

# Unveiling the Magic: Investigating Attention Distillation in Retrieval-Augmented Generation

Zizhong Li Haopeng Zhang Jiawei Zhang

IFM Lab, University of California, Davis

{zzoli, hapzhang, jiwzhang}@ucdavis.edu

## Abstract

Retrieval-augmented generation framework addresses the limitations of large language models by enabling real-time knowledge updates for more accurate answers. An efficient way in the training phase of retrieval-augmented models is attention distillation, which uses attention scores as supervision signals instead of manually annotated query-document pairs. Despite its growing popularity, the detailed mechanisms behind the success of attention distillation remain unexplored, particularly the specific patterns it leverages to benefit training. In this paper, we address this gap by conducting a comprehensive investigation of attention distillation workflow and identifying key factors influencing the learning performance of retrieval-augmented language models. We further propose several insightful indicators for optimizing models' training methods and avoiding ineffective training.

## 1 Introduction

Large language models (LLMs) have showcased remarkable capabilities across various natural language processing tasks (Min et al., 2023; OpenAI, 2023; Ouyang et al., 2022; Zhang et al., 2023a,b). However, in the inference phase, their frozen parameters limit their ability to update knowledge in real-time, resulting in the hallucination problem during generation (Zhang et al., 2022, 2023c). Additionally, these models also lack protection for sensitive training data (Nasr et al., 2023; Lin et al., 2021). One promising method to overcome these limitations is using retrieval-augmented language models (Ram et al., 2023; Shi et al., 2022; Izacard et al., 2022b; Guu et al., 2020; Karpukhin et al., 2020; Khandelwal et al., 2019). Retrieval-augmented language models typically comprise two essential components: (1) *the retriever*, which selects relevant information, and (2) *the reader*, which incorporates this information into the gener-

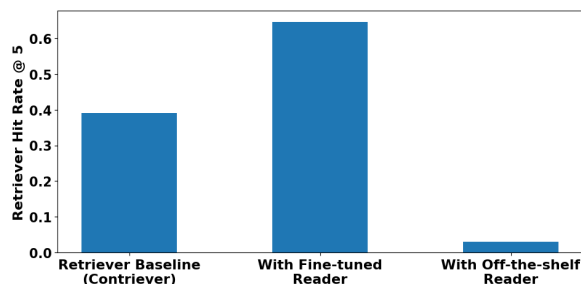


Figure 1: Training *Contriever* on *NaturalQuestions* for the QA task with attention distillation shows an improved Hit Rate @ 5 with a fine-tuned reader but a significant decline with an off-the-shelf reader.

ation process. The integration of these two components allows retrieval-augmented language models to enhance accuracy and reliability by dynamically utilizing external knowledge, while also reducing training costs due to fewer trainable parameters (Shi et al., 2023; Shuster et al., 2021).

The performance of retrieval-augmented language models may significantly depend on the effective synergy between the retriever and the reader. To this end, various methods have been proposed to improve the coordination between these two components (Karpukhin et al., 2020; Jiang et al., 2023). Among these, attention score-based knowledge distillation stands out due to its notable effectiveness in question-answering (QA) tasks (Izacard and Grave, 2020a), outperforming other established methods (Karpukhin et al., 2020; Lewis et al., 2020; Izacard and Grave, 2020b). In this process, the attention scores from the reader are captured and conveyed to the retriever as the supervisory signal (i.e., the retriever uses the attention scores as the basis for assessing the relevance of retrieved information), enabling the retrieval model to identify information candidates more effectively that can significantly improve the language model's responses. This efficient strategy reduces the need for manual annotation of the knowledge corpus, leading to resource savings while still achieving

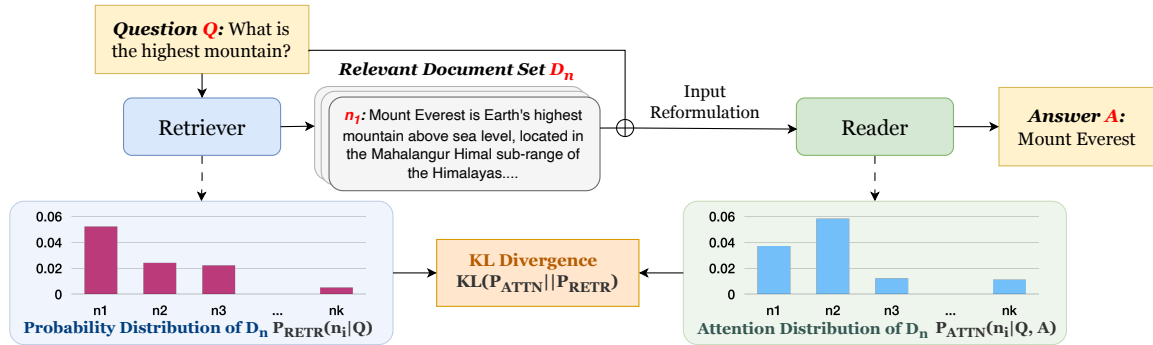


Figure 2: The framework of the Retrieval-augmented Language Model of our experiment.

satisfactory results (Hu et al., 2023; Wang et al., 2023).

