

# Quantifying reliance on external information over parametric knowledge during Retrieval Augmented Generation (RAG) using mechanistic analysis

Anonymous ACL submission

## Abstract

Retrieval Augmented Generation (RAG) is a widely used approach for leveraging external context in several natural language applications such as question answering and information retrieval. Yet, the exact nature in which a Language Model (LM) leverages this non-parametric memory or retrieved context isn't clearly understood. This paper *mechanistically* examines the RAG pipeline to highlight that LMs demonstrate a "shortcut" effect and have a strong bias towards utilizing the retrieved context to answer questions, while relying minimally on model priors. We propose (a) Causal Mediation Analysis; for proving that parametric memory is minimally utilized when answering a question and (b) Attention Contributions and Knockouts for showing the last token residual stream do not get enriched from the subject token in the question, but gets enriched from tokens of RAG-context. We find this pronounced "shortcut" behaviour to be true across both LLMs (e.g., LLaMa) and SLMs (e.g., Phi).

## 1 Introduction

Retrieval Augmented Generation (RAG) (Lewis et al., 2021) is a popular method to enhance a Language Model's (LLM) capability to reason and execute tasks by leveraging additional context provided during inference time (Shao et al., 2023)(Singh et al., 2023)(IngestAI, 2023). Additionally, researchers have also explored shortcomings of RAG systems, such as inconsistent responses (Liu et al., 2023) and only (Wu et al., 2024) delved into the balance between a model's internal knowledge and externally retrieved information, examining their practical value.

Several research papers have proposed the approaches for editing knowledge in language model, including techniques like ROME (Meng et al., 2022a), MEMIT (Meng et al., 2022b) to update or correct facts. On the flip side, with the popularity of LLM integration for various tasks leveraging

proprietary, enterprise, and private data, the use of RAG framework has increased to tackle *hallucinations* while reasoning on **new never seen before** (out of distribution) knowledge tasks. However, a comprehensive study mechanistically probing of Language Model's behavior of choosing between information from RAG-generated context over intrinsic parametric knowledge has not been conducted to the best of our knowledge.

## 2 Probing Methods

To mechanistically interpret the knowledge contributions towards factual reasoning by LLMs and SLMs, we use three methods for causal mediation, described as follows: **Causal Tracing** (Meng et al., 2022a) identifies specific hidden states that significantly influence factual predictions. The approach involves a clean run, corrupted run and a corrupted-with-restoration run. The corrupted run involves corrupting a certain span of the text, and running the forward pass of the model. In the restoration run, activations from the clean run are patched one by one into the corrupted run, and the increase in answer probability is observed; the most crucial activations are thus causally determined. The causal importance of a certain activation is quantified using Indirect Effect, which is defined as the difference between the corrupted run and the corrupted-with-restoration run probabilities:  $IE(h_i^{(l)}) = P_{clean}^*(h_i^{(l)})[y] - P^*[y]$ . The Average Indirect Effect of a hidden node is an average of IE over all the prompts in the dataset.

The **Attention Contribution** (Yuksekgonul et al., 2024), focuses on the role of attention mechanisms in shaping the output of language models. This approach investigates how attention weights, particularly from the subject token in a query to the last token position, contribute to the model's predictions. By examining the norm of these attention weights  $\|a_{i,T}^{(\ell)}\|$ , we observe what tokens

081 the last token pays the most attention to, during  
082 the generation process. The **Attention Knockout**  
083 mechanism (Geva et al., 2023) identifies critical at-  
084 tention edges in transformer-based models that are  
085 essential for maintaining prediction quality. The  
086 process involves identifying critical edges whose  
087 removal significantly degrades the model’s predic-  
088 tion quality, by means of setting the attention from  
089 position  $i$  to  $j$  in the attention matrix to  $-\infty$ .

### 090 3 Data and Models

091 For examining model activations for causal tracing,  
092 patching and inspecting AIE, while systematically  
093 analyzing attention contributions we choose open  
094 source LMs like LLaMa-2 (7B) and Phi-2 (2.7B)  
095 models. And for understanding the behavior in the  
096 non-RAG setting, we leverage the *Knowns 1000*  
097 *dataset*, a dataset of 1209 prompts (Meng et al.,  
098 2022a). For the RAG setting, we augment the  
099 *Knowns 1000 dataset* with added context gener-  
100 ated synthetically using GPT-4. We use GPT-4  
101 generated context to control the length of each seg-  
102 ment within the RAG-context and also the presence  
103 of *attribute* or *object*.

### 104 4 Results

105 Experimenting with LLaMa and Phi-2 family of  
106 models on 1209 samples from the knowns fact  
107 dataset for vanilla-case and RAG-scenario, demon-  
108 strate that both models exhibit a strong bias towards  
109 utilizing external knowledge provided by RAG.

