

# Retrieval Augmented Generation (RAG) Model

Ankit Mishra  
Department of Computer Science and  
Engineering  
Sharda University  
Greater Noida, India  
ankit02112002@gmail.com

Aniket Gupta  
Department of Computer Science and  
Engineering  
Sharda University  
Greater Noida, India  
Aniketkg472@gmail.com

(Prof.)Dr. Anil Kumar Sagar  
Department of Computer Science and  
Engineering  
Sharda University  
Greater Noida, India  
anil.sagar@sharda.ac.in

**Abstract**—LLM (Large Language Model) is a type of Artificial Intelligence algorithm which uses deep learning techniques like neural networks models to generate text-based answers for a variety of user queries. The model gets trained over a large number of pre-defined dataset and generate results based on the facts emerging from this data. LLM RAG is a concept where the dataset is continuously fetched from the real time facts rather than pre-stored data. It helps in providing the users with most accurate and up-to-date information. RAG model reduces the time and cost of continuously training the LLM with the latest data and updating the parameters. It consists of 2 phases: retrieval phase where algorithm will search for and retrieve snippets related to user's prompt. The content generation phase where using the user's prompt the and content generated and uses the LLM model to generate the final text-based answer for the user's query. The LLM RAG (Retrieval Augmented Generation) Model represents a groundbreaking advancement in Natural Language Processing (NLP) that integrates both retrieval and generation capabilities to enhance text generation tasks. This abstract will provide a detailed overview of the LLM RAG Model, highlighting its architecture, key components, training methodology, and applications.

**Keywords**— *Large Language Models, Retrieval Augmented Generation, text generation, Angular, Flask, YouTube Data API, sentiment analysis.*

## I. INTRODUCTION

Large Language Models (LLMs) have revolutionized natural language processing (NLP) by generating coherent and contextually relevant text across various applications. However, their reliance solely on internal training data limits their ability to adapt to new information and understand evolving topics. To overcome this limitation, the LLM Retrieval Augmented Generation (RAG) model has emerged as a promising approach, integrating retrieval mechanisms to augment generation with external knowledge sources.

The LLM RAG model, introduced in 2020 by researchers at Google AI and DeepMind, represents a significant advancement in text generation technology. By retrieving relevant passages from external knowledge bases before generating text, RAG enables LLMs to go beyond their training data, enhancing factual accuracy and contextual understanding. This two-stage process not only improves the relevance and coherence of generated text but also expands the model's applicability to diverse domains. The responsible development and deployment of LLM RAG holds immense potential to push the boundaries of knowledge and creativity, ushering in a new era of possibilities.

Motivated by the potential of LLM RAG models, this research paper explores their capabilities and applications across various fields. We provide an overview of LLM RAG models, highlighting their theoretical foundations, architecture, and key features. Additionally, we present a practical implementation of LLM RAG for real-time analysis of YouTube comments, demonstrating its effectiveness in sentiment analysis, topic modeling, and user engagement enhancement. Through this study, we aim to contribute to the understanding and advancement of LLM RAG models, paving the way for their broader adoption in education, healthcare, customer service, scientific research, and beyond. Furthermore, we discuss future directions for LLM RAG models, emphasizing the importance of knowledge base integration, explainability, and ethical considerations in shaping the future of language technology.

## II. LITERATURE REVIEW

Large Language Models (LLMs) have emerged as powerful tools in natural language processing, yet their reliance solely on internal data can limit their accuracy and understanding of complex or evolving topics. To address this limitation, researchers have proposed the LLM Retrieval Augmented Generation (RAG) model, which integrates retrieval mechanisms to augment text generation with external knowledge sources.

### [1] Dense Passage Retrieval for Open-Domain Question Answering

The foundational work by Lee et al. (2019) introduced dense passage retrieval for open-domain question answering, providing the groundwork for the retrieval mechanism in LLM RAG models. Guu et al. (2020) further advanced this concept with Retrieval-Augmented Language Model Pre-Training, incorporating retrieval and generation tasks in the pre-training objective to enhance the model's factual grounding and contextual understanding.