However, the success of attention distillation heavily relies on the quality of the reader model. As shown in Figure 1, reader models of low quality yield ineffective supervision signals, detrimentally impacting the retriever’s performance. Given the critical nature of the issue, it becomes imperative to delve into the mechanism of attention distillation and identify characteristics of low-quality readers to avoid ineffective training.

A fundamental hypothesis underpinning attention distillation is that more attention to certain tokens suggests that these tokens are likely to be of greater relevance in answering questions (Izacard and Grave, 2020a). However, this correlation between attention scores and tokens has not to be clearly established yet, as existing works lack a quantitative analysis in attention scores’ impacts within the training process of retrieval-augmented language models. Therefore, our research seeks to understand which text segments receive more attention and how to assess the attention distillation quality.

In this paper, we first experimentally confirm that *attention scores are not always effective as training supervisors across different experimental settings under question-answering tasks*. Motivated by this observation, we conduct an in-depth token-level quantitative analysis, seeking to uncover patterns within attention scores that correspond to high-quality supervision. We analyze the attention scores from reader models of various qualities and identify a clear, stable correlation between these scores and supervisory quality, especially in tokens related to answers and questions. Building on these insights, we derived two key indicators to measure the distillation quality based on the commonalities. Our main contributions are as follows:

- We conduct an extensive analysis of attention scores in large language models, mainly focusing on the prevalent decoder-only structure, to understand their impacts on retriever model training and the overall performance of retrieval-augmented language models, thereby identifying key factors that significantly influence the model’s performance.
- We introduce novel metrics to evaluate the reader model’s proficiency in attention distillation, aiming to improve training performance by leaning on effective training sessions.

## 2 Method

In our experiment, we adapt the ATLAS architecture (Izacard et al., 2022b) but use a decoder-only language model structure for our empirical analysis, focusing on question-answering tasks to study attention score distillation mechanisms. Specifically, for a given question  $Q$ , we supply models with a knowledge base  $D = \{d_1, d_2, \dots, d_m\}$ , where each  $d_i$  is a unique document. The objective of the models is to find the question-relevant documents  $D_n = \{n_1, n_2, \dots, n_k\} \subseteq D$  using the retriever, and then incorporate the query and  $D_n$  as the input for the reader to generate the answer  $A$  for the given question.

The attention distillation approach uses attention scores to gauge the importance of each input document  $d_i$  during the answer generation process. To accommodate changes in the reader model’s structure, we utilize the *self-attention scores* related to the output tokens as a measure of document relevance, rather than relying on *cross-attention scores* between input documents and output that ATLAS uses. In addition, the attention level of a token  $t$  is not only evaluated from the self-attention score  $\alpha_t$  but also the norm of the value vector  $\mathbf{v}_t$  should be taken into account (Izacard et al., 2022b). Af-

terwards, the *Softmax* operator is applied to obtain the attention score distribution over the question-relevant documents  $D_n$ :

$$P_{ATTN}(n_i|Q, A) = \text{softmax}\left(\sum_{t=1}^T \alpha_t \|\mathbf{v}_t\|_2\right) \quad (1)$$

where  $T$  represents the total token count in  $n_i$ .

On the other hand, the retriever’s probability distribution  $p_{RETR}$  over  $D_n$  can be defined as:

$$P_{RETR}(n_i|Q) = \frac{\exp(s(n_i, Q)/\theta)}{\sum_{k=1}^K \exp(s(n_k, Q)/\theta)} \quad (2)$$

where  $s$  denotes the dot-product between the representation vectors of the input question  $Q$  and document candidate  $n_i$ , and  $\theta$  is the temperature hyper-parameter.

During training, the attention score’s distribution is distilled into the retriever by minimizing KL-divergence between  $P_{ATTN}(n_i)$  and  $P_{RETR}(n_i)$ , which aligns the retriever’s behavior more closely with the insights derived from attention scores. Figure 2 visually illustrates the retrieval process and the utilization of attention scores during training.

### 3 Experiments

We chose *Falcon-1b* (Penedo et al., 2023) as our primary decoder-only reader model for its performance and flexibility, and we follow ATLAS (Izacard et al., 2022b) in using *Contriver* as the retriever model. During the retrieval process, we set the number of retrieved documents  $D_n$  to a fixed size  $k = 5$  to balance training costs with the amount of information retrieved, thereby avoiding inefficiencies of either extreme.

#### 3.1 Experiment Setup

**Dataset** We assess the model’s performance using the *NaturalQuestions* (Kwiatkowski et al., 2019) and the *TriviaQA* (Joshi et al., 2017) benchmarks, which are the two most popular dataset in the QA task. For the knowledge base, we utilize data from Wikipedia as of December 20, 2018.

**Experimental Settings** Specifically, we use the following settings for our experiments.

**1) Off-the-shelf Distillation Training:** We synchronously train the model using the initial *Falcon-1b* (Penedo et al., 2023) as the reader and *Contriver* (Izacard et al., 2022a) as the retriever.

