

Adaptive-RAG: Learning to Adapt Retrieval-Augmented Large Language Models through Question Complexity

Anonymous ACL submission

Abstract

Retrieval-Augmented Large Language Models (LLMs), which incorporate the non-parametric knowledge from external knowledge bases into LLMs, have emerged as a promising approach to enhancing response accuracy in several tasks, such as Question-Answering (QA). However, even though there are various approaches dealing with queries of different complexities, they either handle simple queries with unnecessary computational overhead or fail to adequately address complex multi-step queries; yet, not all user requests fall into only one of the simple or complex categories. In this work, we propose a novel adaptive QA framework, that can dynamically select the most suitable strategy for (retrieval-augmented) LLMs from the simplest to the most sophisticated ones based on the query complexity. Also, this selection process is operationalized with a classifier, which is a smaller LM trained to predict the complexity level of incoming queries with automatically collected labels, obtained from actual predicted outcomes of models and inherent inductive biases in datasets. This approach offers a balanced strategy, seamlessly adapting between the iterative and single-step retrieval-augmented LLMs, as well as the no-retrieval methods, in response to a range of query complexities. We validate our model on a set of open-domain QA datasets, covering multiple query complexities, and show that ours enhances the overall efficiency and accuracy of QA systems, compared to relevant baselines including the adaptive retrieval approaches.

1 Introduction

Recent Large Language Models (LLMs) (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023; Anil et al., 2023) have shown overwhelming performances across diverse tasks, including question-answering (QA) (Yang et al., 2018; Kwiatkowski et al., 2019). However, they still generate factually incorrect answers since their knowledge solely

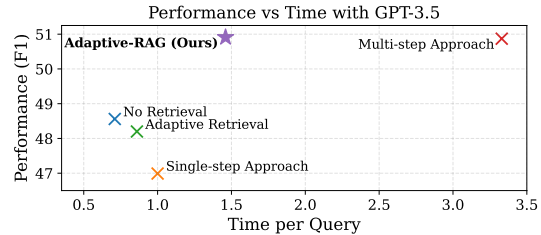


Figure 1: QA performance (F1) and efficiency (Time/Query) for different retrieval-augmented generation approaches. We use the GPT-3.5-Turbo-Instruct as the base LLM.

relies on their parametric memory (Kasai et al., 2022; Mallen et al., 2023). Meanwhile, memorizing all the (ever-changing) world knowledge may not be possible. To address this problem, retrieval-augmented LLMs (Borgeaud et al., 2022; Izacard et al., 2023; Shi et al., 2023), which incorporate non-parametric knowledge into LLMs with additional retrieval modules, have gained increasing attention. Specifically, these models access a knowledge base, serving as an extensive repository of information across various subjects and disciplines, to retrieve relevant information to the given input, and then incorporate this retrieved knowledge into LLMs, which enables the models to stay accurate and current with the world knowledge.

A particularly salient application of retrieval-augmented LLMs is in handling QA tasks, whose goal is to provide correct answers in response to user queries, especially those of high complexity. Early work on retrieval-augmented LLMs focuses primarily on single-hop queries (Lazaridou et al., 2022; Ram et al., 2023), where answers are typically found within a single document; therefore, this approach involves retrieving a relevant document based on the query and subsequently integrating this information into QA models to formulate a response. However, unlike this single-hop QA, some queries require connecting and aggregating multiple documents, which are, furthermore, often not answerable through a single-step process of retrieval-and-response. An example query is ‘When did the people who captured Malakoff

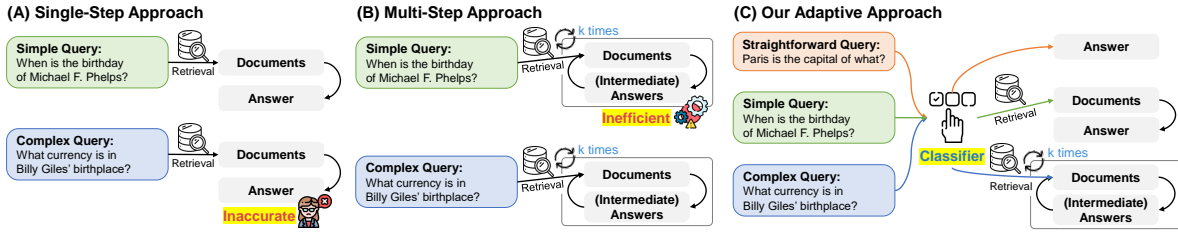


Figure 2: A conceptual comparison of different retrieval-augmented LLM approaches to question answering. (A) In response to a query, this single-step approach retrieves relevant documents and then generates an answer. However, it may not be sufficient for complex queries that require multi-step reasoning. (B) This multi-step approach iteratively retrieves documents and generates intermediate answers, which is powerful yet largely inefficient for the simple query since it requires multiple accesses to both LLMs and retrievers. (C) Our adaptive approach can select the most suitable strategy for retrieval-augmented LLMs, ranging from iterative, to single, to even no retrieval approaches, based on the complexity of given queries determined by our classifier.

075 come to the region where Philipsburg is located?',
 076 which requires four reasoning steps to solve. There-
 077 fore, to effectively handle such complex queries,
 078 recent studies have concentrated largely on multi-
 079 step and multi-reasoning QA, which requires itera-
 080 tive accesses to both LLMs and retrievers multiple
 081 times (Press et al., 2023; Trivedi et al., 2023), at
 082 the cost of heavy computational overheads.

083 Yet, we should rethink that: In a real-world sce-
 084 nario, are all the requests from users complex?
 085 Instead, users might often ask simple and straight-
 086 forward questions, while only occasionally asking
 087 complex ones. Specifically, a query such as ‘Paris
 088 is the capital of what?’ is likely to be asked more
 089 frequently, compared to the aforementioned multi-
 090 step query, and this simpler query might also be
 091 easily answerable by the LLMs themselves, with-
 092 out accessing external knowledge. In other words,
 093 a multi-step QA approach could result in unneces-
 094 sary computational overheads for simple queries,
 095 even though it is vital for complex queries (see
 096 Figure 2 (A)). On the other hand, handling com-
 097 plex queries with single-step-retrieval or even non-
 098 retrieval strategies would be largely insufficient
 099 (Figure 2 (B)). This suggests a need for an adaptive
 100 QA system, which can dynamically adjust the op-
 101 erational strategies of retrieval-augmented LLMs
 102 based on the query complexity. While some recent
 103 approaches are capable of doing this based on the
 104 frequency of entities in queries (Mallen et al., 2023)
 105 or the generated outputs from models for multi-step
 106 QA (Trivedi et al., 2023), they are suboptimal: for-
 107 mer methods are overly simplistic, failing to con-
 108 sider multi-hop queries; meanwhile, later ones are
 109 excessively complex, terminating answer solving
 110 steps after several rounds of module accesses.

111 In this work, considering the diverse complexity
 112 levels of real-world queries, we argue that previ-
 113 ous one-size-fits-all approaches might be inade-
 114 quate to cover all of these variations. Instead, we
 115 propose to select the most suitable strategy from

116 a range of (retrieval-augmented) LLMs, each of
 117 which is tailored to the specific complexity of the
 118 input query. Notably, a critical step in this process
 119 is pre-defining the query complexity, which is in-
 120 strumental in determining the most fitting model
 121 to it. In this work, we operationalize this process
 122 with a novel classifier, which is a smaller model
 123 trained to predict the complexity level of incoming
 124 queries (see Figure 2 (c)). Moreover, we automati-
 125 cally collect its training datasets without human
 126 labeling, by leveraging the predicted outcomes (i.e.,
 127 which models accurately respond to which queries)
 128 as well as by capitalizing on the inherent biases in
 129 existing datasets (i.e., samples in the datasets are de-
 130 signed for either single-step or multi-step QA sce-
 131 narios). This proposed approach can offer a robust
 132 middle ground among the iterative LLM augmen-
 133 tation methods for complex queries, the single-step
 134 methods for simpler queries, and even no-retrieval-
 135 augmented methods for the most straightforward
 136 queries (answerable by LLMs themselves), thereby
 137 significantly enhancing the overall efficiency and
 138 accuracy of QA systems, as shown in Figure 1.
 139 We refer to our framework as Adaptive Retrieval-
 140 Augmented Generation (Adaptive-RAG).

