undue uncertainty—is a crucial step towards mitigating both over-search and under-search, thereby enhancing the overall efficiency and reliability of agentic RAG systems. Our approach is motivated by this principle, aiming to train agents to better assess and reduce uncertainty at each search decision.

### 3 Approach

Current RL powered agentic RAG methods (Jin et al., 2025a; Song et al., 2025a; Chen et al., 2025) do not explicitly model the knowledge self-awareness during the training process, resulting in generations with low confidence, which are not desired and shown to easily contain wrong answer compared to generations with higher confidence (Table 1). To this end, we propose a simple yet effective variant of GRPO (Shao et al., 2024),  $\beta$ -GRPO, which leverages the uncertainty of the search query spans for more effective rewarding and training.

**Agentic RAG with RL (Search-R1 (Jin et al., 2025a))** Given a question, we prompt the policy model to explicitly reason enclosed within `<think></think>` tags about whether to use an off-the-shelf search tool, and, if so, to generate a search query within `<search></search>` tags. The search tool then returns relevant documents inside `<information></information>` tags. Once obtaining new information, the policy model can either continue searching for addi-

tional information or provide a final answer within `<answer></answer>` tags. The instruction given to the policy model could be found in Appendix A.4. If the final answer match the groundtruth, the response will be given a reward 1, otherwise 0. And the policy are updated via policy gradient methods like GRPO (Shao et al., 2024).

**$\beta$ -GRPO** Motivated by the observation that rollouts with low-confidence search calls are more likely to be incorrect, we incorporate model confidence into the RL reward process. Specifically, for each rollout containing search calls (enclosed within `<search></search>` tags), we extract the probabilities of the search tokens including the tags and use the minimum probability among them as a measure of the model confidence for the search calls within a rollout (Jiang et al., 2023). We then set a confidence threshold  $\beta$ : only rollouts with the confidence of search calls (if exist) above  $\beta$  and correct answers receive a reward of 1, otherwise 0.

## 4 Experiments

**Datasets** We follow Search-R1 (Jin et al., 2025a) using a mixture of the NQ (Kwiatkowski et al., 2019) and HotpotQA (Yang et al., 2018) training sets for model training. For evaluation, we consider seven QA benchmarks, including general QA datasets, NQ, TriviaQA (Joshi et al., 2017), and PopQA (Mallen et al., 2023), as well as multi-hop QA datasets: HotpotQA, 2WikiMultiHopQA (Ho et al., 2020), Bamboogle (Press et al., 2023), and MuSiQue (Trivedi et al., 2022). Exact match (EM) is used as our main evaluation metric.

**Baselines** We compare our method with several baselines: methods that do not use a retriever including direct prompting, Chain-of-Thought (CoT) (Wei et al., 2022) prompting, supervised fine-tuning (SFT) (Chung et al., 2022), and reinforcement learning-based fine-tuning (R1) (DeepSeek-AI et al., 2025); methods that use a retriever but do not perform agentic retrieval, such as Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) and IRCOT (Trivedi et al., 2023); and finally, agentic retrieval methods, including Search-o1 (Li et al., 2025) and Search-R1 (Jin et al., 2025a).

Based on our preliminary experiments, we found that training the policy model from scratch using our confidence-based rewards prevents it from learning effective search behavior. Therefore, we use Qwen2.5-3B (Qwen et al., 2025) and initialize it with the parameters from Search-R1. Then

