

ACCELERATING INFERENCE OF RETRIEVAL-AUGMENTED GENERATION VIA SPARSE CONTEXT SELECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) augmented with retrieval exhibit robust performance and extensive versatility by incorporating external contexts. However, the input length grows linearly in the number of retrieved documents, causing a dramatic increase in latency. In this paper, we propose a novel paradigm named Sparse RAG, which seeks to cut computation costs through sparsity. Specifically, Sparse RAG encodes retrieved documents in parallel, which eliminates latency introduced by long-range attention of retrieved documents. Then, LLMs selectively decode the output by only attending to highly relevant caches auto-regressively, which are chosen via prompting LLMs with special control tokens. It is notable that Sparse RAG combines the assessment of each individual document and the generation of the response into a single process. The designed sparse mechanism in a RAG system can facilitate the reduction of the number of documents loaded during decoding for accelerating the inference of the RAG system. Additionally, filtering out undesirable contexts enhances the model’s focus on relevant context, inherently improving its generation quality. Evaluation results on four datasets show that Sparse RAG can be used to strike an optimal balance between generation quality and computational efficiency, demonstrating its generalizability across tasks.

1 INTRODUCTION

Large language models (LLMs) have attracted increasing attention and exhibited impressive abilities to understand instructions and generate fluent outputs in natural language (Brown et al., 2020; Ouyang et al., 2022; Touvron et al., 2023; Team et al., 2023). Nevertheless, LLMs inevitably manifest hallucinations (Ji et al., 2023) due to their struggle with factual errors and inability to secure the accuracy of generated text solely by the parametric knowledge they encapsulate (Zhang et al., 2023; Muhlgray et al., 2023). Feeding the source of truth to LLMs in the format of retrieved context segments (Reid et al., 2024) alleviates this problem. The technique is widely known as Retrieval-Augmented Generation (RAG) (Lewis et al., 2020c; Li et al., 2022; Guu et al., 2020a).

Although the RAG framework is empirically shown to be effective, it can be expensive to scale up. This is because it requires prepending relevant documents retrieved from an external knowledge corpus to the queries (Guu et al., 2020a). As a result, the input length grows linearly in the number of documents, causing a dramatic increase in latency when using a standard Transformer whose latency scales quadratically with the input length. Some prior works such as Fusion-in-Decoder (FiD) (Izacard & Grave, 2020) and Parallel Context Windows (PCW) (Ratner et al., 2022) have proposed to alleviate this issue. Yet these methods fail to strike an optimal balance between generation quality and computational efficiency. FiD was originally designed for the encoder-decoder architecture, and thus is not compatible with currently prevalent decoder-only architectures without significant changes. While PCW can be applied to decoder-only LLMs, it only speeds up the model pre-filling and still incurs high latency since the whole context window cache is still being attended to when decoding each token. Moreover, the heavy reliance of generation on the retrieved knowledge raises significant concerns about the model’s behavior and performance in scenarios where retrieval may fail or return inaccurate results (Shi et al., 2023). A typical approach for mitigating this issue is to rely on an external classifier to rank or filter the documents before prepending them to the input (Yan et al., 2024), but this process requires extra model calls which adds new complexity to inference.

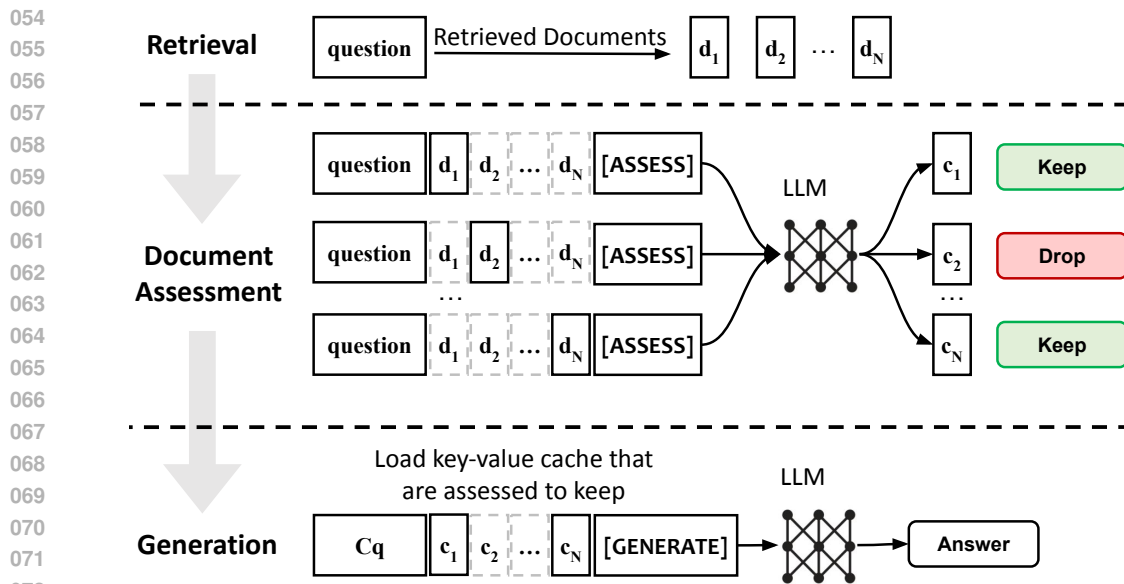


Figure 1: An overview of Sparse RAG at inference. Each of the retrieved documents is assessed for relevance by the LLM and irrelevant documents are dropped. Then, the remaining documents are used for generation.

In light of the issues above, we propose a novel paradigm called Sparse RAG. It operationalizes through massive pre-filling, where the key-value cache is generated by a single forward pass of the input tokens, and selective decoding, where the output is generated by attending to only highly relevant tokens auto-regressively. Previous works where the length of the retrieved contexts during pre-filling are equal to that during decoding are called *dense*-RAG in this paper. Sparse RAG, on the other hand, causes the decoding context to be significantly shorter than the pre-filled context, where retrieved documents that are not highly relevant to the input query have been dynamically dropped. Furthermore, Sparse RAG combines the assessment of each individual context and the generation of the response into a single process, in which special control tokens are used to prompt the LLM to assess the relevance of each retrieved context, and then only the key-value caches of the most relevant contexts are loaded for decoding using another control token.