**2) Fine-tuned Distillation Training:** This experiment involves two steps:

Table 1: Model’s Performance of Different Experimental Settings

Method	Dataset	Evaluation Metrics		
		EM $\uparrow$	F1 $\uparrow$	HR@5 $\uparrow$
Off-the-shelf Distillation	NQ	27.24	33.62	0.030
	TriviaQA	30.55	35.24	0.022
Fine-tuned Distillation (Step1)	NQ	31.76	38.72	0.391
	TriviaQA	44.62	50.79	0.516
Fine-tuned Distillation (Step2)	NQ	35.22	43.44	0.645
	TriviaQA	54.59	61.04	0.643

**Step1.** We start with the initial *Falcon-1b* checkpoint as a reader and *Contriver* as a retriever, only fine-tuning the reader model while keeping the retriever model’s parameters fixed.

**Step2.** We continue training the retriever model using the fine-tuned reader checkpoint from Step1, updating the knowledge base index periodically.

**Evaluation Metrics:** We assess the model performance in terms of retrieval quality and question-answering correctness, given the involvement of both retriever and reader models. We use the *top-5* retrieval Hit Rate (HR@5), which is the proportion of retrieved documents  $D_n$  containing at least one answer  $A$ , to measure the retriever’s effectiveness. For the reader’s QA performance, we employ the standard Exact Match (EM) metric and F1-Score.

#### 3.2 Results and Discussion

In this section, we empirically analyze the effectiveness of attention distillation training by answering the following research questions:

**RQ1: When does the attention distillation work?**

As shown in Table 1, the *Fine-tuned Distillation Training* after Step2 shows the best performance in both EM, F1 and HR@5. In contrast, *Off-the-shelf Distillation Training* performs the worst, with its retriever even underperforming the initial *Contriver* model (i.e., the retriever model of *Fine-tuned Distillation Training* Step1). Notice that the critical difference lies in the quality of the reader models: *Off-the-shelf Distillation Training* uses the initial *Falcon-1b* model, whereas *Fine-tuned Distillation Training* employs a well-tuned *Falcon-1b*. These experimental results strongly suggest that the quality of attention scores is pivotal: **attention scores from the high-quality readers enhance training, whereas low-quality ones lead to poor interaction between the retriever and the reader.**

**RQ2: Are there any commonalities in attention scores from the high-quality readers?**

We sample 1000 data instances from each exper-

Table 2: Average values of attention scores and Spearman correlation in *answer-related* and *question-related* tokens

Experiment	Dataset	Answer-related				Question-related			
		90 <sup>th</sup> percentile		95 <sup>th</sup> percentile		90 <sup>th</sup> percentile		95 <sup>th</sup> percentile	
		Attn.	Corr.	Attn.	Corr.	Attn.	Corr.	Attn.	Corr.
Off-the-shelf Checkpoint	NQ	0.033	0.227	0.039	0.196	0.023	0.103	0.024	0.092
	TriviaQA	0.027	0.218	0.032	0.206	0.021	0.103	0.023	0.067
Off-the-shelf Attention Distillation	NQ	0.017	0.145	0.017	0.076	0.027	0.139	0.039	0.153
	TriviaQA	0.031	0.160	0.035	0.172	0.047	0.144	0.063	0.260
Fine-tuned Attention Distillation (Step1)	NQ	0.039	0.308	0.052	0.282	<b>0.035</b>	<b>0.343</b>	<b>0.045</b>	<b>0.333</b>
	TriviaQA	0.058	0.259	0.074	0.258	0.058	0.349	<u>0.078</u>	<u>0.372</u>
Fine-tuned Attention Distillation (Step2)	NQ	<b>0.049</b>	<b>0.316</b>	<b>0.066</b>	<b>0.350</b>	0.032	0.310	0.039	0.225
	TriviaQA	<u>0.069</u>	<u>0.290</u>	<u>0.089</u>	<u>0.320</u>	<u>0.060</u>	<u>0.367</u>	<u>0.078</u>	0.326

iment to obtain reliable analysis results. We focus on the attention score characteristics **at token level** to identify which tokens receive more attention from high-quality signals. Our analysis firstly finds that in the high-quality readers, the tokens most related to *answer* and *nouns in question* receive the most attention. Based on our initial observations, we secondly focus on studying the distribution of attention scores for *answer-related* and *question-related*<sup>1</sup> tokens. We use token embedding’s *cosine similarity* to measure its proximity to targets (i.e., answer or nouns in question), selecting the top 5% and top 10% of closest tokens and analyzing their average *attention scores* and *Spearman correlation with similarity to target tokens*, as shown in Table 2<sup>2</sup>. We also include the *Off-the-shelf Checkpoint* as a baseline to observe attention score evolution in different settings. This analysis identifies the key commonalities in high-quality attention scores.