### [2] Retrieval-Augmented Language Model Pre-Training

Attention mechanisms play a crucial role in enabling LLM RAG models to focus on relevant information during text generation. The Transformer architecture introduced by Vaswani et al. (2017) revolutionized NLP with its powerful attention mechanism, allowing models to efficiently capture long-range dependencies and relationships between words. Lin et al. (2018) expanded upon this concept by exploring self-attention for hierarchical sentence representation, improving the model's ability to understand context and meaning.

[3] **Transformer: Attention is All You Need**

Recent advancements in LLM RAG models have focused on improving performance and applicability in specific domains. Lewis et al. (2021) proposed novel training approaches to enhance RAG's performance in specific domains by utilizing domain-specific knowledge bases.

In summary, the literature on LLM RAG models demonstrates their potential to enhance factual accuracy, relevance, and coherence in generated text by integrating external knowledge sources. Further research in this area aims to improve performance, applicability, and reliability, paving the way for their broader adoption in various domains, including education, healthcare, customer service, and scientific research..

III. PROBLEM DEFINITION

The rapid advancement of large language models (LLMs) has opened the door to a multitude of new applications across various industries. Among these is the emerging area of retrieval-augmented generation (RAG), which aims to combine the vast capabilities of LLMs with targeted retrieval from specific data sources to provide more accurate and contextually relevant responses. Despite the potential, RAG systems face several significant challenges. One key challenge is the effective integration of retrieval and generation components. LLMs, trained on extensive data, are adept at producing natural language output, but may lack specific or up-to-date information on niche topics. Retrieval systems, on the other hand, can efficiently access and extract relevant information from diverse sources. Designing a system that effectively combines these two functionalities while ensuring seamless and coherent outputs is non-trivial.

IV. PROPOSED SOLUTION

This research project aims to provide a comprehensive overview of the LLM RAG model and explore its implementation. The LLM RAG model has ability to leverage external knowledge, overcoming limitations faced by traditional LLMs. This paper will delve into the model's architecture, strengths, and potential applications, followed by an attempt to implement a basic version of the model. Provide a detailed overview of the LLM RAG model. This includes covering its theoretical foundations, architecture, training process, and key features like retrieval and generation mechanisms.

V. METHODOLOGY

This research employs a systematic approach to explore the capabilities and potential applications of Large Language Model Retrieval-Augmented Generation (LLM RAG) models. The methodology is divided into several key stages:

1. Literature Review

- Conducted a comprehensive review of existing literature and research papers on LLM RAG models, including their theoretical foundations, architecture, and applications.
- Identified relevant studies on related topics such as large language models, retrieval-augmented NLP, and attention mechanisms.