The design of Sparse RAG has two additional unique advantages. First, by reducing the number of key-value cache loads during the decoding process, the LLM can achieve lower latency where it is typically constrained by memory usage. Second, filtering out undesirable contexts enhances the model’s focus on relevant contexts, inherently improving the quality of the generated output. To demonstrate the effectiveness and efficiency of the proposed method, we evaluate on four datasets: PopQA (Mallen et al., 2023), QMSum (Min et al., 2023), TriviaQA (Joshi et al., 2017), and HotpotQA (Yang et al., 2018). Experimental results show that Sparse RAG can achieve similar or better quality and much better latency compared with standard *dense*-RAG or PCW-RAG approaches. Moreover, the choice of the four datasets, which include short and long-form generation tasks, question answering and summarization tasks, and multi-hop reasoning cases, demonstrates the generalizability of the Sparse RAG approach.

2 RELATED WORK

Retrieval-Augmented Generation RAG is a family of techniques for generating output while using retrieved nearest-neighbor context data as a reference. It typically involves two stages: retrieval and generation. Retrieval finds most similar contexts based on BM25 or learned embeddings, where the context can be represented as token embeddings (Khandelwal et al., 2019; Yogatama et al., 2021), dense embeddings (De Jong et al., 2021) or raw text (Guu et al., 2020b; Izacard & Grave, 2020; Lewis et al., 2020b). Once those contexts are retrieved, different architectures are leveraged to incorporate them into the model. Popular approaches include concatenation (Izacard & Grave, 2020; Lewis et al., 2020b) and cross-attention (Borgeaud et al., 2022; Lewis et al., 2020a).

Table 1: Comparisons with existing RAG-related works.

Approach	Corrective	No extra model	Prefill efficiency	Decode efficiency
RAG (Lewis et al., 2020b)	No	Yes	No	No
Corrective RAG (Yan et al., 2024)	Yes	No	No	No
PCW RAG (Ratner et al., 2022)	No	Yes	Yes	No
Sparse RAG (Ours)	Yes	Yes	Yes	Yes

In recent years, LLM architectures have evolved towards decoder-only models with significantly larger sizes. To this end, concatenation of raw text (Lewis et al., 2020b) is becoming popular for its simplicity and practicality, and many advanced approaches have been developed on top of it. Yoran et al. (2023) designed an NLI model to identify irrelevant contexts and improve robustness. Jiang et al. (2023b) actively anticipate future content and decide when and what to retrieve in long-form generation. Self-RAG (Asai et al., 2023) is proposed to selectively retrieve knowledge on an as-needed basis, by introducing a separate critic model. The critic model generates "reflection" tokens to indicate whether to retrieve information. It runs inference on each document once and uses additional "reflection" tokens to select excerpts from the documents to use for generating the response. In contrast, we unify the generation of the special control tokens and regular vocabulary tokens with one single model, eliminating the additional model and computational overhead. CRAG (Yan et al., 2024) explores and designs corrective strategies for RAG to improve its robustness of generation. Specifically, an external T5 model is trained and used to determine the usefulness of the retrieved context. Generally, these approaches explore retrieval as a useful tool to augment generation and whether retrieval is necessary.

Efficiency in RAG The efficiency of LLM inference is a widely explored research area, where different categories of approaches have been studied, often targetting LLM inference in general rather than RAG specifically. Some works focus on architecture-level acceleration; examples include efficient attention (Shazeer, 2019), Mixture of Experts (Fedus et al., 2022), Transformer-alternative architectures (Gu & Dao, 2023), etc. Other works explore algorithm-level acceleration like quantization (Lin et al., 2023) or speculative decoding (Leviathan et al., 2023).

Recently, RAG-specific methods have been explored. RAG Cache (Jin et al., 2024), for example, was proposed as a multilevel dynamic caching system tailored for RAG from the system perspective. Another approach used in FiD (Izacard & Grave, 2020) and PCW (Ratner et al., 2022) parallelizes processing of individual documents and eliminates cross-document attention computations. FiD encodes each retrieved passage independently from other passages and decodes by attending over the concatenation of the resulting representations of all the retrieved passages. PCW carves a long context into chunks ("windows"), restricting the attention mechanism to apply only within each window, and re-uses the positional encodings across the windows.

Comparison with previous works that are the most relevant to our work is illustrated in Table 1. This work aims to strike an optimal balance between generation quality and computational efficiency. It is notable that the extra classifier in CRAG requires maintaining an extra model with more complex serving infrastructure; when there are N contexts retrieved, there are $N + 1$ model runs in total. Our work also relies on classification to refine the retrieved documents, but it is handled by an "internal" classification process that is aligned with the generation process, so the total number of model runs in our case is 1.

3 SPARSE RAG

Sparse RAG is designed for the decoder-only model architecture, which is the default case for most popular LLMs. Figure 1 presents an overview of Sparse RAG inference, in which document relevance assessment is used to improve the robustness of generation. The key hypothesis of our approach is that the RAG task and per context assessment are similar tasks and the model can handle both in one shot using simple and effective training and inference techniques.

3.1 TRAINING PROCESS

Our work assumes that a certain amount of RAG training data is accessible, which allows us to effectively tailor and adapt existing LLMs to our specific needs. In the training phase, we integrate an

162 additional Per Context Assessment (PCA) task into the training mixture. This PCA task is included
 163 alongside the primary RAG task to enhance the model’s ability to assess and respond accurately in
 164 different RAG scenarios. By incorporating the PCA task, we aim to improve the overall performance
 165 and contextual understanding of the LLMs, ensuring they are more adept at handling a diverse range
 166 of scenarios presented during their use.

167 **Data Augmentation with LLMs** For typical RAG data, one question-answer pair can be mapped
 168 to multiple retrieved contexts using either BM25 or an existing stand-alone retriever. However, there
 169 are cases where no golden labels indicating the quality of every retrieved context is available.

170 To collect these missing labels, we leveraged two off-the-shelf LLMs (Anil et al., 2023; Team et al.,
 171 2023)–PALM2 and Gemini–to assess each context. We observe empirically that a second round of
 172 prompting for critique, especially using a different model from the initial round, ensures the best
 173 quality labels. We provide our prompts in Table 10 in the Appendix. We compare different model
 174 combinations for labeling to human ground truth labels in Section 4.