**Commonality1. Higher attention to answer tokens in higher-quality models.** In all training settings, tokens closer to answer tokens (i.e., from a similarity higher than 90<sup>th</sup> percentile to a similarity higher than 95<sup>th</sup> percentile) receive increasingly higher attention scores. It can be observed that for both two measure metrics, the *Off-the-shelf Distillation Training* results are lower compared to the *Off-the-shelf Checkpoint*, while *Fine-tuned Distillation Training* shows improvement in both Step1 and Step2. The results suggest that in *Off-the-shelf Distillation*, the reader’s attention does not effectively "highlight" key information, leading to suboptimal training. In contrast, *Fine-tuned Distillation* after Step1 and Step2 both indicate that high-quality readers focus more on relevant answer

tokens, thereby enhancing both the retriever’s performance and the relevance of attention allocated to these tokens, which is also revealed in Figure 3.

**Commonality 2. Tokens similar to question nouns receive more attention in high-quality models.** Table 2 also indicates that tokens closer to the nouns in question tokens receive higher attention scores. The *Fine-tuned Distillation* experiments exhibit much higher values in both metrics compared to *Off-the-shelf Checkpoint* and *Off-the-shelf Attention Distillation*, aligning with their superior performance. However, unlike Commonality 1, the Spearman correlation between attention to question-related tokens and model performance isn’t consistent: while *Fine-tuned Attention Distillation* Step2 surpasses Step1, its metric values do not consistently align with this improvement, suggesting a more complex relationship.

**RQ3: How do we evaluate the quality of attention distillation on decoder-only readers based on the analysis results?**

**Indicator1.** Focusing on the attention scores of the nearest tokens to answer  $A$ , denoted as  $M_A = \{ma_1, \dots, ma_k\}$ . Higher average  $P_{ATTN}(ma_i)$  values indicate better attention distillation quality. Additionally, a higher average Spearman correlation between the  $P_{ATTN}(ma_i)$  and their semantic similarity to  $A$  also signifies better quality.

**Indicator2.** Examining the attention scores of tokens closest to nouns in question  $Q$ , denoted as  $M_Q = \{mq_1, \dots, mq_k\}$ . An increase in average  $P_{ATTN}(mq_i)$  suggests better quality. Moreover, if the average Spearman correlation between the attention scores of  $M_Q$  and their similarity to  $Q$  is above the threshold for a weak monotonic relationship (i.e., value > 0.3), the attention distillation quality is considered good.

**RQ4: Can we extend the proposed indicators to encoder-to-decoder structure readers?**

<sup>1</sup>We only focus on the nouns in the question in selecting *question-related* tokens.

<sup>2</sup>The highest values in the table are highlighted in bold on the NQ Dataset and underlined on the TriviaQA Dataset.



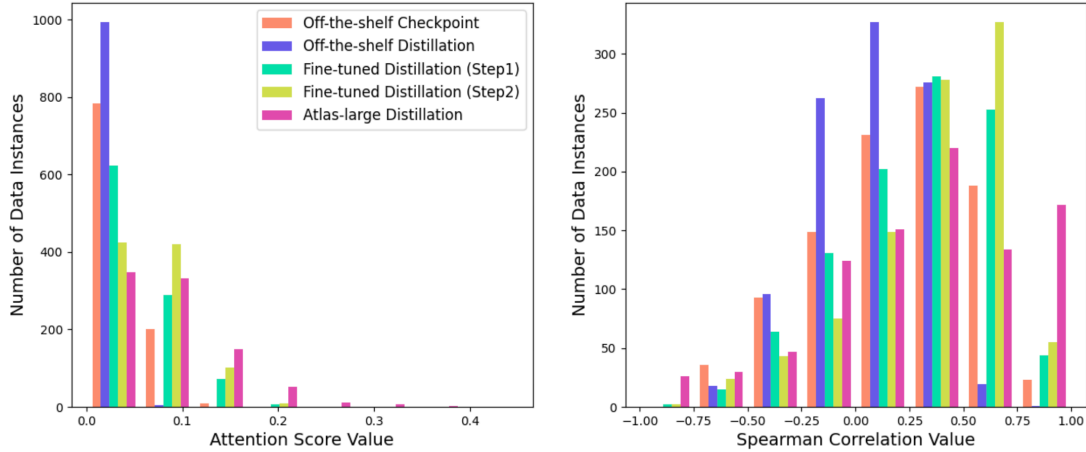


Figure 3: The attention score distribution histogram (left) and Spearman correlation distribution histogram of 95<sup>th</sup> percentile *answer-related* tokens under NQ dataset.

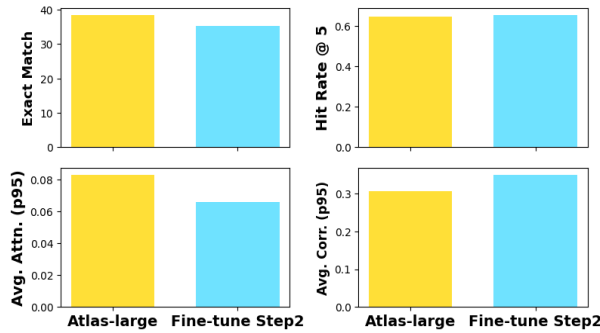


Figure 4: Model performance (top) and their attention distillation analysis (bottom) of *Atlas-large* model (yellow) for the *answer-related* tokens, comparing with *Fine-tuned Distillation Training (Step2)* (blue).