Methods	General QA				Multi-Hop QA			
	NQ <sup>†</sup>	TriviaQA <sup>*</sup>	PopQA <sup>*</sup>	HotpotQA <sup>†</sup>	2wiki <sup>*</sup>	Musique <sup>*</sup>	Bamboogle <sup>*</sup>	Average
Direct Prompting	0.106	0.288	0.108	0.149	0.244	0.020	0.024	0.134
CoT Prompting	0.023	0.032	0.005	0.021	0.021	0.002	0.000	0.015
IRCoT	0.111	0.312	0.200	0.164	0.171	0.067	0.240	0.181
Search-o1	0.238	0.472	0.262	0.221	0.218	0.054	<b>0.320</b>	0.255
RAG	0.348	0.544	0.387	0.255	0.226	0.047	0.080	0.270
SFT	0.249	0.292	0.104	0.186	0.248	0.044	0.112	0.176
R1	0.226	0.455	0.173	0.201	0.268	0.055	0.224	0.229
Search-R1	0.406	0.587	0.435	0.284	0.273	0.049	0.088	0.303
Search-R1-GRPO	0.432	0.578	0.413	0.294	0.271	0.067	0.112	0.309
Search-R1- $\beta$ -GRPO (ours)	<b>0.468</b>	<b>0.625</b>	<b>0.449</b>	<b>0.334</b>	<b>0.304</b>	<b>0.086</b>	0.144	<b>0.344</b>

Table 2: Main results. The best performance is set in bold. <sup>†</sup>/<sup>\*</sup> represents in-domain/out-domain datasets.

Search-R1-GRPO	Search-R1- $\beta$ -GRPO (ours)
<b>Question:</b> In what Country is Sul America Esporte Clube in? <b>Generation:</b> <think> The question asks for the country hosting Sul America Esporte Clube, a team in the Brazilian C1 and C2 league, thus the direct reasoning is likely to be Brazilian football team residing in Brazil.	<b>Question:</b> In what Country is Sul America Esporte Clube in? <b>Generation:</b> <search> in what Country is Sul America Esporte Clube in </search> <information> Doc 1: ... Doc 2: ... Doc 3: ... </information> <answer> Brazil </answer> (Confidence: 0.99)

Table 3: Case Study. The retrieved documents are omitted for clarity.

we continue training using GRPO with different reward functions: one using the original answer-based reward (Search-R1-GRPO), and the other using our proposed confidence-based reward (Search-R1- $\beta$ -GRPO). We set the value of  $\beta$  as 0.4 according to the analysis in Section 5. Detailed training configurations could be found in Appendix A.5.

**Results** As shown in Table 2, agentic search with RL training (Search-R1\*) significantly outperforms other baselines, indicating that incorporating search through autonomous reasoning and RL training is more effective than non-agentic or prompting methods. Our model, Search-R1- $\beta$ -GRPO, achieves the highest overall average EM score across the datasets. Figure 3 in Appendix A.5 shows the training rewards for Search-R1-GRPO and Search-R1- $\beta$ -GRPO. We observe that the rewards for Search-R1-GRPO fluctuate and do not show clear improvement over training steps. In contrast, Search-R1- $\beta$ -GRPO achieves higher and more stable rewards. This improved performance suggests that our proposed reward assignment based on the confidence of search calls within a rollout is effective.

## 5 Analysis

**Ablation on  $\beta$  & Case Study** Following Jiang et al. (2023), we experiment with three confidence threshold values: 0.2, 0.4, and 0.6. The average EM scores are 0.341, 0.344 and 0.336 with a threshold of 0.4 yields the best result. Moreover, we find 115 test cases from the multi-hop

QA datasets where Search-R1- $\beta$ -GRPO produces a correct answer with higher confidence, while Search-R1-GRPO gives an incorrect answer. These cases clearly benefit from the increased model confidence enabled by the proposed  $\beta$ -GRPO. An example is shown in Table 3: Search-R1-GRPO lacks confidence and fails to provide a definite answer, whereas Search-R1- $\beta$ -GRPO generates a confident search query and produces the correct answer.