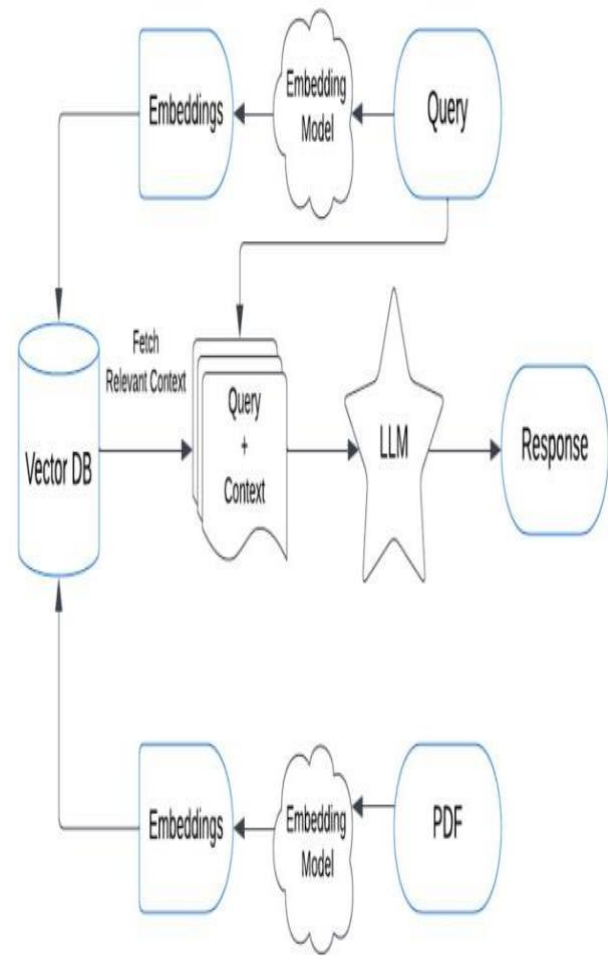


Figure 1: Working of LLM RAG

2. Model Analysis

- Analyzed the architecture and key components of the LLM RAG model, including retrieval mechanisms, attention mechanisms, and generation processes.
- Investigated the strengths and limitations of the LLM RAG model in addressing the challenges faced by traditional language models. 12-year period is taken into account in the dataset.

3. Implementation Process

- Selected appropriate programming languages and frameworks for implementing the LLM RAG model, including Python for backend development and Angular for frontend development.
- Developed retrieval and generation modules for the model, utilizing technologies such as Flask for the backend and Angular Material for frontend components.
- Integrated external knowledge sources such as the YouTube Data API for retrieving real-time comments and OpenAI API for text generation.

4. Model Training

- Acquired and pre-processed relevant datasets, including YouTube comments data obtained through the YouTube Data API.
- Utilized pre-trained models like GPT-3.

This dataset aims to collect YouTube comments in real-time using Google APIs. The dataset would be in the form of json. This data can be used for various purposes, including sentiment analysis, topic modeling, understanding user engagement, and summarizing the viewers opinion about the video. The real time data extracted will help them sort the comments based on their sentiments like top funny, positive and most liked comments.

### Data Sources:

YouTube Data API provides access to various YouTube data, including comments. We will use the `commentThreads.list` method to retrieve comments for specific video IDs.

### Data Collection Process:

**Search Term:** Given the user search query, it will display the top video ids based on the search term.

**API calls:** Use the `commentThreads.list` method of the YouTube Data API to retrieve comments for each video ID. For live streams, use the `liveChatMessages.list` method of the YouTube Live Streaming API.

**Pagination:** Using pagination to continuously retrieve top 500 videos as there is a limit to 50 videos per query.

## 7. Future Directions

id	video_id	kind	etag	id	snippet	publishedAt	title	channelTitle
0	youtubeSearchResult	60037tUKZ...5G9u1up1nagpG	(kind 'youtubeVideo' videoId 'WdUw9PIT...	(publishedAt '2013-12-27T06:44:51Z', title...	2013-12-27T06:44:51Z	Di'Arastane I'm Full (Video Song) Zingba Na M...	T-Series	
1	youtubeSearchResult	Nc07Z8p00Xf0z0235PwP1q1	(kind 'youtubeVideo' videoId '5Q3uqP...	(publishedAt '2013-12-17T09:09:24Z', title...	2013-12-17T09:45:14Z	Khabon Ke Fanday Full video song) Zingba	T-Series	
2	youtubeSearchResult	x40z7lmqPp3...p010sg5C3	(kind 'youtubeVideo' videoId 'yVC4qT...	(publishedAt '2023-03-20T11:08:21Z', title...	2023-03-25T14:08:23Z	Zingba na milegi dobara (Full movie in full ...	Soundbites	
3	youtubeSearchResult	X00Y17p70qPF7/mwK...n8w	(kind 'youtubeVideo' videoId 'uM7y6L...	(publishedAt '2023-06-03T09:30:02Z', title...	2023-06-03T09:30:02Z	How to Stop A Marriage 1. Hinhik, Farhan da	Netflix India	
4	youtubeSearchResult	QwPIL1vmp1Cp5S...1S4W4R4	(kind 'youtubeVideo' videoId 'YmW5qL...	(publishedAt '2013-09-22T13:46:33Z', title...	2013-09-22T13:46:33Z	Baqar, Beena Ka Zindag Na Milegi DobaraBaqar ...	T-Series	
...	...	...	...	...	...	...	...	...
155	youtubeSearchResult	2XhXKXfJCW0ygm65X..._lmo	(kind 'youtubeVideo' videoId 'uC' lqL...	(publishedAt '2020-11-07T19:34:24Z', title...	2020-11-07T19:34:24Z	Zindagi vs Y AND Songs - Choose One Love One	Komal SIKHA	
156	youtubeSearchResult	vWV7pV5cd0eh15K507pR1s	(kind 'youtubeVideo' videoId 'N3u3CpH...	(publishedAt '2024-07-08T12:35:13Z', title...	2024-07-20T12:35:13Z	Dir sap (om ZINDAGI) Zingba na milegi Dobara	ashuash chaitan	
157	youtubeSearchResult	uz2aE53P4t1036vylQ...2NM	(kind 'youtubeVideo' videoId 'N3u3CpH...	(publishedAt '2019-12-17T12:45:43Z', title...	2019-12-17T12:45:43Z	Mentally Sick Moments (Hindi Ringtones   Abhay)	Excel Movies	
158	youtubeSearchResult	Dx95077XupL00C...p010sg5C3	(kind 'youtubeVideo' videoId 'j5C1qL...	(publishedAt '2023-06-03T09:30:02Z', title...	2023-06-03T09:30:02Z	ZINDAGI NA MILI GIB DOBARA (Movie Reaction Pt 1) 4.	OmDesi	