An analysis with the fine-tuned encoder-to-decoder structure *Atlas-large* model is presented in Figure 4. The results show that the performance of *Atlas-large* surpasses *Fine-tuned Distillation Training (Step2)*. However, only the average  $P_{ATTN}(ma_i)$  trend from Indicator1 applies to this encoder-to-decoder structure model, while *Atlas-large* exhibits a polarized distribution for the Spearman correlation values, as shown in Figure 3 and Appendix A.

**RQ5: Can we extend the proposed indicators to perplexity distillation training?**

Finally, we want to determine if our indicators can apply to perplexity distillation, another popular knowledge distillation method used in training the retriever model. We fine-tune *Atlas-large* model with the perplexity distillation method and find that the perplexity distribution does not align with either Commonality 1 or Commonality 2, saying that our indicators are not suitable for perplexity distillation (details in Appendix A and B).

## 4 Related Work

The concept of using attention scores for knowledge distillation was introduced by (Izcard and Grave, 2020a), and the following research has mainly focused on independently optimizing the reader and the retriever. Previous studies have explored improving large language model performance within the retriever-then-read framework by addressing issues like hallucination (Shuster et al., 2021) and dependency on pre-training data (Kandpal et al., 2023), or enhancing retriever efficiency through techniques like specific data sampling (Hofstätter et al., 2021). Only one study has examined the reader-retriever interaction within a neural-retrieval-in-the-loop architecture, noting that imperfect retrievers can degrade reader performance, though it lacked quantitative analysis (BehnamGhader et al., 2022).

Our study offers a comprehensive quantitative analysis of how the reader and the retriever interact during the neural-retrieval-in-the-loop training architecture under the attention distillation mechanism. We introduce novel metrics to evaluate the efficacy of the training process across all general reader-retriever framework.

## 5 Conclusion

In this paper, we comprehensively evaluate attention distillation for training retrieval-augmented language models, emphasizing the importance of attention to answer and question-related tokens. We further introduce novel metrics for assessing language models’ attention distillation ability to optimize the training process.

## Acknowledgement

This work is partially supported by NSF through grants IIS-1763365 and IIS-2106972. We thank the anonymous reviewers for their helpful feedback.

## Limitation

This paper analyzes the attention score-based knowledge distillation quality in training retrieval-augmented language models under various experimental settings in QA tasks. Furthermore, based on our findings, we have developed two indicators to assess the quality of attention score supervision. However, our exploration is conducted based on lightweight language models (i.e., language models with about one billion parameters) due to their flexibility and have yet to extend to larger-scale language models. In the future work, we will extend the study to larger-scale language models, focusing on validating the accuracy of our analysis on them to enhance the generalizability and applicability of our proposed methods.

## References

- Parishad BehnamGhader, Santiago Miret, and Siva Reddy. 2022. Can retriever-augmented language models reason? the blame game between the retriever and the language model. *arXiv preprint arXiv:2212.09146*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.
- Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently teaching an effective dense retriever with balanced topic aware sampling. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 113–122.
- Linmei Hu, Zeyi Liu, Ziwang Zhao, Lei Hou, Liqiang Nie, and Juanzi Li. 2023. A survey of knowledge enhanced pre-trained language models. *IEEE Transactions on Knowledge and Data Engineering*.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022a. [Unsupervised dense information retrieval with contrastive learning](#).
- Gautier Izacard and Edouard Grave. 2020a. Distilling knowledge from reader to retriever for question answering. *arXiv preprint arXiv:2012.04584*.
- Gautier Izacard and Edouard Grave. 2020b. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022b. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*.
- Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Active retrieval augmented generation. *arXiv preprint arXiv:2305.06983*.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, pages 15696–15707. PMLR.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2019. Generalization through memorization: Nearest neighbor language models. *arXiv preprint arXiv:1911.00172*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*.
- Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz,

- Eneko Agirre, Ilana Heintz, and Dan Roth. 2023. Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Computing Surveys*, 56(2):1–40.
- Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. 2023. [Scalable extraction of training data from \(production\) language models](#).
- OpenAI. 2023. [Gpt-4 technical report](#).
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. [The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only](#). *arXiv preprint arXiv:2306.01116*.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *arXiv preprint arXiv:2302.00083*.
- Weijia Shi, Julian Michael, Suchin Gururangan, and Luke Zettlemoyer. 2022. knn-prompt: Nearest neighbor zero-shot inference, 2022b. URL <https://arxiv.org/abs/2205.13792>.
- 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. *arXiv preprint arXiv:2301.12652*.
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. Retrieval augmentation reduces hallucination in conversation. *arXiv preprint arXiv:2104.07567*.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2023. A survey on large language model based autonomous agents. *arXiv preprint arXiv:2308.11432*.
- Haopeng Zhang, Xiao Liu, and Jiawei Zhang. 2023a. [Extractive summarization via ChatGPT for faithful summary generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3270–3278, Singapore. Association for Computational Linguistics.
- Haopeng Zhang, Xiao Liu, and Jiawei Zhang. 2023b. [SummIt: Iterative text summarization via ChatGPT](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10644–10657, Singapore. Association for Computational Linguistics.
- Haopeng Zhang, Semih Yavuz, Wojciech Kryscinski, Kazuma Hashimoto, and Yingbo Zhou. 2022. [Improving the faithfulness of abstractive summarization via entity coverage control](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 528–535, Seattle, United States. Association for Computational Linguistics.
- Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A Smith. 2023c. How language model hallucinations can snowball. *arXiv preprint arXiv:2305.13534*.