**Under-searches & Over-searches** We also measure the rate of over-search and under-search of our Search-R1- $\beta$ -GRPO and the baseline Search-R1-GRPO trained based on Qwen2.5-3B with the methods in Section 2.1. Compared with Search-R1-GRPO, which has overall 21.10% over-search rate and 42.04% under-search rate%, our Search-R1- $\beta$ -GRPO achieves 19.89% over-search rate and 34.71% under-search rate, which are lower than the baseline method. This shows that our method effectively reduces both types of sub-optimal searches.

## 6 Conclusion

In this work, we formally define and quantify sub-optimal search behaviors, over-search and under-search, in agentic RAG systems, revealing their prevalence and impact. By introducing  $\beta$ -GRPO, a confidence-aware policy gradient method, we enable a 3B model with better agentic RAG ability than strong baselines.



## 7 Limitations

We formally define and quantify sub-optimal search behaviors in agentic RAG systems and propose  $\beta$ -GRPO to train agentic RAG models with improved self-knowledge awareness. However, we acknowledge that sub-optimal search behaviors, over-search and under-search, are persistent challenges that require further investigation, especially in more open-ended tasks like deep research (Alzubi et al., 2025). Additionally, due to limited computational resources, we are unable to train larger models and leave it for future work.

## References

- Salaheddin Alzubi, Creston Brooks, Purva Chiniya, Edoardo Contente, Chiara von Gerlach, Lucas Irwin, Yihan Jiang, Arda Kaz, Windsor Nguyen, Sewoong Oh, Himanshu Tyagi, and Pramod Viswanath. 2025. Open deep search: Democratizing search with open-source reasoning agents. *arXiv [cs.LG]*.
- Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Fan Yang, Zenan Zhou, Weipeng Chen, Haofen Wang, Jeff Z Pan, Wen Zhang, and Huajun Chen. 2025. ReSearch: Learning to reason with search for LLMs via reinforcement learning. *arXiv [cs.AI]*.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. *arXiv [cs.LG]*.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z F Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J L Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojuan Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qishi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R J Chen, R L Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S S Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W L Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X Q Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y K Li, Y Q Wang, Y X Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y X Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z Z Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. 2025. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv [cs.CL]*.
- Xinyan Guan, Jiali Zeng, Fandong Meng, Chunlei Xin, Yaojie Lu, Hongyu Lin, Xianpei Han, Le Sun, and Jie Zhou. 2025. DeepRag: Thinking to retrieval step by step for large language models.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. *arXiv [cs.CL]*.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. 2025a. Search-R1: Training LLMs to reason and leverage search engines with reinforcement learning. *arXiv [cs.CL]*.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025b. Search-r1: Training llms to reason and leverage search engines with reinforcement learning.