141 We validate our Adaptive-RAG using benchmark
 142 open-domain QA datasets, covering a wide range
 143 of query complexity from single-hop (Rajpurkar
 144 et al., 2016; Joshi et al., 2017; Kwiatkowski et al.,
 145 2019) to multi-hop (Yang et al., 2018; Ho et al.,
 146 2020; Trivedi et al., 2022b) queries. The exper-
 147 imental results show that our model significantly
 148 improves the overall accuracy and efficiency of QA,
 149 compared to the prior adaptive retrieval strategies,
 150 on multiple LLMs, such as GPT-3.5 (Brown et al.,
 151 2020) and FLAN-T5 series (Chung et al., 2022).

152 Our contributions and findings are threefold:

- 153 • We pointed out the realistic scenario of queries of
 154 varying complexities, and discovered that exist-
 155 ing retrieval-augmented generation approaches
 156 tend to be overly simple or complex for them.

- We propose to adapt retrieval-augmented LLMs to the query complexity assessed by the classifier, which can enable the utilization of the most suitable approach tailored to each query.
- We show that our Adaptive-RAG is highly effective and efficient, balancing between the complexity and the simplicity for diverse queries.

2 Related Work

Open-domain QA Open-domain QA is the task of accurately answering a query by sourcing for the query-relevant documents, and then interpreting them to provide answers (Chen et al., 2017; Zhu et al., 2021), which, thus, generally involves two modules: a retriever (Karpukhin et al., 2020; Xiong et al., 2021) and a reader (Yang et al., 2019; Izacard and Grave, 2021). Along with the recent emergence of LLMs with superior reasoning capabilities thanks to their billion-sized parameters (Wei et al., 2022a), the synergy between the LLMs and retrievers has led to significant advancements (Lazaridou et al., 2022; Ram et al., 2023). Specifically, this integration has been shown to enhance Open-domain QA by mitigating the hallucination problem from LLMs through strengthened reasoning abilities of the reader, as well as utilizing the retrieved, external documents (Cho et al., 2023). Despite these advancements for single-hop retrieval-augmented LLMs, however, the complexity of some queries necessitates a more complex strategy.

Multi-hop QA Multi-hop QA is an extension of conventional Open-domain QA, which additionally requires the system to comprehensively gather and contextualize information from multiple documents (often iteratively), to answer more complex queries (Trivedi et al., 2022a; Yang et al., 2018). In the realm of multi-hop QA, the approach to iteratively access both LLMs and the retrieval module is generally employed. Specifically, Khattab et al. (2022), Press et al. (2023), Pereira et al. (2023) and Khot et al. (2023) proposed to first decompose the multi-hop queries into simpler single-hop queries, repeatedly access the LLMs and retriever to solve these sub-queries, and merge their solutions to formulate a complete answer. In contrast to this decomposition-based approach, other recent studies, such as Yao et al. (2023) and Trivedi et al. (2023), explored the interleaving of Chain-of-Thought reasoning (Wei et al., 2022b) — a method where a logical sequence of thoughts is generated — with document retrieval, repeatedly applying this

process until the reasoning chain generates the answer. In addition, Jiang et al. (2023) introduced an approach for repeatedly retrieving new documents if the tokens within generated sentences have low confidence. However, the aforementioned methods overlooked the fact that, in real-world scenarios, queries are of a wide variety of complexities. Therefore, it would be largely inefficient to iteratively access LLMs and retrievers for every query, which might be simple enough to be solved with a single retrieval step or even only with an LLM itself.

Adaptive Retrieval To handle queries of varying complexities, the adaptive retrieval strategy aims to dynamically decide whether to retrieve documents or not, based on each query’s complexity. In this vein, Mallen et al. (2023) proposed to decide the query’s complexity level based on the frequency of its entities and suggested using the retrieval modules only when the frequency falls below a certain threshold. However, this approach, focusing solely on the binary decision to retrieve or not, may not be sufficient for solving more complex queries that require multiple reasoning steps. Concurrent to our work, Asai et al. (2023) suggested training a sophisticated model to dynamically retrieve, critique, and generate the text. However, we argue that all the aforementioned adaptive retrieval methods that rely on a single model might be suboptimal in handling a variety of queries of a range of different complexities since they tend to be either overly simple or complex for all the input queries, which gives rise to the new approach that can select the most suitable strategy of retrieval-augmented LLMs tailored to the query complexity.

3 Method

In this section, we describe our approach to adapting retrieval-augmented LLMs, by pre-determining the query complexity and then selecting the most fitting strategies for retrieval-augmented LLMs.

3.1 Preliminaries

We begin with preliminaries, formally introducing different strategies of retrieval-augmented LLMs.

Non Retrieval for QA Let us first define an LLM as a model LLM , which takes a sequence of tokens $\mathbf{x} = [x_1, x_2, \dots, x_n]$ as an input and then generates a sequence of tokens $\mathbf{y} = [y_1, y_2, \dots, y_n]$ as an output, which is formalized as follows: $\mathbf{y} = \text{LLM}(\mathbf{x})$. Then, in our problem setup for QA, \mathbf{x} and \mathbf{y} become the input query (q) from the user and the

generated answer (\bar{a}) from the LLM, respectively: $q = x$ and $\bar{a} = y$. Also, subsequently, the most naïve LLM-powered QA model can be represented as follows: $\bar{a} = \text{LLM}(q)$. Ideally, \bar{a} should match the actual correct answer a . This non-retrieval-based QA method is highly efficient and could be a somewhat promising approach to handling easy queries, as the size of LLMs becomes extremely large with its effect on storing a large amount of knowledge. However, this approach is largely problematic on queries that require precise or concurrent knowledge of specific people, events, or any subjects beyond the LLMs’ internal knowledge.

Single-step Approach for QA To address the aforementioned scenarios where LLM may struggle with queries that are not answerable by LLM itself, we can utilize the external knowledge d , which includes useful information for queries, retrieved from the external knowledge source \mathcal{D} that could be an encyclopedia (e.g., Wikipedia) consisting of millions of documents. Specifically, to obtain such d from \mathcal{D} , a specific retrieval model is necessary, which returns documents based on their relevance with the given query. This process can be formulated as follows: $d = \text{Retriever}(q; D)$, where Retriever is the retrieval model, with $d \in \mathcal{D}$. Here, we can use any off-the-shelf retriever (Robertson et al., 1994; Karpukhin et al., 2020).

After the retrieval step is done, we now have a pair of query q and its relevant documents d . Then, in order to augment LLMs with this retrieved external knowledge, we can incorporate it into the input of LLMs, represented as follows: $\bar{a} = \text{LLM}(q, d)$. This process allows LLMs to gain access to external information contained in d , which can provide the supplementary context that the internal knowledge of LLM lacks, which can subsequently improve the accurateness and concurrency of LLMs for QA.

Multi-step Approach for QA Even though the aforementioned single-step approach offers significant improvements over non-retrieval for q that requires external knowledge, it encounters notable limitations, particularly when dealing with complex queries that necessitate synthesizing information from multiple source documents and reasoning over them. This is where a multi-step approach and reasoning for QA become essential.