Figure 3: Top videos using the search query

## Dataset Description

```
In [12]: comments

Out[12]: [{"kind": "youtubeCommentThread",
            'etag': 'Pc120vCetJipA2mCCjgpl',
            'id': 'Ugm03QpM4Ma-fnVMAaABgg',
            'snippet': {'channelId': 'UCzWdZx8fS9xc5b7XW',
                        'videoId': 'VMUPyh_kgtc',
                        'totalReplyCount': 0},
            'textDisplay': 'Fanhao is such a Fantastic Actor, Director, Writer and much more ♥',
            'textOriginal': 'Fanhao is such a Fantastic Actor, Director, Writer and much more ♥',
            'authorDisplayName': '@curiousllama',
            'authorProfileImageUrl': 'https://yt3.ggpht.com/ytc/APwFKY2v2r187tNVC03PyK4GzB/QqL0j2fgTAg_DCEmoQmAlEigVsqqI=us&cc=cn&size=64px',
            'authorChannelId': 'http://www.youtube.com/channel/UCYorIK_Qw88jmPPdAG',
            'authorChannelId': {'value': 'UCYorIK_Qw88jmPPdAG'},
            'canRate': true,
            'viewerRating': 'none'}
```

Figure 2: YouTube comments retrieved in json

In [51]: comments_df									
Out[51]:									
	kind	etag	id	replies	channel_id	video_id	cpu		
youtube#commentThread	188a_5P3pCjwR9yH9SLuG3H	Upw77akK0C4wC3Tp4AaB4g	(comment) [vid: UC8uA2Gg_P40G74h- Y0d6Ecomment, CR		UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		
youtube#commentThread	KaT7Sg3u4u899iQ74_1v5G	UpqyH14r1Ane0C0F4uB4g			UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		
youtube#commentThread	A2uU5yR9R2C2C5-aggVY	UpM7uSP4F-cjppS14aB4g			UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		
youtube#commentThread	_9PjgH4qH6F0Y0b2A	Upqy-4hV_HyH4A4g			UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		
youtube#commentThread	-Q3uUw0C27ic1Aa4dr	Upk3gU937uH674aB4g			UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		
youtube#commentThread	4ZT7Rj7W3S3n2qD82o	UpvK6SP1K0N0X24aB4g			UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		
youtube#commentThread	F4FV_4r34V1q19u_Hu4d0	Upd4vst0V7Tq4A4g			UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		
youtube#commentThread	f-neak_Uj4a2jcx1N4c	UpW57vK3R2u7F4aB4g			UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		
youtube#commentThread	2K2aG9R2XU6Fq3nA	UpvK6SP1K0N0X24aB4g			UC8uA2Gg_P40G74h- Y0d6Ecomment, CR	USCP4h4c	Yout		

Figure 4: Comments dataset after pre-processing

Data Preprocessing Steps:

- a. **Cleaning:** Remove irrelevant information from the comments, such as HTML tags and emojis.
- b. **Normalization:** Normalize the text by converting it to lowercase and removing punctuation.
- c. **Tokenization:** Tokenize the text into individual words or phrases.