414	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke	Chak Li, Chan Jun Shern, Channing Conger, Char-	472
415	Zettlemoyer. 2017. TriviaQA: A large scale distantly	lotte Barette, Chelsea Voss, Chen Ding, Cheng Lu,	473
416	supervised challenge dataset for reading comprehen-	Chong Zhang, Chris Beaumont, Chris Hallacy, Chris	474
417	sion.	Koch, Christian Gibson, Christina Kim, Christine	475
418	Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick	Choi, Christine McLeavey, Christopher Hesse, Clau-	476
419	Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and	dia Fischer, Clemens Winter, Coley Czarnecki, Colin	477
420	Wen-Tau Yih. 2020. Dense passage retrieval for open-	Jarvis, Colin Wei, Constantin Koumouzelis, Dane	478
421	domain question answering. <i>arXiv [cs.CL]</i> .	Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy,	479
422	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	David Carr, David Farhi, David Mely, David Robin-	480
423	field, Michael Collins, Ankur Parikh, Chris Alberti,	son, David Sasaki, Denny Jin, Dev Valladares, Dim-	481
424	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	itris Tsipras, Doug Li, Duc Phong Nguyen, Duncan	482
425	ton Lee, Kristina Toutanova, Llion Jones, Matthew	Findlay, Edede Oiwoh, Edmund Wong, Ehsan As-	483
426	Kelcey, Ming-Wei Chang, Andrew M Dai, Jakob	dar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow,	484
427	Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural	Eric Kramer, Eric Peterson, Eric Sigler, Eric Wal-	485
428	questions: A benchmark for question answering re-	lace, Eugene Brevdo, Evan Mays, Farzad Khorasani,	486
429	search. <i>Trans. Assoc. Comput. Linguist.</i> , 7:453–466.	Felipe Petroski Such, Filippo Raso, Francis Zhang,	487
430	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	Fred von Lohmann, Freddie Sulit, Gabriel Goh,	488
431	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	Gene Oden, Geoff Salmon, Giulio Starace, Greg	489
432	rich Küttler, Mike Lewis, Wen-Tau Yih, Tim Rock-	Brockman, Hadi Salman, Haiming Bao, Haitang	490
433	täschel, Sebastian Riedel, and Douwe Kiela. 2020.	Hu, Hannah Wong, Haoyu Wang, Heather Schmidt,	491
434	Retrieval-augmented generation for knowledge-	Heather Whitney, Heewoo Jun, Hendrik Kirchner,	492
435	intensive NLP tasks. <i>arXiv [cs.CL]</i> .	Henrique Ponde de Oliveira Pinto, Hongyu Ren,	493
436	Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang,	Huiwen Chang, Hyung Won Chung, Ian Kivlichan,	494
437	Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng	Ian O’Connell, Ian O’Connell, Ian Osband, Ian Sil-	495
438	Dou. 2025. Search-o1: Agentic search-enhanced	ber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya	496
439	large reasoning models. <i>arXiv [cs.AI]</i> .	Kostrikov, Ilya Sutskever, Ingmar Kanitscheider,	497
440	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,	Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub	498
441	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.	Pachocki, James Aung, James Betker, James Crooks,	499
442	When not to trust language models: Investigating	James Lennon, Jamie Kiros, Jan Leike, Jane Park,	500
443	effectiveness of parametric and non-parametric mem-	Jason Kwon, Jason Phang, Jason Teplitz, Jason	501