175 **Multitasking Data Format** The LLM is trained on a mixture of two types of tasks: Per Context
 176 Assessment (relevance rating) and answer generation. Specifically, we format the inputs and outputs
 177 of the two task types as

- 179 • Per Context Assessment: {Question}{Context}{Control_Assessment}{Rating}
- 180 • Generation: {Question}{Context_1}...{Context_N}{Control_Generation}{Answer}

181 where {Rating} and {Answer} are targets for the generative tasks while all tokens before them are
 182 inputs. {Control_Assessment} and {Control_Generation} are fixed tokens to ensure the LLM can
 183 differentiate the two tasks.

184 **Parallel Contexts** Since each context is rated independently in the PCA task, we introduce this
 185 independence in the primary RAG task as well so that the two tasks can reuse the KV cache during
 186 inference time. Thus for the generation task, we enforce no cross-attention between different retrieved
 187 contexts, similar to Parallel Context Windows (Ratner et al., 2022). Specifically, we modify two
 188 things in the standard LM training process. First, we change the attention masks to be block-wise,
 189 and restrict that {Context_i} and {Context_j} will not attend to one another. {Question}, {Context_i}
 190 and {Control_Generation} still follow the default attention mechanism, i.e. causal attention, in which
 191 the latter tokens can attend to the previous ones. Second, we adjust the positional encoding to “mimic”
 192 the situation in which {Context_i} is not visible to {Context_j} so that each context will have its own
 193 incremental position id. However, the position id of {Control_Generation} is the normal position id
 194 accounting for all the lengths of contexts.

195 3.2 INFERENCE PROCESS

196 Given the question and retrieved contexts, Sparse RAG handles the assessment task and generation
 197 task in one single pass.

198 **Per Context Assessment** Similar to the training process, when pre-filling the KV cache, each
 199 retrieved context is treated independently. Then, these KV cache will be used to score the input,
 200 which is the concatenation of the {Control_Assessment} token and the {Relevant} token. The
 201 relevance score is the probability of this token as the next token.

202 **Generation** The generation uses a filtered KV cache, where only K out of N cached values are
 203 loaded. We use a simple threshold-based filtering approach: we drop the context when its score is
 204 less than σ . Once the cached KV vectors are loaded, Control_Generation prompts the model to
 205 generate the answer.

206 4 EVALUATION OF PER CONTEXT ASSESSMENT

207 4.1 NEW BENCHMARK: NATURAL QUESTION PER CONTEXT ASSESSMENT

208 We isolated a subset of 50 questions, each with 10 retrieved contexts, from the Natural Questions
 209 dataset. We assigned 3 raters to each question-context pair from a pool of 7 raters and provided the
 210 instructions in Section A.2.

Table 2: Auto-rater comparison to ground truth.

Auto-labeling method		Average F1	F1 Label 0	F1 Label 1
Rater model	Critic model			
PALM2 XL	n/a	0.729	0.765	0.694
PALM2 XL	PALM2 XL	0.781	0.820	0.741
Gemini Ultra	n/a	0.761	0.807	0.716
Gemini Ultra	Gemini Ultra	0.704	0.747	0.660
PALM2 XL	Gemini Ultra	0.728	0.776	0.680
Gemini Ultra	PALM2 XL	0.821	0.861	0.781

We aggregated responses for all 3 raters for each context, selecting the majority decision 0 or 1 for each context. We found that raters unanimously agreed on 351 out of 500 context, with 30% of the documents considered relevant. For questions where raters were not unanimously decided, a specialist rater was assigned to investigate more carefully and set the best label to correct mistakes of the other raters. This resulted in 6 additional documents considered relevant out of the entire dataset, boosting the portion of relevant documents to 31% and slightly increasing alignment with the auto-rater approaches (average F-score increase of 1.4% across auto-rater methods using these corrections as the ground truth).

4.2 LLM RATER COMPARISONS

We tested several different LLM-based automatic labeling methods—different combinations of models and prompts—for creating training data for the classifier in Sparse RAG. We compared several of these auto-rater approaches by creating a ground-truth relevance dataset using human labeling. The auto-rater comparison using the revised human labels as the ground-truth is shown in Table 2. We find that combining two different models in two rounds—initial prompting and critique—provides the labels that are most closely aligned with the human labels. We hypothesize that the different representations learned by two different models are able to capture the most nuance in the input sequences, leading to better relevance judgements. We also observe that Gemini Ultra appears slightly less effective at critiquing model outputs than PALM2 XL.

5 EVALUATION OF SPARSE RAG

5.1 BENCHMARKS AND METRICS

PopQA is a large-scale open-domain question answering (QA) dataset, consisting of 14k entity-centric QA pairs. Each question is created by converting a knowledge tuple retrieved from Wikidata using a template. We follow the setup from (Yan et al., 2024) and use Contriever (Izacard et al., 2022) to retrieve the related contexts. Since PopQA does not include per-context assessment relevance labels, we adopted the “Gemini + PALM2” combination to create training labels. We split the dataset into training, validation and test sets with 8:1:1 ratio. Since the answer is usually short, we report **Exact Match (EM)** and **F1** scores.

QMSum (Zhong et al., 2021) is a human-annotated benchmark for a query-based multi-domain meeting summarization task, which consists of 1,808 query-summary pairs over 232 meetings in multiple domains. To adapt it to the RAG domain, we divide each conversation into different contexts where each turn in the conversation is a context and the average context contains 300 words. Note that this dataset has human labeled per-context assessments that we leverage during training. Eventually we get 250 training examples (one per meeting), 70 validation examples and 77 test examples. The targets for this dataset are longer and we report **RougeLSum** and **F1** scores.

TriviaQA (Joshi et al., 2017) is a realistic text-based question answering dataset that includes 950K question-answer pairs from 662K documents collected from Wikipedia and the web. Similar to PopQA, we used the “Gemini + PALM2” combination to create relevance training labels. We randomly selected 8k training examples and 500 validation and test examples each. We report **Exact Match (EM)** and **F1** scores.

HotpotQA (Yang et al., 2018) is a question answering dataset containing about 113K crowd-sourced questions that are constructed to require the introduction paragraphs of at least two Wikipedia articles to answer (multi-hop reasoning). We sample 6k training examples and 600 validation and 600 test examples. We report **Exact Match (EM)** and **F1** scores.

5.2 BASELINES

RAG We evaluated the performance of standard concatenation-based RAG where an LLM generates output given the query prepended with all the top-ranked documents using the same retriever as Sparse RAG system. RAG is finetuned with the training data.

Off-the-shelf We report a variant of concatenation-based RAG where the model is not finetuned with training data.