In this multi-step approach, LLM interacts with Retriever in several rounds, progressively refining its understanding of q , until it formulates the fi-

nal answer from findings accumulated across those multiple steps. Specifically, the process begins with the initial query q , and at every retrieval step i , new documents d_i are retrieved from \mathcal{D} and then incorporated into the input of LLMs, as follows: $\bar{a}_i = \text{LLM}(q, d_i, c_i)$, where the additional context c_i can be composed of previous documents and outcomes $(d_1, d_2, \dots, d_{i-1}, \bar{a}_1, \bar{a}_2, \dots, \bar{a}_{i-1})$, and $d_i = \text{Retriever}(q, c_i; D)$ ¹. We would like to note that this iterative, multi-step process enables LLM to construct a more comprehensive and extensive foundation to solve queries effectively, specifically adept at complex multi-hop queries where answers depend on interconnected pieces of information. However, it is important to recognize that this multi-step approach can be resource-intensive due to the repeated accesses to Retriever and LLM, which entail substantial computational costs.

3.2 Adaptive-RAG: Adaptive Retrieval-Augmented Generation

We now introduce our adaptive retrieval-augmented LLMs, which are built upon three different strategies described in the previous section, and which are designed to select the most suitable strategy according to the complexity of queries.

Adapting Retrieval-Augmented LLMs Note that in real-world scenarios, not all q from users have the same level of complexity, necessitating tailored strategies for handling each query. In other words, employing the most basic, non-retrieval-based approach $\text{LLM}(q)$ to respond to the complex query q would be also ineffective (Figure 2, A); conversely, using a more elaborate, multi-step approach $\text{LLM}(q, d, c)$ for simple q would be inefficient (Figure 2, B). Therefore, our adaptive framework is designed to dynamically adjust the query-handling strategy of retrieval-augmented LLMs, which is achieved by pre-determining the complexity of each query prior to attempting a solution. Notably, this framework can offer a robust middle ground with a range of solutions, from the simplest approach for the most straightforward queries, to the one-step approach for moderate queries, and up to the most comprehensive and rigorous approach for complex queries. In addition, since the operations of LLM and Retriever remain consistent re-

¹It is worth noting that implementations of the LLM and retriever vary across different multi-step retrieval-augmented LLM approaches (Trivedi et al., 2023; Press et al., 2023; Yao et al., 2023); therefore, the context c_i may incorporate none, some, or all of previous documents and answers.

353 regardless of inputs to them, our method can seeming- 401
354 lessly go back and forth across queries of different 402
355 complexities, without changing the internal model 403
356 architectures or parameters during adaption. 404

357 **Query Complexity Assessment** To operational- 405
358 ize our adaptive retrieval-augmented LLM frame- 406
359 work, we should determine the query complexity, 407
360 and to achieve this, we propose to model a com- 408
361 plexity classifier, whose goal is to return the appro- 409
362 priate complexity level of the given query. Specif- 410
363 ically, given the query q , our classifier can be for- 411
364 mulated as follows: $o = \text{Classifier}(q)$, where 412
365 Classifier is a smaller Language Model that is 413
366 trained to classify one of three different complexity 414
367 levels and o is its corresponding class label. In our 415
368 classifier design, there are three class labels: ‘A’, 416
369 ‘B’, and ‘C’, where ‘A’ indicates that q is straight-
370 forward and answerable by LLM(q) itself, ‘B’ in-
371 dicates that q has the moderate complexity where
372 at least a single-step approach LLM(q, d) is needed,
373 and ‘C’ indicates that q is complex, requiring the
374 most extensive solution LLM(q, d, c)².

375 **Training Strategy** The remaining step is to train 422
376 the smaller Language Model for Classifier, to 423
377 accurately predict its complexity o in response to 424
378 the given query q . Yet, there is no annotated dataset 425
379 available for query-complexity pairs. Hence, we 426
380 propose to automatically construct the training 427
381 dataset with two particular strategies. 428

382 To be specific, we first aim at labeling the query 429
383 complexity based on the results from three different 430
384 retrieval-augmented LLM strategies, in order to 431
385 determine the label by its needs. For example, if 432
386 the simplest non-retrieval-based approach correctly 433
387 generates the answer, the label for its corresponding 434
388 query is assigned ‘A’. Also, to break the tie between 435
389 different models in providing the label to the query, 436
390 we provide a higher priority to a simpler model. 437
391 In other words, if both single-step and multi-step 438
392 approaches produce the same correct answer while 439
393 the non-retrieval-based approach fails, we assign 440
394 label ‘B’ to its corresponding query. 441

395 However, this labeling strategy has a limita- 442
396 tion in that not all the queries are assigned the la- 443
397 bels, since all three retrieval-augmented approaches 444
398 may fail to generate the correct answer. On the 445
399 other hand, the benchmark datasets may already 446
400 have meaningful inductive biases about the most

²We consider three levels of query complexity, and leave the exploration of more fine-grained complexities as future work.

appropriate retrieval-augmented LLM strategies 401
for their queries, considering the ways they are 402
created (e.g., QA datasets that require sequential 403
reasoning usually necessitate the multi-step ap- 404
proach; while queries of those with the labeled 405
single documents can be ideally answerable with 406
the single-step approach). Therefore, for those 407
queries that remain unlabeled after the first labeling 408
step, we assign ‘B’ to queries in single-hop datasets 409
and ‘C’ to queries in multi-hop datasets. Finally, 410
we train Classifier with these automatically- 411
collected query-complexity pairs³, by using a cross- 412
entropy loss. Then, at inference, we can deter- 413
mine the complexity of the query, which is one of 414
{‘A’, ‘B’, ‘C’}, by forwarding it to Classifier: 415
 $o = \text{Classifier}(q)$. 416

4 Experimental Setups 417

In this section, we explain datasets, models, met- 418
rics, and implementation details. We further pro- 419
vide the additional details in Appendix A. 420

4.1 Datasets 421

In order to simulate a realistic scenario, where dif- 422
ferent queries have varying complexities, in the 423
unified experimental setting, we use both the single- 424
hop and multi-hop QA datasets simultaneously. 425

Single-hop QA For simpler queries, we use three 426
benchmark single-hop QA datasets, which consist 427
of queries and their associated documents contain- 428
ing answers, namely **1) SQuAD v1.1** (Rajpurkar 429
et al., 2016), **2) Natural Questions** (Kwiatkowski 430
et al., 2019), and **3) TriviaQA** (Joshi et al., 2017). 431

Multi-hop QA To consider more complex query 432
scenarios, we use three benchmark multi-hop QA 433
datasets, which require sequential reasoning over 434
multiple documents, namely **1) MuSiQue** (Trivedi 435
et al., 2022a), **2) HotpotQA** (Yang et al., 2018), 436
and **3) 2WikiMultiHopQA** (Ho et al., 2020). 437

4.2 Models 438

We compare our Adaptive-RAG against relevant 439
models, including three retrieval-augmented LLM 440
strategies (in Section 3.1) and the adaptive re- 441
trieval approaches (Mallen et al., 2023; Asai et al., 442
2023), which can be grouped into one of three cat- 443
egories: Simple, Adaptive, and Complex. Specif- 444
ically, Simple approaches include the **1) No Re-** 445
trieval and **2) Single-step Approach**-based meth- 446

³As we automatically assign classifier labels, there might be errors in labeling and might be more advanced strategies to automatically assign labels, which we leave as future work.

Table 1: Averaged results on a collection of benchmark datasets for open-domain question answering including the single-hop and multi-hop queries, with different LLMs. Self-RAG* is trained with the different base LLM, namely LLaMA2 (Touvron et al., 2023); therefore, we compare the results of FLAN-T5-XL (3B) with the results from Self-RAG with LLaMA2 (7B) and the results of others with the results from Self-RAG with LLaMA2 (13B). We emphasize our results in bold, for easy comparisons.