444	ories. In <i>Proceedings of the 61st Annual Meeting of</i>	Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Var-	502
445	<i>the Association for Computational Linguistics (Vol-</i>	avva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui	503
446	<i>ume 1: Long Papers)</i> , Stroudsburg, PA, USA. Asso-	Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang,	504
447	ciation for Computational Linguistics.	Joaquin Quinonero Candela, Joe Beutler, Joe Lan-	505
448	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher,	ders, Joel Parish, Johannes Heidecke, John Schul-	506
449	Adam Perelman, Aditya Ramesh, Aidan Clark,	man, Jonathan Lachman, Jonathan McKay, Jonathan	507
450	AJ Ostrow, Akila Welihinda, Alan Hayes, Alec	Uesato, Jonathan Ward, Jong Wook Kim, Joost	508
451	Radford, Aleksander Mądry, Alex Baker-Whitcomb,	Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross,	509
452	Alex Beutel, Alex Borzunov, Alex Carney, Alex	Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao,	510
453	Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex	Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai	511
454	Renzin, Alex Tachard Passos, Alexander Kirillov,	Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kevin	512
455	Alexi Christakis, Alexis Conneau, Ali Kamali, Allan	Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu,	513
456	Jabri, Allison Moyer, Allison Tam, Amadou Crookes,	Kenny Nguyen, Keren Gu-Lemberg, Kevin Button,	514
457	Amin Tootoochian, Amin Tootoonchian, Ananya	Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle	515
458	Kumar, Andrea Vallone, Andrej Karpathy, Andrew	Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lau-	516
459	Braunstein, Andrew Cann, Andrew Codispoti, An-	ren Workman, Leher Pathak, Leo Chen, Li Jing, Lia	517
460	drew Galu, Andrew Kondrich, Andrew Tulloch, An-	Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lil-	518
461	drey Mishchenko, Angela Baek, Angela Jiang, An-	ian Weng, Lindsay McCallum, Lindsey Held, Long	519
462	toine Pelisse, Antonia Woodford, Anuj Gosalia, Arka	Ouyang, Louis Feувrier, Lu Zhang, Lukas Kon-	520
463	Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver,	draciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz,	521
464	Barret Zoph, Behrooz Ghorbani, Ben Leimberger,	Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine	522
465	Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin	Boyd, Madeleine Thompson, Marat Dukhan, Mark	523
466	Zweig, Beth Hoover, Blake Samic, Bob McGrew,	Chen, Mark Gray, Mark Hudnall, Marvin Zhang,	524
467	Bobby Spero, Bogo Giertler, Bowen Cheng, Brad	Marwan Aljubeih, Mateusz Litwin, Matthew Zeng,	525
468	Lightcap, Brandon Walkin, Brendan Quinn, Brian	Max Johnson, Maya Shetty, Mayank Gupta, Meghan	526
469	Guarraci, Brian Hsu, Bright Kellogg, Brydon East-	Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao	527
470	man, Camillo Lugaresi, Carroll Wainwright, Cary	Zhong, Mia Glaese, Mianna Chen, Michael Jan-	528
471	Bassin, Cary Hudson, Casey Chu, Chad Nelson,	ner, Michael Lampe, Michael Petrov, Michael Wu,	529
		Michele Wang, Michelle Fradin, Michelle Pokrass,	530
		Miguel Castro, Miguel Oom Temudo de Castro,	531
		Mikhail Pavlov, Miles Brundage, Miles Wang, Mi-	532
		nal Khan, Mira Murati, Mo Bavarian, Molly Lin,	533
		Murat Yesildal, Nacho Soto, Natalia Gimelshein, Na-	534