LLMLingua In this approach an external LLM was called to compress the prompt (Jiang et al., 2023a). In our comparison, we chose the compression ratio to be the same as Sparse RAG for fairness.

PCW RAG We applied Parallel Context Windows (Ratner et al., 2022) to the RAG process, where no cross-attention is applied between documents. The model is finetuned with the training data.

Corrective RAG We evaluate CRAG using an external T5-XXL classifier trained using heuristic labels (Yan et al., 2024). This classifier is used to process all the documents and decide the rank. Note to facilitate a fair comparison, we did not adopt the "web search" feature of this paper.

5.3 EXPERIMENTAL CONFIGURATION

The base LLMs used in the paper were Gemini (Team et al., 2023). Although our approach could be applied at different training stages of the model, we apply LoRA tuning (Hu et al., 2021) to enforce alignment on top of the foundation LLMs due to its low resource requirements and wide usage. Note that the same LoRA tuning on the training data is applied to Sparse RAG and all baselines. In all our experiments, we apply LoRA in self-attention and use the default rank as 4. By default, we use the XXS size of Gemini which can run on-device.

During training, we use 64 Tensor Processing Units (TPU) V3 chips for PopQA while use 128 Units for the other datasets. The batch size is 64. We use the Adafactor optimizer (Shazeer & Stern, 2018) with a learning rate of 0.003. The training dropout rate is 0.05. We leverage the metrics of the validation set to pick the best checkpoint. During inference, the temperature is set to 0.5. Unless specifically noted, we use sampling decoding with sample number 1 for our experiments.

5.4 INFERENCE SETUP AND METRICS

Evaluation of Sparse RAG was conducted on a Samsung S21 Ultra, utilizing the device's CPU to assess real-world performance on a relatively mid-tier smartphone compared to the latest flagship models. Inference configuration consisted of fixed token lengths for queries, contexts and generated responses. This setup allows for evaluating the system's efficiency and effectiveness under resource constraints typical of mobile devices, providing insights into its practical applicability for on-device question-answering tasks. Specifically, the overall inference process considers two stages.

Prefill stage For the baseline RAG model, we measure the total time taken to process all input tokens (question and all contexts). For PCW RAG and Sparse RAG models, we take advantage of these models' ability to cache the question KV vectors. We first measure the time to process the question alone. Then, we measure the time to process different contexts with the pre-processed KV-cached question. We use **Encoding Speed (ES)** (tokens per second (t/s)) to measure the efficiency of the prefill stage.

Decoding stage To assess the decoding speed comprehensively, we generate output sequences using the same question but varying the amount of relevant context data considered, ranging from the top 1 most relevant document to the top K . For each context size, we produce output sequences of different lengths. This systematic approach allows us to evaluate the impact of both context size and response length on decoding speed. We use **Decoding Speed (DS)** (tokens per second (t/s)) to measure the efficiency of the decoding stage.

Table 3: Quality & efficiency tradeoff for both short-form and long-form generation tasks; Sparse RAG achieves both higher quality and efficiency compared to “dense” RAG approaches.

Dataset	Metrics	Off-the-shelf	LLMLingua	RAG	PCW-RAG	CRAG	Sparse RAG
-	ES	56.28	-	56.28	147.58	-	147.58
PopQA	EM	0.33	1.96	65.43	65.04	66.52	67.71
	F1	12.76	12.15	69.99	69.54	70.99	71.16
	K	20.00	7.84	20.00	20.00	8.9	7.84
	DS	6.65	-	6.65	6.65	-	12.28
QMSum	F1	20.37	22.28	21.43	20.18	-	23.96
	RougeSum	12.67	18.37	18.20	16.95	-	20.10
	K	20.00	4.45	20.00	20.00	-	4.45
	DS	6.65	-	6.65	6.65	-	16.05
TriviaQA	EM	0.00	2.6	46.20	46.00	-	47.50
	F1	12.21	16.30	53.03	53.20	-	55.10
	K	20.00	9.90	20.00	20.00	-	9.90
	DS	6.65	-	6.65	6.65	-	10.18
HotpotQA	EM	0.00	1.33	43	38.83	-	43.50
	F1	12.21	14.67	55.85	50.03	-	55.36
	K	10.00	6.50	10.00	10.00	-	6.50
	DS	10.31	-	10.31	10.31	-	13.00

5.5 MAIN RESULTS

We report the results in Table 3, where both quality and latency metrics are shared. “K” is the number of chosen contexts. Note that the CRAG approach relies on its classifier, which is exclusively trained on the PopQA dataset. Thus we only compare its performance on PopQA. Additionally, both LLMLingua and CRAG leveraged external classifiers, for which we cannot effectively measure ES and DS. We discuss the end-to-end latency of the external classifiers in Table 7.

Notably, our proposed approach achieves the best quality while being the most efficient during inference compared to other approaches. It can be seen that Sparse RAG shares the same pre-filling efficiency with PCW, due to the parallel context encoding, but it achieves significantly better decoding efficiency than standard RAG and PCW RAG. To illustrate, out of 20 retrieved contexts, Sparse RAG has an average of 7.84 contexts for PopQA and 4.45 contexts for QMSum. This leads to almost **double** or even **triple** the decoding speed. Meanwhile, Sparse RAG achieves higher quality metrics than the dense counterparts, demonstrating that Sparse RAG effectively filters noisy and irrelevant contexts.

It is interesting to compare Corrective-RAG (CRAG) with Sparse RAG on PopQA since Corrective-RAG trained their own external classifier with T5 XXL based on the PopQA data. Results suggest that our approach using an “in-place” classifier is outperforming CRAG with an external classifier. Note the encoding and decoding speed of CRAG are not comparable because it involves multiple model runs through the classifier.

For HotpotQA, PCW-RAG has significantly lower quality than RAG. This suggests that the parallel context window approach is not well designed for multi-hop reasoning cases. Sparse RAG, on the other hand, boosts the performance to a level that is similar to RAG while maintaining lower latency, demonstrating the power of the context selection process.

5.6 ANALYSIS

Impact of Confidence Threshold Table 4 illustrates how our metrics vary with different quality thresholds for Sparse RAG. As the threshold gradually increases, the system filters out more contexts, reducing the latency during inference. The response quality metrics increase with increasing threshold up to a certain point, showing the effectiveness of filtering out irrelevant contexts. Then, the performance is stable and eventually drops slightly, possibly because some relevant contexts are accidentally filtered out.

