Deciphering the Interplay of Parametric and Non-parametric Memory in Retrieval-augmented Language Models

Mehrdad Farahani, Richard Johansson Chalmers University of Technology and University of Gothenburg





Output (O)	P_O
Stockholm	0.84
Milan	0.95
Stockholm	0.51
Milan	0.98
Stockholm	0.67
Stockholm	0.64
	Output (O) Stockholm Milan Stockholm Stockholm





RAG Models

- Combine parametric memory (model's weights) and nonparametric memory (external knowledge)
- Support specialized knowledge tasks
- Improve factual accuracy through data retrieval

Atlas the candidate mode!

- Trains jointly on the Language Model (LM) and retriever
- Balances parametric and non-parametric knowledge efficiently

How well does Atlas balance between what it already **knows** and what it **fetches**? To understand <u>when</u> and <u>why</u> Atlas favors one source of knowledge over the other?

RQ¹: Which aspect of the model representation impacts the output in copying mode? **RQ²: What specific parts of the model trigger copying?**

"What is the capital of Sweden?" (s: Sweden, o: Stockholm, r: capital of)







 $TE = Y(X \leftarrow 1) - Y(X \leftarrow 0)$ $IE = Y(X \leftarrow 0, M(X \leftarrow 1)) - Y(X \leftarrow 0)$ *P*(counterfactul|context) $Y = log \frac{1}{P(\text{true answer}|\text{context})}$

Experiments

Experiment 1: What is the balance between parametric and non-parametric behavior?

- a. Modifying object representations with counterfactual
- b. Keeping context unchanged $(X \leftarrow 0)$

Experiment 2: What makes the model decide to rely on the context?

- a. Replacing **object tokens** with **counterfactuals** in the context
- b. Adding **noise** to **subject/relation** tokens
- c. Affecting subject/relation tokens with noise $(X \leftarrow 0)$

Experimental Design

- Utilizing two datasets: PopQA and PrincetonEntityQuestion (PEQ)
- Focusing only on the LM of Atlas
- Applying synthetic context
- Retaining queries with correct answers, with and without synthetic context
- Ensuring the true answer is a substring of the actual answer
- Removing relations with few data points

Object

Using Y as the log ratio stabilizes values and highlights subtle shifts

Relation

Subject

Query: What is the capital of Sweden? Stockholm is the capital of Sweden. Context:

Relations	Query Template	Context Template	#
capital	What is the capital of [subj]?	The capital of [subj] is [obj].	101
capital of	What is [subj] the capital of ?	[subj] is the capital of [obj].	26
color	What color is [subj] ?	The color of [subj] is [obj].	4
composer	Who was the composer of [subj]?	[obj] was the composer of the musical work [subj].	4
country	In what country is [subj] ?	The [subj] is located in [obj].	101

Sample queries with synthetic context templates



AIE and PSE Results

Experiment 1:

Experiment 2:



- Highlights the impact of object tokens in copying mode
- The model performs a form of relevance evaluation
- MLP contributes to translating representations from the encoder to the decoder.
- Severing MLP reduces reliance on object tokens
- The lower impact of Attention suggests it is less involved in this process
- MLP shows a similar effect for relation tokens, indicating that both object and relation tokens undergo similar representation transformations for the decoder.

ATE Results

Reveal significant differences, with greater variability observed in the non-parametric subset. Comparatively, the overall behavior of the model aligns more closely with the non-parametric subset



Behavior

relations







-50

-100

- In the early layers, the MLP indicates that the model focuses on subject and relation tokens
- As processing continues, the focus gradually shifts toward object tokens
- Attention focuses on the entire context and maintains coherence in this process
- MLP's role expands to help with the object extraction step
- MLP and Attention are key in transitioning from relevance evaluation to object extraction
- Both MLP and Attention work closely with subject and relations tokens in the early layers
- Changes in object tokens reveal that MLP must collaborate with Attention to extract object tokens accurately

- Exploring copying vs. recalling decisions through two experiments
- **Copying** from context by **assessing relevance** in the Atlas model
- Contextualizing relevance and translating information with early-to-middle MLP
- Integrating context and extracting object tokens with later Attention



Link to the code and paper

Dataset Specificity: Using PopQA and PEQ datasets, limiting generalizability due to differing memory behaviors

subjects

- **Context Manipulation:** Employing counterfactuals, which may not capture noisy or ambiguous contexts fully
- Model Generalization: Evaluating Atlas's adaptability to other RAG models, especially in varied contexts
- **Temporal Relevance:** Balancing parametric and nonparametric memory during temporal changes remains a challenge