## A Quantitative Analysis of Answer-Related Tokens

We present a detailed analysis of *answer-related* tokens’ attention score distribution (or perplexity distribution of *Perplexity Distillation Training*) shown in Table 3. In addition to the histograms visualized for the attention score distribution of 95<sup>th</sup> percentile *answer-related* tokens under the NQ dataset (i.e., Figure 3), we also plot the corresponding attention score distribution figure under the TriviaQA dataset, as shown in Figure 5.

## B Quantitative Analysis of Question-Related Tokens

We present a detailed analysis of *question-related* tokens’ attention score distribution (or perplexity distribution of *Perplexity Distillation Training*) shown in Table 4, Figure 6, and Figure 7.

## C The Internal Relationship within Indicators

We also find a strong positive correlation (i.e., Pearson correlation  $> 0.5$  in most cases) between our two proposed indicators in decoder-only language model structure. In contrast, this correlation does not appear in encoder-to-decoder language model structure, which also indicates that Indicator2 is inapplicable to this language model structure.

## D Dataset Statistics

For the *NaturalQuestions* dataset, we split it according to the number of 79,168/8,757/3,610 to form the train/validation/test dataset; for the *TriviaQA* dataset, we split it according to the number of 78,785/8,837/11,313 to form the train/validation/test dataset.

## E Implementation Details

We conducted all computations on a Nvidia A100 GPU. For the *Off-the-shelf Distillation Training* and the *Fine-tuned Distillation Training*, we use *Falcon-1b* as the initial reader model and *Contriever* as the initial retriever model, which has about 1 billion and 110 million training parameters respectively. For the *Atlas-large Distillation Training* and *Perplexity Distillation Training*, we use *T5-large* as the initial reader model and *Contriever* as the initial retriever model, which has about 770 million and 110 million training parameters respectively.

**Off-the-shelf Distillation Training** We set the batch size to 1, the maximum length of the input prompt to 128, and limit the generation max length to 32. We set the learning rate to 1e-5 and used the Adam optimizer. For *NaturalQuestions* dataset, we set the total training steps to 160,000 with approximately 2000 warmup steps, training for about 40 hours. For *TriviaQA* dataset, we set the total training steps to 320,000 with approximately 4000 warmup steps, training for about 60 hours.

**Fine-tuned Distillation Training** For Step 1, we set the batch size to 1, the maximum length of the input prompt to 128, and limit the generation max length to 32. We set the learning rate to 1e-5 and used the Adam optimizer. For *NaturalQuestions* dataset, we set the total training steps to 160,000 with approximately 2000 warmup steps, training for about 30 hours. For *TriviaQA* dataset, we set the total training steps to 320,000 with approximately 4000 warmup steps, training for about 45 hours.

For Step 2, we set the batch size to 1, the maximum length of the input prompt to 128, and limit the generation max length to 32. We set the learning rate to 5e-7 and used the Adam optimizer. For *NaturalQuestions* dataset, we set the total training steps to 6,000 with approximately 300 warmup steps, training for about 2 hours. For *TriviaQA* dataset, we set the total training steps to 32,000 with approximately 600 warmup steps, training for about 3 hours.

**Atlas-large Distillation Training** We set the batch size to 1, the maximum length of the input prompt to 128, and limit the generation max length to 32. We set the learning rate to 4e-5 and used the Adam optimizer. For *NaturalQuestions* dataset, we set the total training steps to 10,000 with approximately 500 warmup steps, training for about 20 hours. For *TriviaQA* dataset, we set the total training steps to 30,000 with approximately 600 warmup steps, training for about 40 hours.

**Perplexity Distillation Training** We set the batch size to 1, the maximum length of the input prompt to 128, and limit the generation max length to 32. We set the learning rate to 4e-5 and used the Adam optimizer. For *NaturalQuestions* dataset, we set the total training steps to 20,000 with approximately 1000 warmup steps, training for about 40 hours. For *TriviaQA* dataset, we set the total training steps to 10,000 with approximately 500 warmup steps, training for about 15 hours.



Table 3: Mean and std. of attention scores (or perplexity distribution in *Perplexity Distillation Training*) and the Spearman correlations of the answer-related tokens.