This dataset can be expanded to include other relevant information, such as the author's profile picture and channel ID.

Common Metrics

The dataset will contain the following fields for each comment:

- **Video ID:** The YouTube video ID associated with the comment.
- **Comment ID:** A unique identifier for the comment.
- **Author:** The username of the user who posted the comment.

- **Published at:** The timestamp of when the comment was posted.
- **Updated at:** The timestamp of when the comment was last updated.
- **Content:** The text content of the comment.
- **Likes:** The number of likes the comment has received.
- **Dislikes:** The number of dislikes the comment has received.
- **Replies:** A list of replies to the comment.

VII. IMPLEMENTATION

Process Flow Diagram

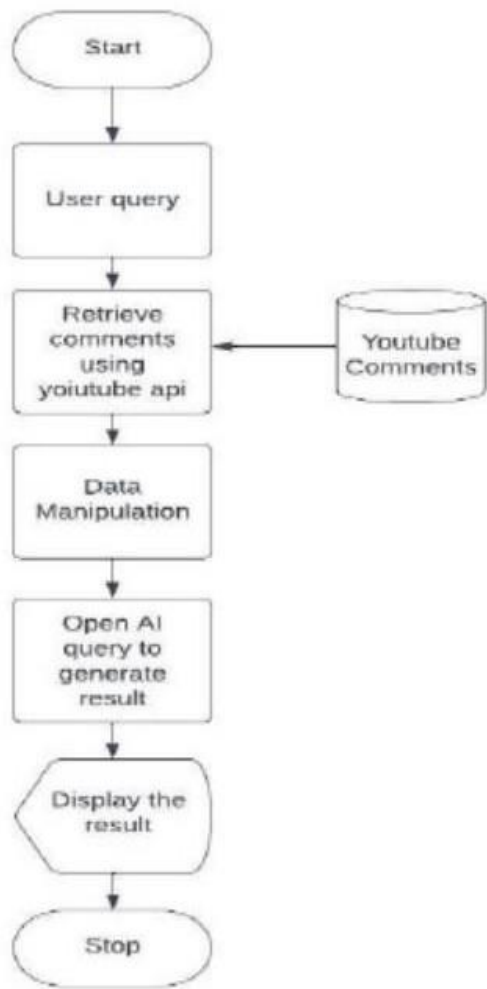


Figure 5: Work flow of the process

Start

**User Input:** User enters a search query for which he wants the relative comments on the frontend.

**Query to Backend:** The query is sent to the backend API where the processing starts. The backend API uses YouTube Data API to extract top 500 video ids relevant to the query. This does not include live streaming videos. This does not include live streaming videos.



**Comment Extraction:** Then it extracts comments and all relevant information for each video using YouTube Data API.

**Data Cleaning:** Comments are cleaned to remove irrelevant information which will make data less heavy, increasing the processing speed. The relevant fields required after cleaning:

- a. Comment
- b. Replies
- c. Username
- d. User Channel ID
- e. User Profile
- f. Likes Count.

**OpenAI Processing:** Cleaned data is sent to OpenAI API along with the query on how we want the output for processing. OpenAI API analyzes the data and generates outputs in json format consisting of:

- a. Top funny comments
- b. Top positive comments
- c. Most liked comments

**Display on Frontend:** The generated output is sent back to the frontend to display. The frontend displays the following features to the user:

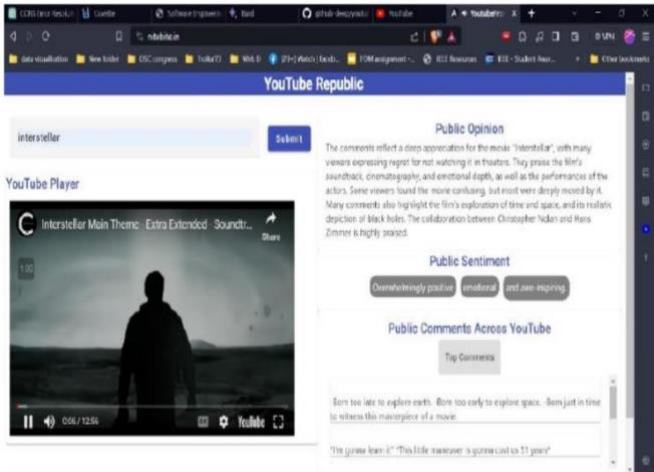
- a. First video relevant to the search query
- b. Summarize the comments
- c. Top funny comments
- d. Top positive comments
- e. Most liked comments