Types	Methods	FLAN-T5-XL (3B)				FLAN-T5-XXL (11B)				GPT-3.5 (Turbo)			
		EM	F1	Step	Time	EM	F1	Step	Time	EM	F1	Step	Time
Simple	No Retrieval	14.87	21.12	0.00	0.11	17.83	25.14	0.00	0.08	35.77	48.56	0.00	0.71
	Single-step Approach	34.83	44.31	1.00	1.00	37.87	47.63	1.00	1.00	34.73	46.99	1.00	1.00
Adaptive	Adaptive Retrieval	23.87	32.24	0.50	0.56	26.93	35.67	0.50	0.54	35.90	48.20	0.50	0.86
	Self-RAG*	9.90	20.79	0.72	0.43	10.87	22.98	0.74	0.23	10.87	22.98	0.74	1.50
	Adaptive-RAG (Ours)	37.17	46.94	2.17	3.60	38.90	48.62	1.35	2.00	37.97	50.91	1.03	1.46
Complex	Multi-step Approach	39.00	48.85	4.69	8.81	40.13	50.09	2.13	3.80	38.13	50.87	2.81	3.33

ods. Adaptive approaches include the **3) Adaptive Retrieval** (Mallen et al., 2023), **4) Self-RAG** (Asai et al., 2023), and our **5) Adaptive-RAG**, which can adaptively perform retrieval based on the question complexity. For the **6) Multi-step Approach**, we use the most sophisticated state-of-the-art method (Trivedi et al., 2023) iteratively accessing both the retriever and LLM with Chain-of-Thought reasoning (Wei et al., 2022b), for every query. Note that models across different categories are not directly comparable. Yet, in the ideal setting, Adaptive approaches should be more effective than those in the Simple category while simultaneously being more efficient than the Complex one.

4.3 Evaluation Metrics

When it comes to evaluating adaptive models, it is essential to simultaneously consider both the task performance and efficiency along with their trade-offs. Thus, we report the results with four metrics, where two of them measure the effectiveness and the other two measure the efficiency. In particular, for effectiveness, we use F1 and EM following the standard evaluation protocol (Mallen et al., 2023; Baek et al., 2023), where F1 measures the number of overlapping words between the predicted answer and the ground truth and EM measures whether they are the same. For efficiency, we measure the number of retrieval-and-generate steps and the average time for answering each query relative to the one-step approach.

4.4 Implementation Details

For a fair comparison and following Mallen et al. (2023) and Trivedi et al. (2023), we use the same retriever, a term-based sparse retrieval model known as BM25 (Robertson et al., 1994), across all different models. For the external document corpus, we use different sources depending on the dataset type: the Wikipedia corpus preprocessed by Karpukhin et al. (2020) for single-hop datasets, and the pre-

processed corpus by Trivedi et al. (2023) for multi-hop datasets. Regarding the LLMs that are used to generate answers, we use the FLAN-T5 series models (Chung et al., 2022) of XL with 3B parameters and XXL with 11B parameters, and the GPT-3.5 model (gpt-3.5-turbo-instruct). For the retrieval-augmented LLM design, we follow the implementation details from Trivedi et al. (2023), which include input prompts, instructions, and the number of test samples for evaluation (e.g., 500 samples per dataset). In our Adaptive-RAG, for the query-complexity classifier, we use and train the T5-Large model (Raffel et al., 2020). Specifically, the classifier is trained using the epoch that shows the best performance until 100 training iterations from the validation set, with the learning rate of $3e-5$ and the AdamW (Loshchilov and Hutter, 2019) as an optimizer. Regarding its training data, we sample and annotate 400 queries from 6 datasets based on its inductive bias (single-hop for one-step approach and multi-hop for multi-step). In addition, we use predicted outcomes of three different strategies over 400 queries sampled from each dataset. Note that those queries used for classifier training do not overlap with the testing queries for QA.

5 Experimental Results and Analyses

In this section, we show the overall experimental results and offer in-depth analyses of our method.

Main Results First of all, Table 1 shows our main results averaged over all considered datasets, which corroborate our hypothesis that simple retrieval-augmented strategies are less effective than the complex strategy, while the complex one is significantly more expensive than the simple ones. In addition, we report the more granular results with FLAN-T5-XL on each of single-hop and multi-hop datasets in Table 2 (and more with different LLMs in Table 6 and Table 7 of Appendix), which are consistent with our hypothesis observed in Table 1.

Table 2: Results on each of a collection of datasets with FLAN-T5-XL (3B) as the LLM. We emphasize our results in bold.

Data	Types	Methods	SQuAD				Natural Questions				TriviaQA			
			EM	F1	Step	Time	EM	F1	Step	Time	EM	F1	Step	Time
Single-step	Simple	No Retrieval	3.60	10.50	0.00	0.11	14.20	19.00	0.00	0.13	25.00	31.80	0.00	0.13
		Single-step Approach	27.80	39.30	1.00	1.00	37.80	47.30	1.00	1.00	53.60	62.40	1.00	1.00
	Adaptive	Adaptive Retrieval	13.40	23.10	0.50	0.55	28.20	36.00	0.50	0.56	38.40	46.90	0.50	0.56
		Self-RAG*	2.20	11.20	0.63	0.50	31.40	39.00	0.63	0.17	12.80	29.30	0.68	0.45
		Adaptive-RAG (Ours)	26.80	38.30	1.37	2.02	37.80	47.30	1.00	1.00	52.20	60.70	1.23	1.54
Complex	Multi-step Approach	24.40	35.60	4.52	9.03	38.60	47.80	5.04	10.18	53.80	62.40	5.28	9.22	

Data	Types	Methods	MuSiQue				HotpotQA				2WikiMultiHopQA			
			EM	F1	Step	Time	EM	F1	Step	Time	EM	F1	Step	Time
Multi-step	Simple	No Retrieval	2.40	10.70	0.00	0.11	16.60	22.71	0.00	0.11	27.40	32.04	0.00	0.10
		Single-step Approach	13.80	22.80	1.00	1.00	34.40	46.15	1.00	1.00	41.60	47.90	1.00	1.00
	Adaptive	Adaptive Retrieval	6.40	15.80	0.50	0.55	23.60	32.22	0.50	0.55	33.20	39.44	0.50	0.55
		Self-RAG*	1.60	8.10	0.73	0.51	6.80	17.53	0.73	0.45	4.60	19.59	0.93	0.49
		Adaptive-RAG (Ours)	23.60	31.80	3.22	6.61	42.00	53.82	3.55	5.99	40.60	49.75	2.63	4.68
Complex	Multi-step Approach	23.00	31.90	3.60	7.58	44.60	56.54	5.53	9.38	49.60	58.85	4.17	7.37	

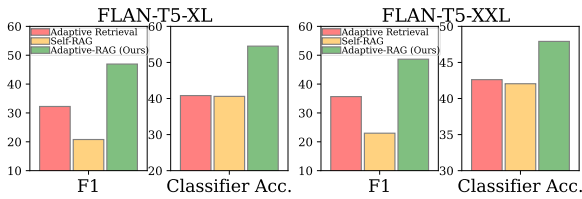


Figure 3: Performance on QA and query-complexity assessment of different adaptive approaches for retrieval-augmented LLMs with FLAN-T5 XL and XXL. For labeling the complexity of queries, we use the silver data annotated from the prediction outcomes of models (described in Section 3.2).

However, in a real-world scenario, not all users ask queries with the same level of complexity, which emphasizes the importance of the need for adaptive strategies. Note that among the adaptive strategies, our Adaptive-RAG shows remarkable effectiveness over the competitors (Table 1). This indicates that merely focusing on the decision of whether to retrieve or not is suboptimal. Also, as shown in Table 2, such simple adaptive strategies are particularly inadequate for handling complex queries in multi-hop datasets, which require aggregated information and reasoning over multiple documents. Meanwhile, our approach can consider a more fine-grained query handling strategy by further incorporating an iterative module for complex queries. Furthermore, in a realistic setting, we should take into account not only effectiveness but also efficiency. As shown in Table 1, compared to the complex multi-step strategy, our proposed adaptive strategy is significantly more efficient across all model sizes. This is meaningful in this era of LLMs, where the cost of accessing them is a critical factor for practical applications and scalability.

Classifier Performance To understand how the proposed classifier works, we analyze its performance across different complexity options. As

Table 3: Results on QA and complexity classification with varying the data annotation strategies for training the classifier.