Experiment	Dataset	Avg. Attn. (p90)	Spearman Corr. (p90)	Avg. Attn. (p95)	Spearman Corr. (p95)
Off-the-shelf Model Checkpoint	NQ	0.033 ± 0.016	0.227 ± 0.259	0.039 ± 0.023	0.196 ± 0.349
	TriviaQA	0.027 ± 0.013	0.218 ± 0.252	0.032 ± 0.019	0.206 ± 0.331
Off-the-shelf Attention Distillation	NQ	0.017 ± 0.008	0.145 ± 0.193	0.017 ± 0.010	0.076 ± 0.254
	TriviaQA	0.031 ± 0.012	0.160 ± 0.174	0.035 ± 0.017	0.172 ± 0.236
Fine-tuned Distillation Training (Step1)	NQ	0.039 ± 0.023	0.308 ± 0.276	0.052 ± 0.036	0.282 ± 0.336
	TriviaQA	0.058 ± 0.031	0.259 ± 0.261	0.074 ± 0.050	0.258 ± 0.331
Fine-tuned Distillation Training (Step2)	NQ	0.049 ± 0.023	0.316 ± 0.280	0.066 ± 0.036	0.350 ± 0.336
	TriviaQA	0.069 ± 0.036	0.290 ± 0.267	0.089 ± 0.061	0.320 ± 0.323
Atlas-large Distillation Training	NQ	0.062 ± 0.036	0.171 ± 0.462	0.083 ± 0.058	0.307 ± 0.471
	TriviaQA	0.072 ± 0.045	0.141 ± 0.379	0.091 ± 0.067	0.217 ± 0.438
Perplexity Distillation Training	TriviaQA	0.072 ± 0.039	0.029 ± 0.142	0.071 ± 0.042	0.013 ± 0.202

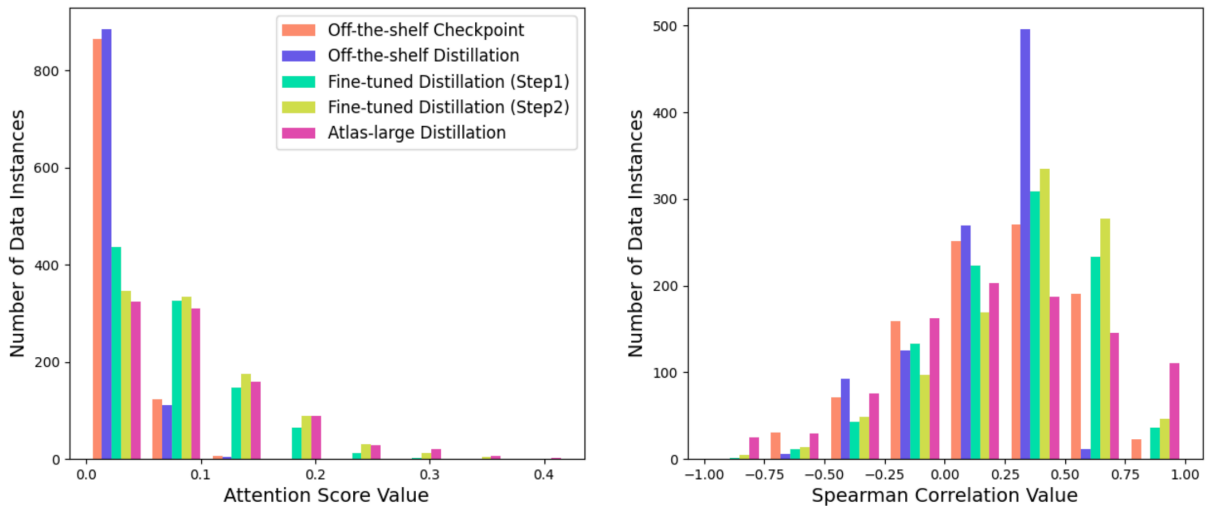


Figure 5: The attention score distribution histogram (left) and Spearman correlation distribution histogram of 95<sup>th</sup> percentile *answer-related* tokens under TriviaQA dataset.

Table 4: Mean and std. of average attention scores (or perplexity distribution in *Perplexity Distillation Training*) and Spearman correlations of the question-related tokens

Experiment	Dataset	Avg. Attn. (p90)	Spearman Corr. (p90)	Avg. Attn. (p95)	Spearman Corr. (p95)
Off-the-shelf Model Checkpoint	NQ	0.023 ± 0.011	0.103 ± 0.253	0.024 ± 0.014	0.092 ± 0.309
	TriviaQA	0.021 ± 0.010	0.103 ± 0.245	0.023 ± 0.013	0.067 ± 0.304
Off-the-shelf Attention Distillation	NQ	0.027 ± 0.010	0.139 ± 0.237	0.039 ± 0.017	0.153 ± 0.341
	TriviaQA	0.047 ± 0.016	0.144 ± 0.220	0.063 ± 0.025	0.260 ± 0.280
Fine-tuned Distillation Training (Step1)	NQ	0.035 ± 0.015	0.343 ± 0.238	0.045 ± 0.023	0.333 ± 0.303
	TriviaQA	0.058 ± 0.024	0.349 ± 0.222	0.078 ± 0.037	0.372 ± 0.285
Fine-tuned Distillation Training (Step2)	NQ	0.032 ± 0.014	0.310 ± 0.256	0.039 ± 0.021	0.225 ± 0.340
	TriviaQA	0.060 ± 0.025	0.367 ± 0.227	0.078 ± 0.037	0.326 ± 0.311
Atlas-large Distillation Training	NQ	0.037 ± 0.027	0.082 ± 0.251	0.038 ± 0.032	0.086 ± 0.345
	TriviaQA	0.047 ± 0.245	0.076 ± 0.249	0.050 ± 0.038	0.081 ± 0.348
Perplexity Distillation Training	TriviaQA	0.063 ± 0.038	-0.012 ± 0.207	0.060 ± 0.042	-0.036 ± 0.297