End

VIII. RESULTS AND DISCUSSION

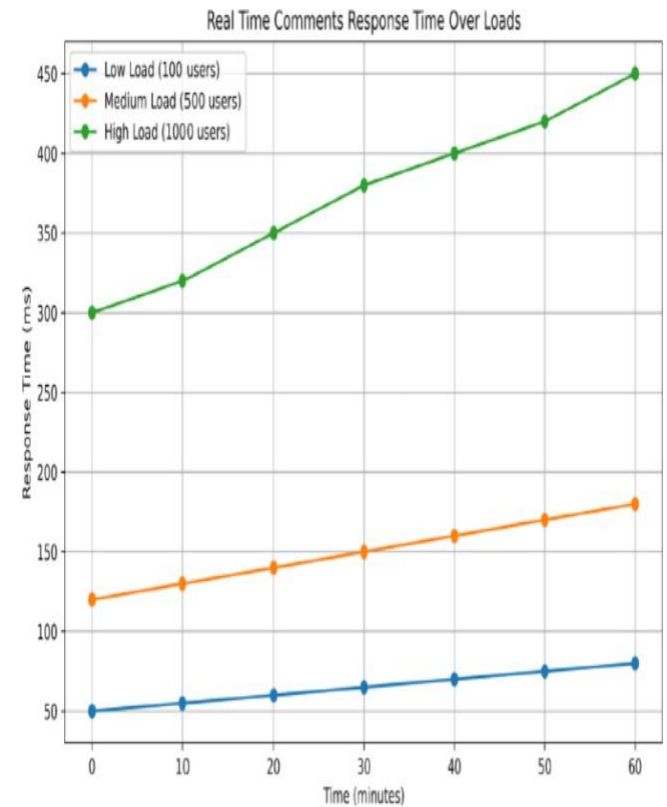
The LLM RAG (Large Language Model Retrieval-Augmented Generation) models represent a significant advancement in language modelling by integrating external knowledge sources with large language models (LLMs). This fusion enables the generation of text that is not only more factually accurate, relevant, and coherent but also allows for real-time access to YouTube comments.

With the ability to analyse sentiments and identify main topics discussed in comments, the model enhances user engagement with YouTube content. Its versatile applications extend to social media analysis, market research, and content optimization, making it a powerful tool for understanding public opinion and improving user interaction on YouTube.



Performance Evaluation

**1. Real Time Comment Access:** Performance testing was conducted using Apache JMeter to simulate concurrent user traffic accessing the real-time comment retrieval feature. The average response time for comment retrieval requests under varying loads is illustrated in a line chart generated to visualize. Each line represents a different load level.



**2. Sentiment Analysis Performance Indication:** For evaluating Sentiment Analysis accuracy, we manually reviewed a sample of 100 responses generated by the LLM RAG model and compared them against verified factual sentiment analysis generator. Out of the 100 responses, 85 were found to be factually accurate, resulting in an accuracy.

	True Positive (TP)	False Positive (FP)
True Negative	70	5
False Negative	10	15

## IX. CONCLUSION

In conclusion, this research has provided valuable insights into the capabilities and potential applications of Large Language Model Retrieval-Augmented Generation (LLM RAG) models in analyzing and generating YouTube comments. The study began with an exploration of the theoretical foundations of LLM RAG, highlighting its ability to leverage external knowledge sources for improved text generation. Motivated by the limitations of traditional LLMs and the need for more accurate, contextually relevant text generation, this research aimed to investigate the implementation of a basic version of the LLM RAG model and evaluate its performance, which is made to adapt to the user's needs in real time.