Table 4: Sampling a few confidence threshold values on PopQA and QMSum. A higher confidence threshold means filtering more memories. This always improves efficiency and improves quality up to a certain point.

Threshold	PopQA				QMSum			
	EM	F1	K	DS	F1	RougeLSum	K	DS
0.05	66.95	70.97	9.75	10.61	22.85	19.49	7.92	11.92
0.1	66.84	70.66	8.72	11.70	23.78	19.98	6.68	12.89
0.15	67.17	71.16	7.84	12.28	23.43	19.66	5.77	13.01
0.2	66.77	70.54	7.13	12.88	23.2	19.79	5.05	14.54
0.25	65.75	69.64	6.56	13.00	23.96	20.1	4.45	16.05
0.3	63.86	68.2	5.98	13.08	23.84	19.99	3.93	16.38

Silver Labels vs LLM Labels In Corrective RAG, the T5 model was trained with silver labels that come from title matching (Yan et al., 2024). We collect the same silver labels to replace the LLM labels and train the Sparse RAG model with this new dataset. From the results shown in Table 5, we observe that the quality of the labels generated by the LLMs is slightly higher than that of the silver labels from Yan et al. (2024). We hypothesize that the superior quality of the LLM-generated labels is attributable to our method, which involved a two-round process of soliciting responses from two different LLMs. By engaging two distinct models, we likely enhanced the robustness and accuracy of the labels through a form of cross-validation, thereby mitigating potential biases or errors that might arise from relying on a single LLM.

Table 5: Comparing CRAG silver labels with our LLM labels on PopQA.

Approach	EM	F1	K	DS
Sparse RAG	67.17	71.16	7.84	12.28
w/ silver label	66.97	71.05	8.26	11.99

Ablation of Prefill Documents To assuage concerns that Sparse RAG quality improvements are the result of “massive” prefilling which is not practical in real scenarios, we compare different numbers of prefill documents (10 and 20) for PopQA. Results are shown in Table 6. Even with fewer documents prefilled, the conclusion of Sparse RAG holds. It achieves better quality with higher decoding speed. It makes sense that the gap between dense RAG and Sparse RAG is relatively small at 10 prefill documents. This is because the headroom is smaller as the top 10 documents naturally has less to filter, making it less “sparse”. Our approach still achieves best latency-quality trade-off, showing the capability to generalize with different numbers of input documents. Meanwhile, only using the top 5/3/1 documents introduces significantly lower EM and F1 scores. The reason is that it is usually difficult to guarantee high-quality retrieval in the first step. This motivates our design to increase the range of retrieved documents and then perform sparse context selection.

Table 6: Ablation on different number of prefill documents for PopQA.

Approach	Prefill Documents	EM	F1	K	ES	DS
RAG	1	146.10	46.01	50.12	1.00	22.74
RAG	3	120.36	55.28	59.32	3.00	18.96
RAG	5	102.51	58.66	63.49	5.00	15.98
RAG	10	64.66	68.67	10	80.74	10.31
PCW RAG	10	63.9	68.58	10	147.48	10.31
Sparse RAG	10	65.86	70.2	7.79	147.48	12.33
RAG	20	65.43	69.99	20	56.28	6.65
PCW RAG	20	65.04	69.95	20	147.58	6.65
Sparse RAG	20	67.17	71.16	7.84	147.58	12.28

Table 7: Latency Decomposition.

End-To-End	K	External Classifier (ms)	Init Time (ms)	Encoding (ms)	Copy (ms)	Decoding (ms)	Total (ms)	Init + Copy Percentage
RAG	20.00	0	120	90962	0	4811	95893	0.13%
PCW-RAG	20.00	0	17	34716	151	4811	39697	0.43%
CRAG	8.9	40200	56	27362	0	2878	70497	0.08%
Sparse RAG	7.84	0	17	34716	56	2605	37396	0.20%

Inference Efficiency Ablations To provide a comprehensive analysis of the inference efficiency of Sparse RAG, we conducted ablation studies focusing on key factors that influence computational cost: memory length, number of memories, and decoding length. These factors are intrinsically linked to the complexity of the retrieval and generation process, and understanding their impact on performance is crucial to optimizing the system for real-world applications.

While it usually makes sense to just compare the total time elapsed, our situation is a bit different. Our approach is much faster at encoding than the original RAG, and it is faster at decoding than PCW-RAG, so it makes sense to compare the 2 critical intermediate steps separately.

To visualize the comparison of latency decomposition across prominent RAG approaches, we present the above Table 7 detailing the computational demands associated with each stage of the retrieval and generation process. While initialization and copying times can vary slightly across RAG approaches, our analysis shows these differences have a minimal impact on overall performance, contributing only around 0.2% to the total retrieval and generation time. For instance, the increased initial memory allocation in traditional RAG for multiple contexts and the time spent copying cached memory in methods like PCW-RAG are dwarfed by the encoding and decoding durations.

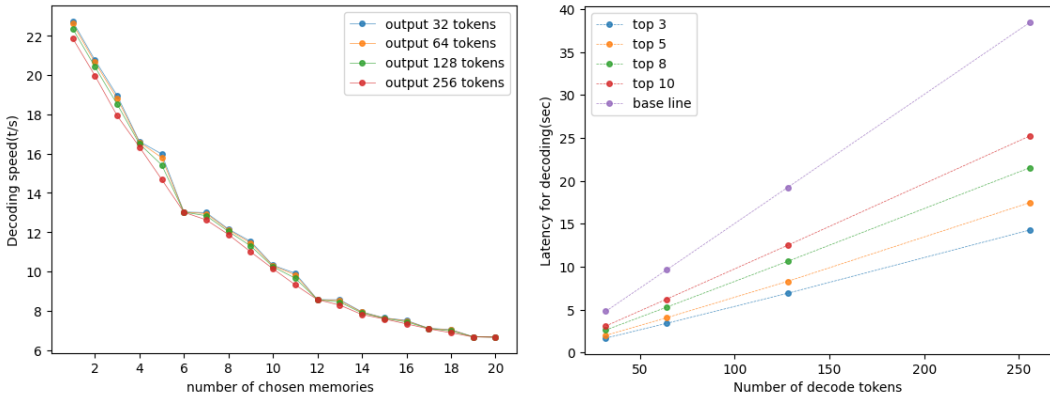
Focusing on the decoding stage, we delve into the performance differences when varying the number of retrieved memories (i.e., top-K documents) and the length of the generated response. This analysis provides valuable insights into the computational overhead associated with incorporating additional contexts and generating more extensive outputs, ultimately shedding light on the practical trade-offs for real-time question answering on resource-constrained devices.