Training Strategies	QA		Classifier (Accuracy)			
	F1	Step	All	No	One	Multi
Adaptive-RAG (Ours)	46.94	1084	54.52	30.52	66.28	65.45
w/o Binary	43.43	640	60.30	62.19	65.70	39.55
w/o Silver	48.79	1464	40.00	0.00	53.98	75.91

Figure 3 shows, the classification accuracy of our Adaptive-RAG is better than that of the other adaptive retrieval baselines, which leads to overall QA performance improvements. This indicates that ours more accurately classifies the options with various granularities, including not retrieving documents, retrieving with a single access, and retrieving with multiple accesses.

Analyses on Training Data for Classifier We have shown that the classifier plays an important role in adaptive retrieval. Here, we further analyze the different strategies for training the classifier by ablating our full training strategy, which includes two approaches: generating silver data from predicted outcomes of models and utilizing inductive bias in datasets (see Section 3.2). As Table 3 shows, compared to the training strategy relying solely on the data derived from inductive bias, ours is significantly more efficient. This efficiency is partly because ours also takes into account the case that does not consider any documents at all, as also implied by the classification accuracy; meanwhile, queries in the existing datasets do not capture the information on whether the retrieval is not required. On the other hand, in the case of only using the silver data annotated from the correct predictions, while its overall classification accuracy is high, the overall QA performance implies that relying on the silver data may not be optimal. This may be

Table 4: Case study with GPT-3.5 (Turbo), where we represent the factual error in red and the accurate information in blue.

Dataset	Question	Adaptive Retrieval	Adaptive-RAG (Ours)
NQ (Single-hop)	Which famous corporate logo changed to a flat colour/color sans serif font in its first major change since 1999?	Query Type: B (Single-step Approach) Reasoning: The famous corporate logo that changed to a flat color sans-serif font in its first major change since 1999 is Microsoft . Microsoft changed its logo in 2012, replacing the italicized, multi-colored font with a simple, flat, grey font. So the answer is: Microsoft .	Query Type: A (Non Retrieval) Reasoning: The famous corporate logo that changed to a flat color sans serif font in its first major change since 1999 is Google . Google changed its logo in 2015. So the answer is: Google .
MuSiQue (Multi-hop)	Who is the child of the Italian navigator who explored the eastern coast of the continent César Gaytan was born in for the English?	Query Type: A (Non Retrieval) Reasoning: César Gaytan was born in the continent of North America . The Italian navigator who explored the eastern coast of North America for the English was Giovanni Caboto , also known as John Cabot. So the answer is: Giovanni Caboto/John Cabot .	Query Type: C (Multi-step Approach) Reasoning: César Gaytan was born in Guadalajara, Jalisco, Mexico . The Italian navigator who explored the eastern coast of the continent for the English is John Cabot . John Cabot’s son is Sebastian Cabot . So the answer is: Sebastian Cabot .

Table 5: Results with varying model sizes for classifiers.

Sizes	QA		Classifier (Accuracy)			
	F1	Step	All	No	One	Multi
Small (60M)	45.83	964	53.48	26.65	70.62	53.18
Base (223M)	45.97	983	53.41	26.42	69.46	56.82
Large (770M)	46.94	1084	54.52	30.52	66.28	65.45

because this silver data does not cover complexity labels on incorrectly predicted queries, which leads to lower generalization effect on queries relevant to them. Meanwhile, by also incorporating complexity labels from dataset bias (single-hop vs multi-hop), the classifier becomes more accurate in predicting multi-hop queries, leading to the better performance. It is worth noting that our automatic labeling strategies are two particular instantiations for training the classifier, and that there could be other instantiations, which we leave as future work.

Analyses on Classifier Size To investigate the sensitivity of our classifier according to its varying sizes, we conducted further experiments. As shown in Table 5, we observe no significant performance differences among classifiers of various sizes, even with reduced complexity and fewer parameters in smaller classifiers. This indicates that our proposed classifier can contribute to resource-efficient settings in real-use cases with smaller sizes without compensating performance.

Case Study We conduct a case study to qualitatively compare our Adaptive-RAG against Adaptive Retrieval. Table 4 shows the classified complexity and the query handling patterns for both simple and complex questions. First, for the simple single-hop question, our Adaptive-RAG identifies that it is answerable by only the LLM’s parametric knowledge about ‘Google’. In contrast, Adaptive Retrieval fetches additional documents, leading to longer processing times and occasionally pro-

ducing incorrect responses due to the inclusion of partially irrelevant information about ‘Microsoft’. Meanwhile, when faced with a complex question, our Adaptive-RAG seeks out relevant information, including details like ‘a son of John Cabot’, which may not have been stored in LLMs, while Adaptive Retrieval fails to request such information from external sources, resulting in inaccurate answers.

6 Conclusion

In this work, we proposed the Adaptive Retrieval-Augmented Generation framework, referred to as Adaptive-RAG, to handle queries of various complexities. Specifically, Adaptive-RAG is designed to dynamically adjust its query handling strategies in the unified retrieval-augmented LLM based on the complexity of queries that they encounter, which spans across a spectrum of the non-retrieval-based approach for the most straightforward queries, to the single-step approach for the queries of moderate complexity, and finally to the multi-step approach for the complex queries. The core step of our Adaptive-RAG lies in determining the complexity of the given query, which is instrumental in selecting the most suitable strategy for its answer. To operationalize this process, we trained a smaller Language Model with query-complexity pairs, which are automatically annotated from the predicted outcomes and the inductive biases in datasets. We validated our Adaptive-RAG on a collection of open-domain QA datasets, covering the multiple query complexities including both the single- and multi-hop questions. The results demonstrate that our Adaptive-RAG enhances the overall accuracy and efficiency of QA systems, allocating more resources to handle complex queries while efficiently handling simpler queries, compared to the existing one-size-fits-all approaches that tend to be either minimalist or maximalist over varying query complexities.

650 Limitations

651 While our Adaptive-RAG shows clear advantages
652 in effectiveness and efficiency by determining the
653 query complexity and then leveraging the most
654 suitable approach for tackling it, it is important to
655 recognize that there still exists potential avenues
656 for improving the classifier from the perspectives
657 of its training datasets and architecture. Specifi-
658 cally, as there are no available datasets for training
659 the query-complexity classifier, we automatically
660 create new data based on the model prediction out-
661 comes and the inductive dataset biases. However,
662 our labeling process is one specific instantiation
663 of labeling the query complexity, and it may have
664 the potential to label queries incorrectly despite its
665 effectiveness. Therefore, future work would create
666 new datasets that are annotated with a diverse range
667 of query complexities, in addition to the labels of
668 question-answer pairs. Also, as shown in Figure 3,
669 there is still room to improve the effectiveness of
670 the classifier. In other words, our classifier design
671 based on the smaller LM is the initial, simplest
672 instantiation for classifying the query complexity,
673 and based upon it, future work may improve the
674 classifier architecture and its performance, which
675 will positively contribute to the QA performance.

676 Ethics Statement

677 The experimental results on Adaptive-RAG vali-
678 date its applicability in realistic scenarios, where a
679 wide range of diverse user queries exists. Nonethe-
680 less, given the potential diversity of real-world user
681 inputs, it is crucial to also consider scenarios where
682 these inputs might be offensive or harmful. We
683 should be aware that such inputs could lead to the
684 retrieval of offensive documents and the genera-
685 tion of inappropriate responses by the retrieval-
686 augmented LLMs. To address this challenge, de-
687 veloping methods to detect and manage offensive
688 or inappropriate content in both user inputs and re-
689 trieved documents within the retrieval-augmented
690 framework is essential. We believe that this is a
691 critical area for future work.