110 Utilizing *Causal Tracing* method and measuring  
111 Average Indirect Effect (AIE) at different positions  
112 of the prompt, such as Last Subject Token (LST),  
113 Last Token (LT), it is found that for the vanilla-case  
114 (non-RAG) LST had high AIE, but it substantially  
115 lowered when RAG-generated context was added.  
116 As concluded in (Meng et al., 2022a), LST has  
117 the largest influence from model priors and low-  
118 ering AIE of LST demonstrates reduced influence  
119 of parametric memory. We specifically observe  
120  $\sim 10X$  decrease in AIE of LST for LLaMA-2 and  
121  $\sim 35X$  decrease in AIE of LST for Phi-2, when  
122 RAG-generated context is added.

123 This finding was further corroborated by utiliz-  
124 ing two other probing methods - Attention Knock-  
125 outs and Attention Contributions. The LT is a cru-  
126 cial component in the LLM decoding process, as it  
127 is projected onto the vocabulary during decoding  
128 time. Thus any information that has to be decoded,  
129 will be propagated by the MLP and attention layers

130 to the LT residual stream. We measured the atten-  
131 tion contributions from the Subject Token (ST) to  
132 the LT and observe a substantial decrease in ST  
133 contributions for the RAG-scenario as compared to  
134 the vanilla non-RAG case where no external con-  
135 text is provided. For LLaMa-2, the mean attention  
136 contribution decreased by  $\sim 1.6X$  for RAG case, in  
137 comparison to non-RAG vanilla case, and for Phi-2  
138 a reduction of  $7x$  was observed for ST contribu-  
139 tion. Conversely, the *answer token* contribution for  
140 RAG setting, increases significantly for LLaMa-2  
141 and Phi-2 in comparison to ST contribution in the  
142 RAG setting. This further confirms our hypothe-  
143 sis of the LLM being less reliant on its parametric  
144 memory and exhibiting a "shortcut" behavior.

145 Using Attention knockouts (Geva et al., 2023)  
146 approach, it is observed that "knocking out" atten-  
147 tion from ST to the LT, reduces the probability  
148 of the answer in the LM’s last token predictions  
149 by 20% in LLaMa-2 and 25% in Phi-2. This is in  
150 sharp contrast to the RAG setting, where knocking  
151 off attention at ST positions leads to  $<5\%$  drop in  
152 the answer probabilities. This finding further rein-  
153 forces the finding that the model takes a "shortcut"  
154 while relying minimally on its parametric memory.

### 155 5 Conclusions and Future Work

156 Using Causal Tracing (Meng et al., 2022a) in over  
157 1200 samples of the *known facts dataset*, in RAG-  
158 scenario for LLaMa-2 and Phi-2, we observe a  
159 reduced AIE on the last subject token, and poten-  
160 tially reduced dependence on parametric memory.  
161 This is further corroborated by our experiments  
162 with attention contributions and attention knock-  
163 outs. Using three mechanistic probing techniques,  
164 we observe 1) reduced reliance on parametric mem-  
165 ory 2) reduced information flow from the subject  
166 token to the last token residual stream 3) a shortcut  
167 behavior where information from the attribute to-  
168 ken flows to the last token residual stream during  
169 factual predictions in the RAG setting.

170 Future work will address the extension to larger  
171 LMs ( $> 13B$  parameters). We also plan to study the  
172 impact of LM behavior in longer context, and in  
173 settings where language models are known to ex-  
174 hibit primacy and recency bias (Liu et al., 2023) in  
175 a future work. Additionally, we aim to replicate our  
176 findings using a conventional RAG pipeline to au-  
177 tomatically create context rather than synthetically  
178 generating it using GPT4.

## References

- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. [Dissecting recall of factual associations in auto-regressive language models](#). *Preprint*, arXiv:2304.14767.
- IngestAI. 2023. [Retrieval-augmented generation \(rag\): Enhancing llms with external knowledge](#).
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#). *Preprint*, arXiv:2005.11401.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. [Lost in the middle: How language models use long contexts](#). ArXiv:2307.03172.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. [Locating and editing factual associations in gpt](#). *Advances in Neural Information Processing Systems*, 35:17359–17372.
- Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2022b. [Mass-editing memory in a transformer](#). *arXiv preprint arXiv:2210.07229*.
- C. Shao, T. Kim, and Z. Gao. 2023. [Eragent: Enhancing retrieval-augmented language models with improved accuracy, efficiency, and personalization](#). *arXiv preprint arXiv:2405.06683*.
- A. Singh, M. Sachan, and K. Guu. 2023. [Improving the domain adaptation of retrieval augmented generation \(rag\) models for open domain question answering](#). *Transactions of the Association for Computational Linguistics*.
- Kevin Wu, Eric Wu, and James Zou. 2024. [How faithful are rag models? quantifying the tug-of-war between rag and llms’ internal prior](#). *Preprint*, arXiv:2404.10198.
- Mert Yuksekgonul, Varun Chandrasekaran, Erik Jones, Suriya Gunasekar, Ranjita Naik, Hamid Palangi, Ece Kamar, and Besmira Nushi. 2024. [Attention satisfies: A constraint-satisfaction lens on factual errors of language models](#). *Preprint*, arXiv:2309.15098.