As illustrated in Fig 2b, we observe that the baseline method requires over 50% (varies from 52% to over 55% depending on number of output tokens) more time to generate outputs of varying lengths compared to the Sparse RAG approach. This significant disparity in generation speed is primarily attributed to the increased computational burden incurred by the baseline method when considering a larger number of memories. As shown in Fig 2a, this heightened demand for computational resources results in a notable slowdown in terms of tokens per second (t/s). This observation underscores the efficiency advantages offered by Sparse RAG, especially in scenarios where a substantial amount of context needs to be considered during the decoding process.

Ablation on Foundation Model Size We applied Sparse RAG to different sizes of LLMs by testing it on Gemini XS and Gemini XXS. The results of these experiments are presented in Table 8. The findings demonstrate that Sparse RAG is compatible with various foundation models, effectively adapting to different model sizes. Notably, with a reduced amount of decoding caches, Sparse RAG is capable of achieving the highest quality results. This indicates that Sparse RAG maintains its efficiency and effectiveness across different foundation models, making it a versatile approach for various LLM configurations.

Using Golden Context Label during inference Since QMSum provides golden per-context labels, we leverage these labels during inference to evaluate the upper bound performance under the condition of perfect per-context assessment. By utilizing the highest quality labels available, we aim to determine the best possible outcomes that our model can achieve. The results of this experiment, which highlight the performance ceiling under ideal labeling conditions, are presented in Table 9. This approach allows us to understand the maximum potential of our model when supplied with optimal input data, thereby offering insights into its ultimate capabilities.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539



(a) Decoding speed with different number of contexts (b) E2E Latency for decoding different number of tokens

Figure 2: Inference Efficiency Comparison.

Table 8: Ablation on different model sizes.

Approach	Model Size	EM	F1	K
RAG	XS	66.52	70.87	20
PCW RAG	XS	65.75	70.37	20
Sparse RAG	XS	68.26	72.26	6.27
RAG	XXS	65.43	69.99	20
PCW RAG	XXS	65.04	69.95	20
Sparse RAG	XXS	67.17	71.16	7.84

Table 9: Trying golden labels on QMSum.

Approach	F1	RougeLSum	K	DS
Sparse RAG	23.96	20.1	4.45	16.05
+ golden label	26.76	21.93	1.13	21.16

6 CONCLUSION

This paper presents Sparse RAG to address the challenges of increased input length and latency. Through a novel approach of massive pre-filling and selective decoding, Sparse RAG efficiently manages the key-value cache of retrieved documents, allowing the LLMs to focus on highly relevant tokens. This selective attention mechanism not only reduces the computational burden during inference but also enhances the generation quality by filtering out irrelevant contexts. Evaluation on four diverse datasets validates Sparse RAG’s ability to achieve a balanced trade-off between high-quality generation and computational efficiency, proving its versatility and effectiveness for both short- and long-form content generation tasks. This innovative paradigm showcases the potential for improving LLM performance in various applications by optimizing context management and inference processes.

Under cases where per-context assessment label is missing, the prompts used for LLM rating may need to be adjusted for different use cases. These adjustments ensure that the model can still function effectively despite the lack of specific context labels. Furthermore, future research will explore Sparse RAG in multimodal contexts, investigating how Sparse RAG can handle and integrate information from multiple types of data to improve its performance and applicability across diverse scenarios.

REFERENCES

- 540
541
542 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Sia-
543 mak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El
544 Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick,
545 Kevin Robinson, Sebastian Ruder, et al. PaLM 2 technical report. *CoRR*, abs/2305.10403, 2023.
546 doi: 10.48550/ARXIV.2305.10403. URL [https://doi.org/10.48550/arXiv.2305.](https://doi.org/10.48550/arXiv.2305.10403)
547 10403.
- 548 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning
549 to retrieve, generate, and critique through self-reflection. *CoRR*, abs/2310.11511, 2023. doi: 10.
550 48550/ARXIV.2310.11511. URL <https://doi.org/10.48550/arXiv.2310.11511>.
- 551 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican,
552 George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al.
553 Improving language models by retrieving from trillions of tokens. In *International conference on*
554 *machine learning*, pp. 2206–2240. PMLR, 2022.
- 555 Tom B Brown, Benjamin Mann, Nick Ryder, et al. Language models are few-shot learners. In
556 *Advances in neural information processing systems*, pp. 1877–1901, 2020.
- 557
558 Michiel De Jong, Yury Zemlyanskiy, Nicholas FitzGerald, Fei Sha, and William Cohen. Mention
559 memory: incorporating textual knowledge into transformers through entity mention attention.
560 *arXiv preprint arXiv:2110.06176*, 2021.
- 561
562 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter
563 models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39,
564 2022.
- 565 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
566 *preprint arXiv:2312.00752*, 2023.
- 567
568 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Retrieval aug-
569 mented language model pre-training. In *Proceedings of the 37th International Conference on*
570 *Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of*
571 *Machine Learning Research*, pp. 3929–3938. PMLR, 2020a. URL [http://proceedings.](http://proceedings.mlr.press/v119/guu20a.html)
572 [mlr.press/v119/guu20a.html](http://proceedings.mlr.press/v119/guu20a.html).
- 573 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented
574 language model pre-training. In *International conference on machine learning*, pp. 3929–3938.
575 PMLR, 2020b.
- 576
577 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
578 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
579 *arXiv:2106.09685*, 2021.
- 580
581 Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open
582 domain question answering. *arXiv preprint arXiv:2007.01282*, 2020.
- 583
584 Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand
585 Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning.
586 *Trans. Mach. Learn. Res.*, 2022, 2022. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=jKN1pXi7b0)
587 [jKN1pXi7b0](https://openreview.net/forum?id=jKN1pXi7b0).
- 588
589 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea
590 Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Comput.*
591 *Surv.*, 55(12):248:1–248:38, 2023. doi: 10.1145/3571730. URL [https://doi.org/10.](https://doi.org/10.1145/3571730)
592 1145/3571730.
- 593
594 Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. LlmLingua: Compressing
595 prompts for accelerated inference of large language models. *arXiv preprint arXiv:2310.05736*,
596 2023a.