692 References

693 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John-
694 son, Dmitry Lepikhin, Alexandre Passos, Siamak
695 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng
696 Chen, Eric Chu, Jonathan H. Clark, Laurent El
697 Shafey, Yanping Huang, Kathy Meier-Hellstern, Gau-
698 rav Mishra, Erica Moreira, Mark Omernick, Kevin

Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, 699
Yuanzhong Xu, Yujing Zhang, Gustavo Hernández 700
Ábrego, Junwhan Ahn, Jacob Austin, Paul Barham, 701
Jan A. Botha, James Bradbury, Siddhartha Brahma, 702
Kevin Brooks, Michele Catasta, Yong Cheng, Colin 703
Cherry, Christopher A. Choquette-Choo, Aakanksha 704
Chowdhery, Clément Crepy, Shachi Dave, Mostafa 705
Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, 706
Nan Du, Ethan Dyer, Vladimir Feinberg, Fangxi- 707
aoyu Feng, Vlad Fienber, Markus Freitag, Xavier 708
Garcia, Sebastian Gehrmann, Lucas Gonzalez, and 709
et al. 2023. [Palm 2 technical report](#). *arXiv preprint* 710
arXiv:2305.10403. 711

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, 712
and Hannaneh Hajishirzi. 2023. [Self-RAG: Self-](#) 713
[reflective retrieval augmented generation](#). In 714
NeurIPS 2023 Workshop on Instruction Tuning and 715
Instruction Following. 716

Jinheon Baek, Soyeong Jeong, Minki Kang, Jong Park, 717
and Sung Ju Hwang. 2023. [Knowledge-augmented](#) 718
[language model verification](#). In *Proceedings of the* 719
2023 Conference on Empirical Methods in Natural 720
Language Processing, EMNLP 2023, Singapore, De- 721
cember 6-10, 2023, pages 1720–1736. Association 722
for Computational Linguistics. 723

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, 724
Trevor Cai, Eliza Rutherford, Katie Millican, George 725
van den Driessche, Jean-Baptiste Lespiau, Bogdan 726
Damoc, Aidan Clark, Diego de Las Casas, Aurelia 727
Guy, Jacob Menick, Roman Ring, Tom Hennigan, 728
Saffron Huang, Loren Maggiore, Chris Jones, Albin 729
Cassirer, Andy Brock, Michela Paganini, Geoffrey 730
Irving, Oriol Vinyals, Simon Osindero, Karen Si- 731
monyán, Jack W. Rae, Erich Elsen, and Laurent Sifre. 732
2022. [Improving language models by retrieving from](#) 733
[trillions of tokens](#). In *International Conference on* 734
Machine Learning, ICML 2022, 17-23 July 2022, Bal- 735
timore, Maryland, USA, volume 162 of *Proceedings* 736
of Machine Learning Research, pages 2206–2240. 737
PMLR. 738

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie 739
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind 740
Neelakantan, Pranav Shyam, Girish Sastry, Amanda 741
Askell, Sandhini Agarwal, Ariel Herbert-Voss, 742
Gretchen Krueger, Tom Henighan, Rewon Child, 743
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, 744
Clemens Winter, Christopher Hesse, Mark Chen, Eric 745
Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, 746
Jack Clark, Christopher Berner, Sam McCandlish, 747
Alec Radford, Ilya Sutskever, and Dario Amodei. 748
2020. [Language models are few-shot learners](#). In *Ad-* 749
vances in Neural Information Processing Systems 33: 750
Annual Conference on Neural Information Process- 751
ing Systems 2020, NeurIPS 2020, December 6-12, 752
2020, virtual. 753

Danqi Chen, Adam Fisch, Jason Weston, and Antoine 754
Bordes. 2017. [Reading wikipedia to answer open-](#) 755
[domain questions](#). In *Proceedings of the 55th Annual* 756
Meeting of the Association for Computational Lin- 757
guistics, ACL 2017, Vancouver, Canada, July 30 - 758

759	<i>August 4, Volume 1: Long Papers</i> , pages 1870–1879.		
760	Association for Computational Linguistics.		
761	Sukmin Cho, Jeongyeon Seo, Soyeong Jeong, and		
762	Jong C. Park. 2023. Improving zero-shot reader by		
763	reducing distractions from irrelevant documents in		
764	open-domain question answering . In <i>Findings of the</i>		
765	<i>Association for Computational Linguistics: EMNLP</i>		
766	<i>2023, Singapore, December 6-10, 2023</i> , pages 3145–		
767	3157. Association for Computational Linguistics.		
768	Hyung Won Chung, Le Hou, Shayne Longpre, Barret		
769	Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang,		
770	Mostafa Dehghani, Siddhartha Brahma, Albert Web-		
771	son, Shixiang Shane Gu, Zhuyun Dai, Mirac Suz-		
772	gun, Xinyun Chen, Aakanksha Chowdhery, Sharan		
773	Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao,		
774	Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav		
775	Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam		
776	Roberts, Denny Zhou, Quoc V. Le, and Jason Wei.		
777	2022. Scaling instruction-finetuned language models.		
778	<i>arXiv preprint arXiv:2210.11416</i> .		
779	Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara,		
780	and Akiko Aizawa. 2020. Constructing A multi-hop		
781	QA dataset for comprehensive evaluation of reason-		
782	ing steps . In <i>Proceedings of the 28th International</i>		
783	<i>Conference on Computational Linguistics, COLING</i>		
784	<i>2020, Barcelona, Spain (Online), December 8-13,</i>		
785	<i>2020</i> , pages 6609–6625. International Committee on		
786	Computational Linguistics.		
787	Gautier Izacard and Edouard Grave. 2021. Leveraging		
788	passage retrieval with generative models for open do-		
789	main question answering . In <i>Proceedings of the 16th</i>		
790	<i>Conference of the European Chapter of the Associ-</i>		
791	<i>ation for Computational Linguistics: Main Volume,</i>		
792	<i>EACL 2021, Online, April 19 - 23, 2021</i> , pages 874–		
793	880. Association for Computational Linguistics.		
794	Gautier Izacard, Patrick S. H. Lewis, Maria Lomeli,		
795	Lucas Hosseini, Fabio Petroni, Timo Schick, Jane		
796	Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and		
797	Edouard Grave. 2023. Atlas: Few-shot learning		
798	with retrieval augmented language models . <i>J. Mach.</i>		
799	<i>Learn. Res.</i> , 24:251:1–251:43.		
800	Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun,		
801	Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie		
802	Callan, and Graham Neubig. 2023. Active retrieval		
803	augmented generation . In <i>EMNLP 2023</i> .		
804	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke		
805	Zettlemoyer. 2017. Triviaqa: A large scale distantly		
806	supervised challenge dataset for reading comprehen-		
807	sion . In <i>Proceedings of the 55th Annual Meeting of</i>		
808	<i>the Association for Computational Linguistics, ACL</i>		
809	<i>2017, Vancouver, Canada, July 30 - August 4, Volume</i>		
810	<i>1: Long Papers</i> , pages 1601–1611. Association for		
811	Computational Linguistics.		
812	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick		
813	S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen,		
814	and Wen-tau Yih. 2020. Dense passage retrieval for		
815	open-domain question answering . In <i>Proceedings of</i>		
	<i>the 2020 Conference on Empirical Methods in Natu-</i>		
	<i>ral Language Processing, EMNLP 2020, November</i>		
	<i>16-20, 2020</i> . Association for Computational Linguis-		
	tics.		816
			817
			818
			819
	Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ro-		
	nan Le Bras, Akari Asai, Xinyan Yu, Dragomir R.		
	Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui.		
	2022. Realtime QA: what’s the answer right now?		
	<i>arXiv preprint arXiv:2207.13332</i> .		820
			821
			822
			823
			824
	Omar Khattab, Keshav Santhanam, Xiang Lisa		
	Li, David Hall, Percy Liang, Christopher Potts,		
	and Matei Zaharia. 2022. Demonstrate-search-		
	predict: Composing retrieval and language mod-		
	els for knowledge-intensive NLP . <i>arXiv preprint</i>		
	<i>arXiv.2212.14024</i> , abs/2212.14024.		825
			826
			827
			828
			829
			830
	Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao		
	Fu, Kyle Richardson, Peter Clark, and Ashish Sab-		
	harwal. 2023. Decomposed prompting: A modular		
	approach for solving complex tasks . In <i>The Eleventh</i>		
	<i>International Conference on Learning Representa-</i>		
	<i>tions, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> .		
	OpenReview.net.		831
			832
			833
			834
			835
			836
			837
	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-		
	field, Michael Collins, Ankur Parikh, Chris Alberti,		
	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-		
	ton Lee, Kristina Toutanova, Llion Jones, Matthew		
	Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob		
	Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natu-		
	ral questions: A benchmark for question answering		
	research . <i>Transactions of the Association for Compu-</i>		
	<i>tational Linguistics</i> , 7:452–466.		838
			839
			840
			841
			842
			843
			844
			845
			846
	Angeliki Lazaridou, Elena Gribovskaya, Wojciech		
	Stokowiec, and Nikolai Grigorev. 2022. Internet-		
	augmented language models through few-shot		
	prompting for open-domain question answering .		
	<i>arXiv preprint arXiv:2203.05115</i> .		847
			848
			849
			850
			851
	Belinda Z. Li, Sewon Min, Srinivasan Iyer, Yashar		
	Mehdad, and Wen-tau Yih. 2020. Efficient one-pass		
	end-to-end entity linking for questions . In <i>Proceed-</i>		
	<i>ings of the 2020 Conference on Empirical Methods in</i>		
	<i>Natural Language Processing, EMNLP 2020, Online,</i>		
	<i>November 16-20, 2020</i> , pages 6433–6441. Associa-		
	tion for Computational Linguistics.		852
			853
			854
			855
			856
			857
			858
	Ilya Loshchilov and Frank Hutter. 2019. Decoupled		
	weight decay regularization . In <i>7th International</i>		
	<i>Conference on Learning Representations, ICLR 2019,</i>		
	<i>New Orleans, LA, USA, May 6-9, 2019</i> . OpenRe-		
	view.net.		859
			860
			861
			862
			863
	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,		
	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.		
	When not to trust language models: Investigating		
	effectiveness of parametric and non-parametric mem-		
	ories . In <i>Proceedings of the 61st Annual Meeting of</i>		
	<i>the Association for Computational Linguistics (Vol-</i>		
	<i>ume 1: Long Papers), ACL 2023, Toronto, Canada,</i>		
	<i>July 9-14, 2023</i> , pages 9802–9822. Association for		
	Computational Linguistics.		864
			865
			866
			867
			868
			869
			870
			871
			872