535	talie Cone, Natalie Staudacher, Natalie Summers,	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	596
536	Natan LaFontaine, Neil Chowdhury, Nick Ryder,	Junxiao Song, Mingchuan Zhang, Y K Li, Y Wu, and	597
537	Nick Stathas, Nick Turley, Nik Tezak, Niko Felix,	Daya Guo. 2024. DeepSeekMath: Pushing the limits	598
538	Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel	of mathematical reasoning in open language models.	599
539	Bundick, Nora Puckett, Ofir Nachum, Ola Okelola,	<i>arXiv [cs.CL]</i> .	600
540	Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins,		
541	Olivier Godement, Owen Campbell-Moore, Patrick	Yuanhao Shen, Xiaodan Zhu, and Lei Chen. 2024.	601
542	Chao, Paul McMillan, Pavel Belov, Peng Su, Pe-	<b>SMARTCAL: An approach to self-aware tool-use</b>	602
543	ter Bak, Peter Bakkum, Peter Deng, Peter Dolan,	<b>evaluation and calibration.</b> In <i>Proceedings of the</i>	603
544	Peter Hoeschele, Peter Welinder, Phil Tillet, Philip	<i>2024 Conference on Empirical Methods in Natural</i>	604
545	Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming	<i>Language Processing: Industry Track</i> , pages 774–	605
546	Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Ra-	789, Miami, Florida, US. Association for Computa-	606
547	jan Troll, Randall Lin, Rapha Gontijo Lopes, Raul	tional Linguistics.	607
548	Puri, Reah Miyara, Reimar Leike, Renaud Gaubert,		
549	Reza Zamani, Ricky Wang, Rob Donnelly, Rob	Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen,	608
550	Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchan-	Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-	609
551	dani, Romain Huet, Rory Carmichael, Rowan Zellers,	Rong Wen. 2025a. R1-searcher: Incentivizing the	610
552	Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan	search capability in LLMs via reinforcement learning.	611
553	Cheu, Saachi Jain, Sam Altman, Sam Schoenholz,	<i>arXiv [cs.AI]</i> .	612
554	Sam Toizer, Samuel Miserendino, Sandhini Agar-		
555	wal, Sara Culver, Scott Ethersmith, Scott Gray, Sean	Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen,	613
556	Grove, Sean Metzger, Shamez Hermani, Shantanu	Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-	614
557	Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shi-	Rong Wen. 2025b. <b>R1-searcher: Incentivizing the</b>	615
558	rong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay,	<b>search capability in llms via reinforcement learning.</b>	616
559	Srinivas Narayanan, Steve Coffey, Steve Lee, Stew-		
560	art Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	617
561	Xu, Tarun Gogineni, Taya Christianson, Ted Sanders,	and Ashish Sabharwal. 2022. <b>MuSiQue: Multi-</b>	618
562	Tejal Patwardhan, Thomas Cunningham, Thomas	<b>hop questions via single-hop question composition.</b>	619
563	Degry, Thomas Dimson, Thomas Raoux, Thomas	<i>Transactions of the Association for Computational</i>	620
564	Shadwell, Tianhao Zheng, Todd Underwood, Todor	<i>Linguistics</i> , 10:539–554.	621
565	Markov, Toki Sherbakov, Tom Rubin, Tom Stasi,		
566	Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	622
567	Walters, Tyna Eloundou, Valerie Qi, Veit Moeller,	and Ashish Sabharwal. 2023. Interleaving retrieval	623
568	Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne	with chain-of-thought reasoning for knowledge-	624
569	Chang, Weiye Zheng, Wenda Zhou, Wesam Manassra,	intensive multi-step questions. In <i>Proceedings of the</i>	625
570	Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian,	<i>61st Annual Meeting of the Association for Computa-</i>	626
571	Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen	<i>tional Linguistics (Volume 1: Long Papers)</i> , Strouds-	627
572	He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury	burg, PA, USA. Association for Computational Lin-	628
573	Malkov. 2024. <b>Gpt-4o system card.</b>	guistics.	629
574			
575	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt,	Hongru Wang, Cheng Qian, Wanjun Zhong, Xiusi Chen,	630
576	Noah Smith, and Mike Lewis. 2023. <b>Measuring and</b>	Jiahao Qiu, Shijue Huang, Bowen Jin, Mengdi Wang,	631
577	<b>narrowing the compositionality gap in language mod-</b>	Kam-Fai Wong, and Heng Ji. 2025a. <b>Otc: Optimal</b>	632
578	<b>els.</b> In <i>Findings of the Association for Computational</i>	<b>tool calls via reinforcement learning.</b>	633
579	<i>Linguistics: EMNLP 2023</i> , pages 5687–5711, Singa-		
580	pore. Association for Computational Linguistics.	Liang Wang, Haonan Chen, Nan Yang, Xiaolong Huang,	634
581		Zhicheng Dou, and Furu Wei. 2025b. <b>Chain-of-</b>	635
582	Cheng Qian, Emre Can Acikgoz, Hongru Wang, Xiusi	<b>retrieval augmented generation.</b>	636
583	Chen, Avirup Sil, Dilek Hakkani-Tür, Gokhan Tur,		
584	and Heng Ji. 2025. <b>Smart: Self-aware agent for tool</b>	Liang Wang, Nan Yang, Xiaolong Huang, Binxing	637
585	<b>overuse mitigation.</b>	Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,	638
586		and Furu Wei. 2022. Text embeddings by weakly-	639
587	Qwen, :, An Yang, Baosong Yang, Beichen Zhang,	supervised contrastive pre-training. <i>arXiv [cs.CL]</i> .	640
588	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,		
589	Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin,	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	641
590	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang,	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and	642
591	Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang,	Denny Zhou. 2022. Chain-of-thought prompting	643
592	Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li,	elicits reasoning in large language models. <i>arXiv</i>	644
593	Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji	<i>[cs.CL]</i> .	645
594	Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang		
595	Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,	646
	Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru	William Cohen, Ruslan Salakhutdinov, and Christo-	647
	Zhang, and Zihan Qiu. 2025. <b>Qwen2.5 technical</b>	pher D. Manning. 2018. <b>HotpotQA: A dataset for</b>	648
	<b>report.</b>	<b>diverse, explainable multi-hop question answering.</b>	649
		In <i>Proceedings of the 2018 Conference on Empiri-</i>	650
		<i>cal Methods in Natural Language Processing</i> , pages	651