- 594 Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang,
595 Jamie Callan, and Graham Neubig. Active retrieval augmented generation. In Houda Bouamor,
596 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in*
597 *Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 7969–7992.
598 Association for Computational Linguistics, 2023b. URL [https://aclanthology.org/](https://aclanthology.org/2023.emnlp-main.495)
599 [2023.emnlp-main.495](https://aclanthology.org/2023.emnlp-main.495).
- 600
601 Chao Jin, Zili Zhang, Xuanlin Jiang, Fangyue Liu, Xin Liu, Xuanzhe Liu, and Xin Jin. Raggache:
602 Efficient knowledge caching for retrieval-augmented generation. *arXiv preprint arXiv:2404.12457*,
603 2024.
- 604
605 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly
606 supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
- 607
608 Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization
609 through memorization: Nearest neighbor language models. *arXiv preprint arXiv:1911.00172*,
610 2019.
- 611
612 Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative
613 decoding. In *International Conference on Machine Learning*, pp. 19274–19286. PMLR, 2023.
- 614
615 Mike Lewis, Marjan Ghazvininejad, Gargi Ghosh, Armen Aghajanyan, Sida Wang, and Luke
616 Zettlemoyer. Pre-training via paraphrasing. *Advances in Neural Information Processing Systems*,
617 33:18470–18481, 2020a.
- 618
619 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
620 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-
621 tion for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:
622 9459–9474, 2020b.
- 623
624 Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Na-
625 man Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel,
626 and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In
627 Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-
628 Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Confer-
629 ence on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,*
630 *2020, virtual*, 2020c. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html)
631 [6b493230205f780e1bc26945df7481e5-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html).
- 632
633 Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. A survey on retrieval-augmented text
634 generation. *CoRR*, abs/2202.01110, 2022. URL <https://arxiv.org/abs/2202.01110>.
- 635
636 Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation-
637 aware weight quantization for llm compression and acceleration. *arXiv preprint arXiv:2306.00978*,
638 2023.
- 639
640 Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi.
641 When not to trust language models: Investigating effectiveness of parametric and non-parametric
642 memories. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the*
643 *61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
644 *ACL 2023, Toronto, Canada, July 9-14, 2023*, pp. 9802–9822. Association for Computational
645 Linguistics, 2023. doi: 10.18653/V1/2023.ACL-LONG.546. URL [https://doi.org/10.](https://doi.org/10.18653/v1/2023.acl-long.546)
646 [18653/v1/2023.acl-long.546](https://doi.org/10.18653/v1/2023.acl-long.546).
- 647
648 Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer,
649 Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual
650 precision in long form text generation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.),
651 *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing,*
652 *EMNLP 2023, Singapore, December 6-10, 2023*, pp. 12076–12100. Association for Computational
653 Linguistics, 2023. URL <https://aclanthology.org/2023.emnlp-main.741>.

- 648 Dor Muhlgay, Ori Ram, Inbal Magar, Yoav Levine, Nir Ratner, Yonatan Belinkov, Omri Abend,
649 Kevin Leyton-Brown, Amnon Shashua, and Yoav Shoham. Generating benchmarks for factuality
650 evaluation of language models. *CoRR*, abs/2307.06908, 2023. doi: 10.48550/ARXIV.2307.06908.
651 URL <https://doi.org/10.48550/arXiv.2307.06908>.
- 652 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,
653 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser
654 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan
655 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In
656 *NeurIPS*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/hash/
657 b1efde53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html).
- 658 Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas,
659 Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. Parallel context windows for large
660 language models. *arXiv preprint arXiv:2212.10947*, 2022.
- 661 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste
662 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini
663 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint
664 arXiv:2403.05530*, 2024.
- 665 Noam Shazeer. Fast transformer decoding: One write-head is all you need. *arXiv preprint
666 arXiv:1911.02150*, 2019.
- 667 Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost.
668 In *International Conference on Machine Learning*, pp. 4596–4604. PMLR, 2018.
- 669 Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H. Chi, Nathanael
670 Schärli, and Denny Zhou. Large language models can be easily distracted by irrelevant context.
671 In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and
672 Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*,
673 volume 202 of *Proceedings of Machine Learning Research*, pp. 31210–31227. PMLR, 23–29 Jul
674 2023. URL <https://proceedings.mlr.press/v202/shi23a.html>.
- 675 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
676 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
677 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 678 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
679 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand
680 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language
681 models. *CoRR*, abs/2302.13971, 2023. doi: 10.48550/ARXIV.2302.13971. URL <https://doi.org/10.48550/arXiv.2302.13971>.
- 682 Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. Corrective retrieval augmented generation.
683 *arXiv preprint arXiv:2401.15884*, 2024.
- 684 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov,
685 and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question
686 answering. *arXiv preprint arXiv:1809.09600*, 2018.
- 687 Dani Yogatama, Cyprien de Masson d’Autume, and Lingpeng Kong. Adaptive semiparametric
688 language models. *Transactions of the Association for Computational Linguistics*, 9:362–373, 2021.
- 689 Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. Making retrieval-augmented language
690 models robust to irrelevant context. *CoRR*, abs/2310.01558, 2023. doi: 10.48550/ARXIV.2310.
691 01558. URL <https://doi.org/10.48550/arXiv.2310.01558>.
- 692 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lema Liu, Tingchen Fu, Xinting Huang, Enbo Zhao,
693 Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming
694 Shi. Siren’s song in the AI ocean: A survey on hallucination in large language models. *CoRR*,
695 abs/2309.01219, 2023. doi: 10.48550/ARXIV.2309.01219. URL [https://doi.org/10.
696 48550/arXiv.2309.01219](https://doi.org/10.48550/arXiv.2309.01219).

702 Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadallah,
703 Asli Celikyilmaz, Yang Liu, Xipeng Qiu, et al. Qmsum: A new benchmark for query-based multi-
704 domain meeting summarization. *arXiv preprint arXiv:2104.05938*, 2021.
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A PROMPTS AND INSTRUCTIONS

A.1 PROMPTS USED FOR LLMs

We share the prompt used for calling LLMs to get per context assessment in Table 10.

Table 10: The zero-shot prompts for LLM labeling and critique.