873	OpenAI. 2023. GPT-4 technical report . <i>arXiv preprint arXiv:2303.08774</i> .	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esioibu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>arXiv preprint arXiv:2307.09288</i> .	930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952
874			
875	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library . In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019</i> , pages 8024–8035.		
876			
877			
878			
879			
880			
881			
882			
883			
884			
885			
886	Jayr Alencar Pereira, Robson do Nascimento Fidalgo, Roberto de Alencar Lotufo, and Rodrigo Frassetto Nogueira. 2023. Visconde: Multi-document QA with GPT-3 and neural reranking . In <i>Advances in Information Retrieval - 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2-6, 2023, Proceedings, Part II</i> , volume 13981 of <i>Lecture Notes in Computer Science</i> , pages 534–543. Springer.		
887			
888			
889			
890			
891			
892			
893			
894			
895	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> .	Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022a. 9835 musique: Multi-hop questions via single-hop question composition . <i>Trans. Assoc. Comput. Linguistics</i> , 10:539–554.	953 954 955 956
896			
897			
898			
899			
900	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>J. Mach. Learn. Res.</i> , 21:140:1–140:67.	Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022b. MuSiQue: Multi-hop questions via single-hop question composition . <i>Transactions of the Association for Computational Linguistics</i> , 10:539–554.	957 958 959 960 961
901			
902			
903			
904			
905	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100, 000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016</i> , pages 2383–2392. The Association for Computational Linguistics.	Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 10014–10037. Association for Computational Linguistics.	962 963 964 965 966 967 968 969 970
906			
907			
908			
909			
910			
911			
912	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models . <i>Transactions 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, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022a. Emergent abilities of large language models . <i>Trans. Mach. Learn. Res.</i> , 2022.	971 972 973 974 975 976 977
913			
914			
915			
916			
917	Stephen E. Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, and Mike Gatford. 1994. Okapi at TREC-3 . In <i>Proceedings of The Third Text REtrieval Conference, TREC 1994, Gaithersburg, Maryland, USA, November 2-4, 1994</i> , volume 500-225 of <i>NIST Special Publication</i> , pages 109–126. National Institute of Standards and Technology (NIST).	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022b. Chain-of-thought prompting elicits reasoning in large language models . In <i>NeurIPS</i> .	978 979 980 981 982
918			
919			
920			
921			
922			
923			
924			
925	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. REPLUG: retrieval-augmented black-box language models . <i>arXiv preprint arXiv:2301.12652</i> .	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin	983 984 985 986 987 988
926			
927			
928			
929			

989 Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos*, pages 38–45. Association for Computational Linguistics.

995 Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang,
996 Jialin Liu, Paul N. Bennett, Junaid Ahmed, and
997 Arnold Overwijk. 2021. [Approximate nearest neighbor negative contrastive learning for dense text retrieval](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

1002 Wei Yang, Yuqing Xie, Aileen Lin, Xingyu Li, Luchen
1003 Tan, Kun Xiong, Ming Li, and Jimmy Lin. 2019.
1004 [End-to-end open-domain question answering with bertserini](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Demonstrations*, pages 72–77. Association for Computational Linguistics.

1011 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,
1012 William Cohen, Ruslan Salakhutdinov, and Christo-
1013 pher D. Manning. 2018. [HotpotQA: A dataset for diverse, explainable multi-hop question answering](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.

1019 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
1020 Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023.
1021 [React: Synergizing reasoning and acting in language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

1025 Fengbin Zhu, Wenqiang Lei, Chao Wang, Jianming
1026 Zheng, Soujanya Poria, and Tat-Seng Chua. 2021.
1027 [Retrieving and reading: A comprehensive survey on open-domain question answering](#). *arXiv preprint arXiv:2101.00774*.

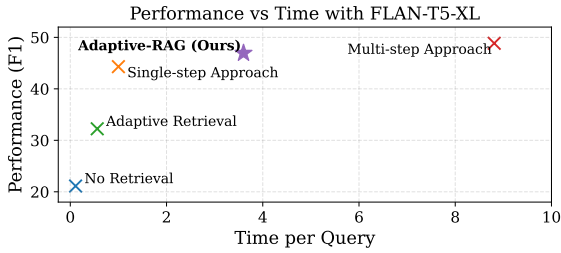


Figure 4: QA performance (F1) and efficiency (Time/Query) for different retrieval-augmented generation approaches. We use the FLAN-T5-XL (3B) as the base LLM.

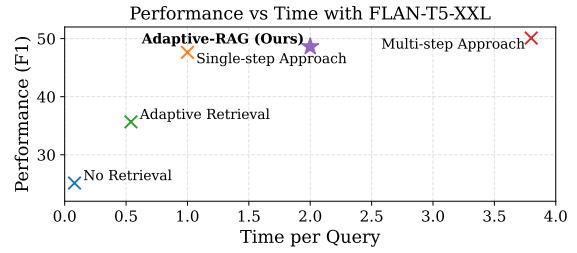


Figure 5: QA performance (F1) and efficiency (Time/Query) for different retrieval-augmented generation approaches. We use the FLAN-T5-XXL (11B) as the base LLM.

A Additional Experimental Setups

A.1 Datasets

We use publicly open datasets for both single-hop and multi-hop QA datasets, referring to [Karpukhin et al. \(2020\)](#) and [Trivedi et al. \(2023\)](#), respectively. We describe the characteristics of each dataset:

1) SQuAD v1.1 ([Rajpurkar et al., 2016](#)) is created through a process where annotators write questions based on the documents they read.

2) Natural Questions ([Kwiatkowski et al., 2019](#)) is constructed by real user queries on Google Search.