652 2369–2380, Brussels, Belgium. Association for Com-  
653 putational Linguistics.



## A Appendix

### A.1 Formal Definition of Under-search & Over-search

Formally, let an LLM agent’s interaction for a question be a sequence of steps  $T = \{s_1, s_2, \dots, s_N\}$ . Each step  $s_t$  comprises a reasoning component  $r_t$ . If the model decides to retrieve information, the retrieval step  $s_t^R = (r_t, q_t, c_t)$  includes a search sub-query  $q_t$  and the retrieved context  $c_t = \text{search}(q_t)$ . The sub-answer  $a_t$  for this step  $s_t^R$  is typically derived using  $c_t$  and reflected in  $r_{t+1}$ . If the model does not retrieve, the non-retrieval step  $s_t^{NR} = (r_t)$  relies on the existing context  $\{s_1, s_2, \dots, s_{t-1}\}$  and the model’s internal knowledge  $M$  to derive  $a_t$  reflected in  $r_t$ . Let  $a_t^*$  be the ground-truth answer step  $s_t$ . Over-search occurs if a retrieval step  $s_t^R$ ’s answer  $a_t$  could have been derived from  $M$  and  $\{s_1, s_2, \dots, s_{t-1}\}$  only. Under-search occurs if a non-retrieval step  $s_t^{NR}$  leads to  $a_t \neq a_t^*$ .

### A.2 Detailed Step-wise Analysis Procedure

To empirically measure the rates of over-search and under-search, we conducted a detailed step-wise analysis of the agent’s decision-making process. The interactions of the agent are logged as a sequence of steps, where each step can involve internal reasoning (thinking), querying a search tool, processing retrieved context, and generating an answer. We define specific procedures to identify and quantify each type of sub-optimal search behavior:

- 1. Step Extraction:** We parse the agent’s interaction log following the definition in appendix A.1. Each distinct thinking process is a decision point and considered a step, typically delineated by `<step>` and `</step>` tags (or a similar structured logging format). A "search step" is identified as any step where all three relevant operations—think (the model’s reasoning), search (the search query issued), context (the information retrieved). A "non-search step" typically only consists of thinking. In this work specifically, the Step Extraction is done by prompting QwQ-32B (Qwen et al., 2025) as we discover that reasoning LLM typically perform better on such task.
- 2. Extraction of Partial Input:** For each identified search step, we reconstruct the input that would have been available to the model before it decided to search. This is achieved by taking the complete output generated by the

agent from the beginning of the interaction up to and including the content of the think field of the current search step.

- 3. Querying with Internal Knowledge for Over-search Analysis:** For over-search rate measurement, the extracted partial output is then appended with a specific instructional prompt: "I will use my own knowledge to answer this query and provide my answer to this query enclosed in `<query_answer>` `</query_answer>` tags." This combined text serves as a new input to the original RL-tuned model (e.g., Search-R1- $\beta$ -GRPO and Search-R1-GRPO), which is tasked with generating an answer without performing any new search. The over-search rate is then measured by computing the percentage of steps that provide equivalent answer (determined by QwQ-32B in our analysis) for both with and without searching, among all "search steps".
- 4. Generation of Reference Answer for Under-search Analysis:** For each identified non-search step, the original query or sub-query that the agent was attempting to answer at that point is presented to a more powerful, state-of-the-art language model (e.g., ChatGPT-4o (OpenAI et al., 2024)) with recent knowledge cutoff date. This model generates a "reference answer," which is assumed to be of high quality. The reference answer obtained is compared with the actual answer generated by the agent for that non-search step. The under-search rate is calculated as the proportion of non-search steps where the agent’s answer does not match (determined by QwQ-32B in our analysis) the reference answer, quantifying how often the agent fails to search when doing so would have likely led to a more accurate or complete answer.