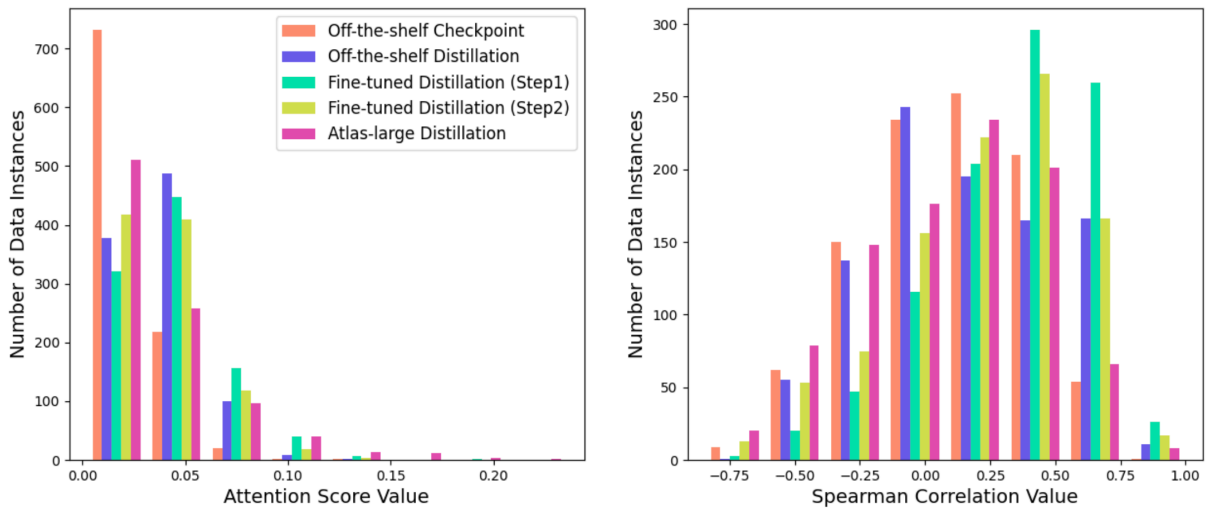


Figure 6: The attention score distribution histogram (left) and Spearman correlation distribution histogram of 95<sup>th</sup> percentile *question-related* tokens under NQ dataset.

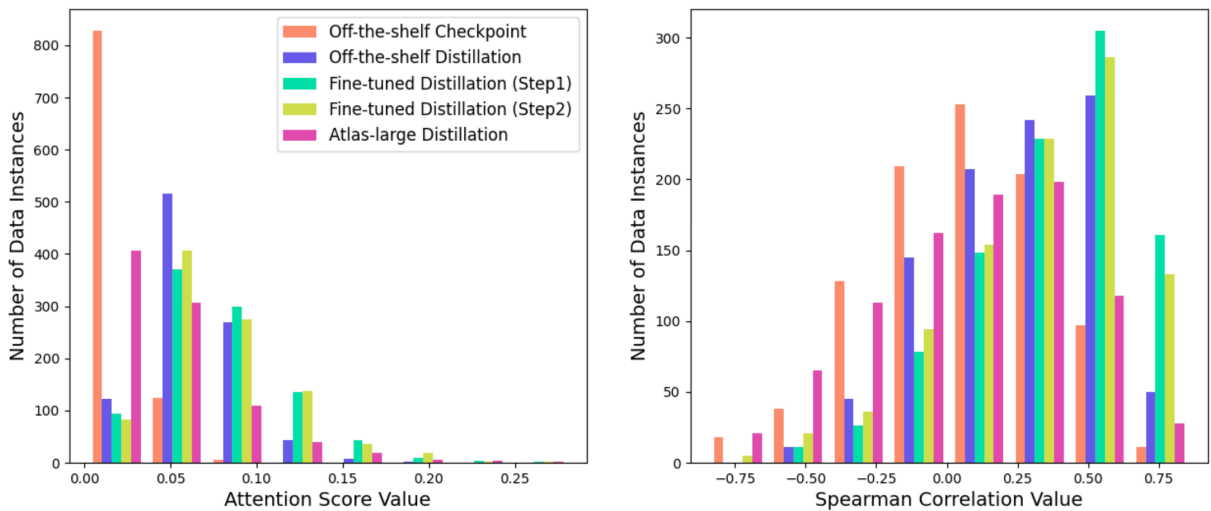


Figure 7: The attention score distribution histogram (left) and Spearman correlation distribution histogram of 95<sup>th</sup> percentile *question-related* tokens under the TriviaQA dataset.