

---

# Enhancing Retrieval in QA Systems with Derived Feature Association

---

Abhishek Goyal   Keyush Shah   Isaac Wasserman

Department of Computer and Information Science

University of Pennsylvania

Philadelphia, PA 19104

{abhi2358, keyush06, isaacrw}@seas.upenn.edu

## Abstract

Retrieval augmented generation (RAG) has become the standard in long-context question answering (QA) systems. However, typical implementations of RAG rely on a rather naive retrieval mechanism, in which texts whose embeddings are most similar to that of the query are deemed most relevant. This has consequences in subjective QA tasks, where the most relevant text may not directly contain the answer. In this work, we propose a novel extension to RAG systems, which we call **Retrieval from AI Derived Documents (RAIDD)**. RAIDD leverages the full power of the LLM in the retrieval process by deriving inferred features, such as summaries and example questions, from the documents at ingest. We demonstrate that this approach significantly improves the performance of RAG systems on long-context QA tasks.

## 1 Introduction

### 1.1 Preliminaries

First introduced by [6], retrieval augmented generation (RAG) allows LLMs with limited context windows to leverage a large corpus of documents during generation [13]. RAG extends LLMs with a retrieval mechanism that takes a query, selects the most relevant texts from a given corpus, and hands them to the generator to inform its answer. Early approaches, were optimized end-to-end, using a jointly learned retriever and generator that communicated through a shared embedding space [6]. However, the requirement that such a system must be trained from scratch for each choice of generator architecture makes this approach expensive and cumbersome given the rapid pace with which new LLMs are developed. However, this paradigm was subverted by [11], which assumes the generator to be black-box, training a generator-agnostic retriever that simply prepends the retrieved text to the generator’s input. In practice, the retriever is often further simplified to score documents based on their cosine similarity in a pretrained embedding space [9]; also popular is the use of BM25, a simple term-frequency based similarity metric [13].

### 1.2 Motivation

RAG systems, especially those which rely on embedding cosine similarity or BM25 to measure relevance, are fast and remarkably effective for answering questions whose answers are explicitly stated in the text. However, from a user’s perspective, this is only marginally more effective than a simple **ctrl+f** search, we expect more from the systems that we call “artificially intelligent”. In particular, we expect them to be able to answer questions whose answers are not explicitly stated in the text, but can be easily inferred from the text. Consider

<b>Question</b>	All of historians speak highly of Picardo’s work, is this true? Why?
<b>Target Text</b>	“...was somewhat frowned upon in the 1960s and 1970s, and over half a century later is seen by archeologists and historians as a matter of significant <b>controversy and regret</b> ...”
<b>Ground Truth</b>	False, because some people believe that Parrado destroyed the part of historical and architectural.
<b>Retrieved Text</b>	“...Picardo’s published architectural drawings were highly regarded. They were described as <b>“magnificent”</b> by the leading Spanish restoration architect ...”
<b>Prediction</b> ✗	Yes, because his architectural drawings were described as “magnificent” ...

Figure 1: Example of a question from the LooGLE [7] dataset answered by a GPT-4 based RAG system. The system identifies the a text that describes how Picardo’s work was regarded by one figure, but it fails to identify the more subtly worded target text which contains the answer.

the example in Figure 1, using cosine similarity between the query and text; the retriever latches onto the text which most explicitly describes commentary on the artist’s work and ignores the text which contains the answer but does not contain words like “regarded”. This is a common failure mode for RAG systems, and it demonstrates how the retriever can be a hindrance to what is otherwise a powerful model for natural language understanding.

### 1.3 Related Work

**TODO:** Write about other attempts to improve dense retrieval.

Zhao et al. (2024) [13] provides a taxonomy of RAG proposed enhancements that acknowledges “data augmentation” as a minor category of prior work, referencing MakeAnAudio [4] which augments an audio retrieval database with synthesized captions, however this is simply a bridge between audio and powerful language models. Our work is more closely related to [1], which demonstrates the utility of transforming RAG documents into more digestible, concise, and explicit forms. They propose “propositional retrieval”, a method quickly adopted by the AI agent community [3] due to its effectiveness, portability, and runtime efficiency. Their method preprocesses documents to extract individual propositions from the text and indexes them instead of chunks of the original text. This greatly improves retrieval of information which is explicitly stated, but it sacrifices nuance that may be necessary for answering more complex questions.