Round 1 prompt
<p>You are now doing a reading comprehension task. It is important that you be as thorough, detail-oriented, and accurate as possible in your response.</p> <p>You are given a question, a set of accepted answers, a document and its title. The document does not necessarily contain the right answer to the question.</p> <p>You should read the title and the document and then check if they provide one of the correct answers to the question.</p> <p>If the title and document together contain the correct answer to the question, output a score of 1.0, otherwise output a score of 0.0.</p> <p>question: ;question; accepted answers: ;answers; title: ;title; document: ;document; output:</p>
Round 2 prompt
<p>Your job is to correct another model’s performance on a reading comprehension task.</p> <p>The model was given a question, a set of accepted answers, a document and its title. The document and title do not necessarily contain the right answer. The model was instructed to output a score of 1.0 if the document contains the answer, and a score of 0.0 otherwise.</p> <p>You will be given the same information as the other model along with its output. You should read the title and document and then check if they provide one of the correct answers to the question.</p> <p>Then check if you agree with the previous model’s output.</p> <p>If you agree, output the same score unchanged.</p> <p>If you disagree, output the corrected score.</p> <p>Your output should be as accurate as possible.</p> <p>question: ;question; accepted answers: ;answers; title: ;title; document: ;document; previous model’s score: ;score; output:</p>

A.2 RATER GUIDELINES

We share the instructions provided to the human labelers in Table 11.

Table 11: Instructions for raters creating ground-truth relevance dataset.

Human Rater Instructions
<p>Please read the question, the answer and the context. Please answer if the context can help answer the question. If it can, select 1. Otherwise select 0.</p> <p>1: good 0: bad</p> <p>Please use the answers as a hint. However, do not use ”is the answer in the context?” as a heuristic for making the decision.</p>

B DATASET ANALYSIS

During the human-labeling process, several raters flagged documents and questions that were difficult to label. In total, 23 out of 500 documents were flagged and 15 out of 50 questions were flagged.

We observed trends in questions and contexts shared in Table 12 that raised concerns about whether and how a human would be able to assess the relevance of the context. These concerns extend to expectations of how well LLMs would be able to do that as well.

Table 12: Overview of trends, datasets, and examples with associated comments.

Trend	Dataset	Example question	Comments
Questions with time-dependent answers	NQ	who is the president of usa right now	Depends on when the question is asked.
	NQ	who is the current director of the us mint	Depends on when the question is asked.
	NQ	when is the next dead-pool movie being released	Depends on when the question is asked.
	NQ	total number of death row inmates in the us	Fluctuates over time.
	PopQA	What is Prague the capital of?	The borders in this region and the name of the country have changed several times in the 20th century.
	PopQA	What is Dennis Rodman’s occupation?	The accepted answers are ”actor, actress, actors, actresses”. He was an actor later in his career, but he rose to prominence as a professional basketball player.
Missing synonyms, abbreviations, and aliases, combined with unclear granularity	NQ	in which regions are most of Africa petroleum and natural gas found	”Region” can refer to different levels of granularity (e.g. Sub-Saharan Africa vs. Ethiopia), but the only accepted answer is ”Nigeria”.
	NQ	what type of car is a jeep	The accepted answers are only ”off-road vehicles”, ”light utility vehicles”, ”sport utility vehicles”, but ”SUV” is clearly a correct answer as well.
Non-exhaustive list of answers	NQ	cast of law & order special victim unit	The accepted answers include 16 cast members, but the show went on for 25 seasons with many cast changes and guest stars not included in the list.
Oddly phrased question	NQ	right to property according to the constitution of India is a	The only correct answer is ”constitutional right”, but that is included in the question. It’s not clear what type of answer would be appropriate here.
Overly specific answers expected	NQ	where does the story the great gatsby take place	The only accepted answer here is ”Long Island of 1922”, but the place is Long Island and the question does not ask about when the story is set.
Question refers to an entity with a common name without disambiguation	PopQA	What genre is Frances?	There is a musician and a film called ”Frances”, and both of those could arguably have a genre associated with them.
	PopQA	Who was the producer of Hurt?	The question is referring to a song performed by Christina Aguilera but there are many other songs, movies, and other entities that share the name and also have a producer.
	PopQA	What is the capital of Cherokee County?	There are many different Cherokee Counties in different states in the USA.

We explored several ways of filtering our human-labeled subsample of Natural Questions to determine how they impacted context assessment F-scores overall and for each auto-rater. We provide two additional filtered versions of the human-labeled RAG relevance dataset as alternatives. See Table 13 for statistics on the filtered datasets and Table 14 for the auto-rater F-scores for each. Both statistical filtering approaches (e.g. removing contexts with non-unanimous labels) and targeted filtering approaches (e.g. removing questions or contexts flagged by human raters) lead to some improvement in F-scores for relevance labels.

Table 13: Summary of Dataset Evaluations

Dataset	# Questions	# Docs	Total % Relevant	Average # Docs per Question	Std Dev # Docs per Question	Average % relevant docs per question	Std Dev % relevant docs per question
Original 3-rater labels	50	500	0.30	10	0	0.30	0.23
Revisited (released)	50	500	0.31	10	0	0.31	0.26
Filter non-unanimous (released)	50	351	0.23	7.02	2.59	0.27	0.29
Filter flagged docs	48	467	0.31	9.729	0.57	0.31	0.27
Filter flagged questions	35	350	0.29	10	0	0.29	0.26
Filter flagged docs and question (released)	34	330	0.29	9.7	0.58	0.29	0.27
Filter 4-rater split	50	444	0.28	8.88	1.49	0.29	0.27

Table 14: Evaluation of Labeling Methods w/ Filtering

Dataset Filters	# Docs	Rater Model	Critic Model	Average F1
Specialized rater corrections	500	PALM2 XL	n/a	0.729
		PALM2 XL	PALM2 XL	0.781
		Gemini Ultra	n/a	0.761
		Gemini Ultra	Gemini Ultra	0.704
		PALM2 XL	Gemini Ultra	0.728
		Gemini Ultra	PALM2 XL 340B	0.821
Filter non-unanimous docs	351	PALM2 XL	n/a/	0.741
		PALM2 XL	PALM2 XL	0.811
		Gemini Ultra	n/a	0.792
		Gemini Ultra	Gemini Ultra	0.739
		PALM2 XL	Gemini Ultra	0.763
		Gemini Ultra	PALM2 XL	0.856
Filter flagged docs and questions	330	PALM2 XL	n/a	0.750
		PALM2 XL	PALM2 XL	0.797
		Gemini Ultra	n/a	0.782
		Gemini Ultra	Gemini Ultra	0.741
		PALM2 XL	Gemini Ultra	0.753
		Gemini Ultra	PALM2 XL	0.833