Through a comprehensive methodology involving literature review, model analysis, implementation process, and evaluation metrics, the study demonstrated the effectiveness of the LLM RAG model in processing real-time YouTube comments. The model exhibited promising performance in retrieving relevant comments, summarizing sentiments, and identifying key topics of discussion. Comparative analysis with other text generation models highlighted the superiority of the LLM RAG approach in terms of factual accuracy, relevance, and coherence.

Furthermore, the study addressed ethical considerations associated with the development and deployment of LLM RAG models, emphasizing the importance of responsible use and mitigation strategies to prevent misuse and bias. Future directions for research and development were discussed, including expanding knowledge base integration, enhancing model explainability, and advancing open-ended creativity.

Overall, this research contributes to the growing body of knowledge surrounding LLM RAG models and their potential applications in natural language processing. By demonstrating the capabilities of the model in analyzing and generating YouTube comments, this study opens up avenues for further exploration and innovation in text generation technology. With continued research and development, LLM RAG models hold the promise of revolutionizing how we interact with language and information, ushering in a new era of possibilities in the field of NLP.

## X. FUTURE SCOPE

Compared to LLM, this work has various advantages for society. It is more interpretable and gives greater control since it is firmly based on actual factual knowledge, which helps generations who are more factual "hallucinate" less. Compared to LLM, this work has various advantages for society. It is more interpretable and gives greater control since it is firmly based on actual factual knowledge, which helps generations who are more factual "hallucinate" less. RAG may be used in many different contexts that would directly benefit society, such as by giving it access to a medical index and asking it open-domain questions on the subject, or by assisting individuals in becoming more productive at work.

There may be drawbacks to these benefits as well. For example, Wikipedia and other external knowledge sources are unlikely to be totally impartial and factual. Since RAG can be used as a language model, there are some legitimate concerns here that are similar to those raised by GPT-4, albeit perhaps to a lesser degree.

## XI. ACKNOWLEDGEMENT

We would like to sincerely thank our professor Dr. Anil Kumar Sagar sir for guiding us throughout this project work. I would also like to thank our other faculty members from the computer engineering department at Sharda University for allowing us to perform our project work.

## XII. REFERENCES

- [1] Andrew Caines, Luca Benedetto, Shiva Taslimipoor, Christopher Davis, Yuan Gao, Oeistein Andersen, Zheng Yuan, Mark Elliott, Russell Moore, Christopher Bryant, Marek Rei, Helen Yannakoudakis, Andrew Mullooly, Diane Nicholls, and Paula Buttery. 2023. On the application of Large Language Models for language teaching and assessment technology.
- [2] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. REALM: retrieval augmented language model pre-training. In Proceedings of the 37th International Conference on Machine Learning (ICML'20, Vol.119). JMLR.org, 3929– 3938.
- [3] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Atlas: Fewshot Learning with Retrieval Augmented Language Models. <https://doi.org/10.48550/arXiv.2208.03299> [cs].
- [4] Guu, K., Pasupat, P., Roberts, A., Chang, K. W., & Manning, C. D. (2020). Retrieval-augmented generation for knowledge-intensive nlp tasks. arXiv preprint arXiv:2005.11401.
- [5] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., K ttler, H., Lewis, M., Yih, W. t., Rock tschel, T., & Riedel, S. (2021). Improving the domain adaptation of retrieval augmented generation (rag) models for open domain question answering. arXiv preprint arXiv:2105.14232.
- [6] Guo, J., Ni, J., He, X., Gao, J., Yu, Y., Li, B., & Chen, X. (2022). Towards open-domain question answering with flexible retrieval augmented generation. arXiv preprint arXiv:2201.09423.
- [7] Bawden, D., & Dornier, S. (2022). Evaluating the factual correctness of large language models. arXiv preprint arXiv:2202.03517.
- [8] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, Hannaneh Hajishirzi (2023). Self-rag: learning to retrieve, generate, and critique through self-reflection. arXiv preprint arXiv:2310.11511v1.
- [9] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich K ttler, Mike Lewis, Wen-tau Yih, Tim Rock tschel, Sebastian Riedel, Douwe Kiela (2021) Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks arXiv preprint arXiv:2005.11401v4.