### 1.4 Contribution

**TODO:** This has to be explained in detail: how ours contribute to the previously stated results

## 2 Method

### 2.1 Derived Document Association

RAIDD generalizes the retrieval paradigm first introduced by Chen et al. (2023) [1], in which input documents are preprocessed to generate new documents whose utility is more transparent. In this original paradigm, these more transparent documents are placed into the question answering context. RAIDD diverges from Chen et al. (2023) in two main ways: (1) rather than simply dividing the original documents into their component propositions,

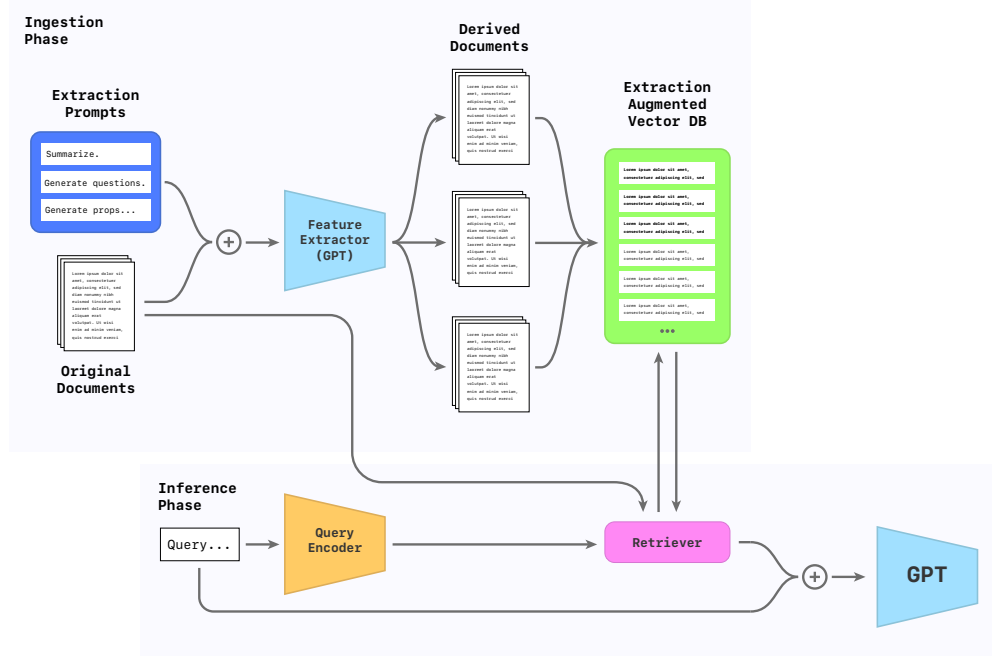


Figure 2: During the document ingest phase, RAIDD derives new documents from the input by prompting an GPT feature extractor to summarize, generate propositions, and generate questions from the original documents. At inference, the retriever identifies the most relevant derived documents and places the corresponding source documents into context from question answering.

we generate summaries and potential questions as our derived document, and (2) rather than placing the derived documents themselves into context, we use their source documents to minimize information loss. We call this practice of using the derived documents as handles for the original text “derived document association”.

More concretely, our method involves two phases: ingest and inference. During ingest, we prompt an LLM to generate derived documents from each input document. These derived documents are either summaries of each chunk, sets of questions from each chunk, or sets of propositions from each chunk. We generate embeddings from these derived documents and store them in a vector database. At inference time, the retriever matches our query against all of the derived documents. For each of the top  $k$  derived documents, we place the corresponding original document into our question answering context.

## 2.2 Implementation Details

Our experiments are implemented using the LlamaIndex [9] RAG library. For generating derived documents, we prompt an instruction tuned Mixtral-8x7B model [5]. We use the frontier-class Mistral Large model [12] as our question answering model. Our retriever uses a simple cosine similarity metric between the query and document embeddings, which are encoded using OpenAI’s text-embedding-ada-002 [2].

## 3 Experiments

### 3.1 RAIDD-S

In this version of our experiment, we generate summaries of the source documents leveraging the power of Language models. The summaries are generated for the source documents and converted into embeddings that are then stored in the Vector DB.