### A.3 Search Frequency vs. Optimal Hops

One indicator of potential over-search is when the number of search queries generated by an agent exceeds the optimal number of reasoning hops required to answer a question. A significantly higher search count often points to redundant information gathering. For this experiment, we only use the test set from Bamboogle (Press et al., 2023) and MuSiQue (Trivedi et al., 2022) as they are the only two datasets providing pre-defined number of hops for each test sample.

Model	Dataset	Search vs. Hops	Correct (%)	Incorrect (%)	Sum (%)
R1-Searcher	Musique	Less	2.8	19	21.8
		Match	21.8	45.8	67.6
		More	1.8	8.8	10.6
R1-Searcher	Bamboogle	Less	0	0	0
		Match	40.8	52.8	93.6
		More	3.2	3.2	6.4
Search-R1	Musique	Less	1.8	7	8.8
		Match	12.4	27.6	40
		More	8.8	42.4	51.2
Search-R1	Bamboogle	Less	0.8	1.6	2.4
		Match	28.8	28	56.8
		More	12	28.8	40.8

Table 4: Comparison of the number of searches generated vs. annotated hops on Bamboogle and Musique datasets. "More" indicates potential over-search as number of searchers exceeds pre-defined optimal hops. "Less" may indicate a potential under-search.

R1-Searcher exhibits a tendency to perform more searches than hops in 10.6% of Musique cases and 6.4% of Bamboogle cases. Search-R1 shows a more pronounced tendency, with 51.2% (Musique) and 40.8% (Bamboogle) of cases issuing more searches than annotated hops. This result suggests that models trained with different methods do not inherently solve over-search and might even exacerbate it under certain configurations if not properly guided. While "Less" searches than hops might indicate efficient reasoning or under-search, the "More" category strongly suggests instances of over-searching.

#### A.4 Instruction

Answer the given question. You must conduct reasoning inside `<think>` and `</think>` first every time you get new information. After reasoning, if you find you lack some knowledge, you can call a search engine by `<search> query </search>`, and it will return the top searched results between `<information>` and `</information>`. You can search as many times as you want. If you find no further external knowledge needed, you can directly provide the answer inside `<answer>` and `</answer>` without detailed illustrations. For example, `<answer> Beijing </answer>`. Question: question.

#### A.5 Training Configuration & Rewards

We train Search-R1-GPRO and Search-R1- $\beta$ -GPRO for 200 steps, with a learning rate of 1e-6 and batch size of 512. For a question, we produce 5 generations with temperature of 1 to form a GPRO group. For the search engine, for fair compari-

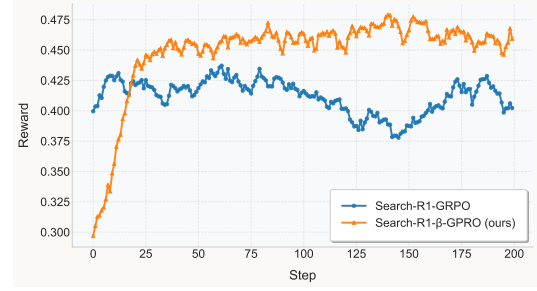


Figure 3: Training Rewards for Search-R1-GPRO and Search-R1- $\beta$ -GPRO.

son, we also use 2018 Wikipedia dump (Karpukhin et al., 2020) as the knowledge source and E5 (Wang et al., 2022) as the retriever as Search-R1 and for each search query, top-3 documents are returned. Our training are conducted on two A100 GPUs.