3) TriviaQA ([Joshi et al., 2017](#)) comprises trivia questions sourced from various quiz websites.

4) MuSiQue ([Trivedi et al., 2022a](#)) is collected by compositing multiple single-hop queries, to form queries spanning 2-4 hops.

5) HotpotQA ([Yang et al., 2018](#)) is constructed by having annotators create questions that link multiple Wikipedia articles.

6) 2WikiMultiHopQA ([Ho et al., 2020](#)) is derived from Wikipedia and its associated knowledge graph path, needing 2-hops.

A.2 Models

We describe the details of models as follows:

1) No Retrieval. This approach uses only the LLM itself, to generate the answer to the given query.

2) Single-step Approach. This approach first retrieves the relevant knowledge with the given query from the external knowledge sources and then augments the LLM with this retrieved knowledge to generate the answer, which iterates only once.

3) Adaptive Retrieval. This baseline ([Mallen et al., 2023](#)) adaptively augments the LLM with the retrieval module, only when the entities appearing in queries are less popular. To extract entities, we use the available entity-linking method ([Li et al., 2020](#)), namely BLINK, for questions.

4) Self-RAG. This baseline ([Asai et al., 2023](#)) trains the LLM to adaptively perform retrieval and

generation, where the retrieval is conducted once it predicts the special retrieval token above a certain threshold, and the answer generation is followed.

5) Adaptive-RAG. This is our model that adaptively selects the retrieval-augmented generation strategy, smoothly oscillating between the non-retrieval, single-step approach, and multi-step approaches⁴ without architectural changes, based on the query complexity assessed by the classifier.

6) Multi-step Approach. This approach ([Trivedi et al., 2023](#)) is the multi-step retrieval-augmented LLM, which iteratively accesses both the retriever and LLM with interleaved Chain-of-Thought reasoning ([Wei et al., 2022b](#)) repeatedly until it derives the solution or reaches the maximum step number.

A.3 Implementation Details

For computing resources, we use A100 GPUs with 80GB memory. In addition, due to the significant costs associated with evaluating retrieval-augmented generation models, we perform experiments with a single run. Finally, we implemented models using PyTorch ([Paszke et al., 2019](#)) and Transformers library ([Wolf et al., 2020](#)).

B Additional Experimental Results

Performance vs Time We further provide a comparison of different retrieval-augmented generation approaches with FLAN-T5-XL and FLAN-T5-XXL models in Figure 4 and Figure 5, respectively, in the context of performance and efficiency trade-offs. Similar to the observation made from the GPT-3.5 model in Figure 1, our proposed Adaptive-RAG is significantly more effective as well as efficient.

Performance per Dataset In addition to detailing the performance of each dataset with the FLAN-T5-XL model, as shown in Table 2, we also present

⁴For the multi-step approach, we use the state-of-the-art question answering strategy from IRCoT ([Trivedi et al., 2023](#)).

Table 6: Results on each of a collection of datasets with FLAN-T5-XXL (11B) as the LLM. We emphasize our results in bold.

Data	Types	Methods	SQuAD				Natural Questions				TriviaQA			
			EM	F1	Step	Time	EM	F1	Step	Time	EM	F1	Step	Time
Single-step	Simple	No Retrieval	7.00	14.40	0.00	0.08	18.80	25.50	0.00	0.08	32.80	39.20	0.00	0.08
		Single-step Approach	28.80	40.80	1.00	1.00	41.40	51.20	1.00	1.00	56.00	64.70	1.00	1.00
	Adaptive	Adaptive Retrieval	15.60	25.60	0.50	0.54	31.00	39.70	0.50	0.54	44.80	52.20	0.50	0.54
		Self-RAG [*]	1.60	11.90	0.59	0.31	39.20	47.10	0.75	0.09	14.60	33.70	0.76	0.22
		Adaptive-RAG (Ours)	27.80	39.80	1.17	1.50	41.20	51.00	1.00	1.00	52.00	60.30	1.03	1.33
Complex	Multi-step Approach	24.60	36.90	2.13	3.83	39.60	49.60	2.16	3.94	52.60	61.10	2.17	4.03	

Data	Types	Methods	MuSiQue				HotpotQA				2WikiMultiHopQA			
			EM	F1	Step	Time	EM	F1	Step	Time	EM	F1	Step	Time
Multi-step	Simple	No Retrieval	4.20	13.40	0.00	0.08	17.40	25.44	0.00	0.09	26.80	32.93	0.00	0.08
		Single-step Approach	16.80	25.70	1.00	1.00	37.60	49.27	1.00	1.00	46.60	54.13	1.00	1.00
	Adaptive	Adaptive Retrieval	8.40	17.80	0.50	0.54	26.60	36.01	0.50	0.54	35.20	42.68	0.50	0.54
		Self-RAG [*]	1.20	8.20	0.68	0.27	5.60	17.86	0.76	0.26	3.00	19.14	0.90	0.25
		Adaptive-RAG (Ours)	20.60	28.50	1.89	3.12	44.20	54.78	1.58	2.53	47.60	57.36	1.46	2.55
Complex	Multi-step Approach	19.40	27.50	2.09	3.66	47.00	57.81	2.08	3.73	57.60	67.65	2.17	3.63	

Table 7: Results on each of a collection of datasets with GPT-3.5 (Turbo) as the LLM. We emphasize our results in bold.

Data	Types	Methods	SQuAD				Natural Questions				TriviaQA			
			EM	F1	Step	Time	EM	F1	Step	Time	EM	F1	Step	Time
Single-step	Simple	No Retrieval	16.00	29.20	0.00	0.62	39.80	55.70	0.00	0.56	64.00	75.60	0.00	0.68
		Single-step Approach	18.00	33.80	1.00	1.00	32.40	46.80	1.00	1.00	55.20	66.50	1.00	1.00
	Adaptive	Adaptive Retrieval	15.40	30.00	0.50	0.81	36.40	51.20	0.50	0.78	62.00	71.90	0.50	0.84
		Self-RAG [*]	1.60	11.90	0.59	1.91	39.20	47.10	0.75	0.52	14.60	33.70	0.76	1.59
		Adaptive-RAG (Ours)	19.80	34.40	0.87	1.21	36.80	52.00	0.68	0.86	62.40	73.80	0.22	0.79
Complex	Multi-step Approach	17.40	31.50	2.50	3.24	35.60	49.70	2.58	3.79	54.80	67.10	2.30	2.65	

Data	Types	Methods	MuSiQue				HotpotQA				2WikiMultiHopQA			
			EM	F1	Step	Time	EM	F1	Step	Time	EM	F1	Step	Time
Multi-step	Simple	No Retrieval	20.40	31.30	0.00	0.81	37.40	51.04	0.00	0.74	37.00	48.50	0.00	0.90
		Single-step Approach	16.40	26.70	1.00	1.00	39.60	50.44	1.00	1.00	46.80	57.69	1.00	1.00
	Adaptive	Adaptive Retrieval	18.80	30.30	0.50	0.90	38.60	50.70	0.50	0.87	44.20	55.11	0.50	0.95
		Self-RAG [*]	1.20	8.20	0.68	1.66	5.60	17.86	0.76	1.67	3.00	19.14	0.90	1.81
		Adaptive-RAG (Ours)	21.80	32.60	1.90	2.29	40.40	52.56	0.93	1.48	46.60	60.09	1.59	2.23
Complex	Multi-step Approach	23.00	32.50	3.41	3.61	45.80	58.36	2.73	3.18	52.20	66.08	3.36	3.35	

1104 the results for each dataset with the FLAN-T5-
1105 XXL and GPT-3.5 models in Table 2 and Table 7,
1106 respectively. The experimental results show that
1107 our Adaptive-RAG consistently balances between
1108 efficiency and accuracy. It is worth noting that
1109 while the GPT-3.5 model performs effectively in
1110 addressing straightforward queries even without
1111 document retrieval, it benefits significantly from
1112 our Adaptive-RAG in terms of effectiveness when
1113 solving complex multi-hop queries.