Note that there could be multiple chunks for a document depending upon the chunk size (we have experimented with different chunk sizes, results are tabulated further in this document). For each subsequent chunk, the previous chunk is provided as input to the model for providing more context. We require deep understanding of context that spans multiple sections or chunks of text, such as in document summarization or detailed question answering from long texts, hence providing previous chunks as context to subsequent chunks can be crucial. This helps the model maintain coherence and understand the document’s narrative or argumentative structure.

### 3.2 RAIDD-Q

TODO:

- Run RAIDD-Q
- Run with derived document ensemble

### 3.3 Configurations

chunkSize=64, chunkOverlap=10, topK=32  
chunkSize=128, chunkOverlap=25, topK=16  
chunkSize=256, chunkOverlap=50, topK=8  
chunkSize=512, chunkOverlap=100, topK=4  
chunkSize=1024, chunkOverlap=200, topK=2  
chunkSize=2048, chunkOverlap=400, topK=1

### 3.4 Dataset and Evaluation

We evaluate long-context question answering performance using the long-dependency QA subset of LooGLE [7]. To minimize cost, we use the first 100 questions of the dataset, which utilize a total of 14 input contexts. Question answering performance is measured using ROUGE [8] and GPT-4 [10] prompted to decide whether the generated answer was sufficiently similar to the ground-truth given the question.

Note that generating summarizations and questions used as the AI derived documents are just the means to an end that is answering the questions.

### 3.5 Improvements over Baseline

### 3.6 Ablation Studies

Our best model will likely use a combination of derived documents. For each feature used, remove it from the base model and evaluate performance difference.

### 3.7 Analysis

Here, we will perform analyses of success and failure modes. We will also look at how RAIDD improves the retrieval of target passages and changes the rank of individual documents. We will also look at a few example questions and how our method compares to the baseline.

Table 1: Performance comparison of various flavors of RAIDD. RAIDD-S uses summary generation, while RAIDD-Q uses question generation.

Method	Chunk Size	Chunk Overlap	TopK	GPT-4 Self-Score	ROUGE-1	ROUGE-L
RAG	64	10	32	0.43	0.216	0.174
	128	25	16	0.46	0.209	0.167
	256	50	8	<u>0.48</u>	<b>0.249</b>	<b>0.206</b>
	512	100	4	<u>0.48</u>	0.223	<u>0.189</u>
	1024	200	2	0.39	0.212	0.167
	2048	400	1	0.35	0.191	0.152
RAIDD-S	64	10	32	0.48	0.214	0.171
	128	25	16	0.41	0.204	0.167
	256	50	8	<b>0.49</b>	<u>0.230</u>	0.183
	512	100	4	0.43	0.224	0.178
	1024	200	2	0.31	0.205	0.178
	2048	400	1	0.31	0.201	0.163
RAIDD-Q	64	10	32			
	128	25	16			
	256	50	8			
	512	100	4			
	1024	200	2			
	2048	400	1			

## 4 Conclusion

## References

- [1] Tong Chen, Hongwei Wang, Sihao Chen, Wenhao Yu, Kaixin Ma, Xinran Zhao, Hongming Zhang, and Dong Yu. Dense X Retrieval: What Retrieval Granularity Should We Use? *arXiv*, 2023.
- [2] Ryan Greene, Ted Sanders, Lilian Weng, and Arvind Neelakantan. New and improved embedding model, Dec 2022.
- [3] Chase Harrison. Langchain – propositional-retrieval, 2024.
- [4] Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. Make-An-Audio: Text-To-Audio Generation with Prompt-Enhanced Diffusion Models. *arXiv*, 2023.
- [5] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th  ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mixtral of Experts. *arXiv*, 2024.
- [6] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich K  ttler, Mike Lewis, Wen-tau Yih, Tim Rockt  schel, Sebastian Riedel, and Douwe Kiela. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *arXiv*, 2020.
- [7] Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. Loogle: Can long-context language models understand long contexts? *arXiv preprint arXiv:2311.04939*, 2023.
- [8] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.

- [9] Jerry Liu. LlamaIndex, 11 2022.
- [10] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Ceron Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.
- [11] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. REPLUG: Retrieval-Augmented Black-Box Language Models. *arXiv*, 2023.
- [12] Mistral AI Team. Au large | mistral ai | frontier ai in your hands, Feb 2024.

- [13] Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang, and Bin Cui. Retrieval-augmented generation for ai-generated content: A survey, 2024.