

LAMPIRAN

29 April 2024

Dear Kurnia Muludi, Kaira Milani Fitria, Joko Triloka, Sutedi,

Congratulations, your submitted paper titled “Retrieval-Augmented Generation Approach: Document Question Answering using Large Language Model” was reviewed and Accepted for Publication in International Journal of Advanced Computer Science and Applications (IJACSA) Volume 15 No 3 March 2024.

Paper Title: Retrieval-Augmented Generation Approach: Document Question Answering using Large Language Model

Authors: Kurnia Muludi, Kaira Milani Fitria, Joko Triloka, Sutedi

The acceptance of your paper for publication in IJACSA reflects the high status of your work by your fellow professionals in the field.

Bibliographic Information

U.S ISSN: 2156-5570(Online)

U.S ISSN: 2158-107X(Print)

Publication Frequency: Monthly

Upon publication of papers, our next steps will be to submit all published papers in International Indexes and University Libraries. Some of the indexes include Web of Science, Scopus, Inspec, Ebescio, Microsoft Academic, WorldCat. IJACSA is also indexed in the Thomson Reuters Emerging Sources Citation Index and is also listed in the Thomson Reuters Master Journal List - <http://science.thomsonreuters.com/cgi-bin/jrnlst/jlresults.cgi?PC=MASTER&ISSN=2158-107X>

All authors are deemed to be individually and collectively responsible for the content of papers published by The Science and Information (SAI) Organization.

Hence, it is the responsibility of each author to ensure that papers submitted to The Science and Information (SAI) Organization attain the highest ethical standards with respect to plagiarism.

Regards,
Managing Editor

IJACSA
U.S ISSN: 2156-5570 (Online)
U.S ISSN: 2158-107X (Print)
The Science and Information (SAI) Organization
editorijacsa@thesai.org
www.ijacsa.thesai.org

IJACSA

Submission date: 20-Feb-2024 03:58PM (UTC+1100)

Submission ID: 2298820201

File name: IJACSA_revised.pdf (1.67M)

Word count: 5939

Character count: 33854

Retrieval-Augmented Generation Approach: Document Question Answering using Large Language Model

17 Kurnia Muludi¹, Kaira Milani Fitria^{2*}, Joko Triloka³, Sutedi⁴

Department of Computer Science, Informatics & Business Institute Darmajaya, Lampung, Indonesia
kurnia@darmajaya.ac.id¹, kairaamilanii@gmail.com^{2*}, joko.triloka@darmajaya.ac.id³, sutedi@darmajaya.ac.id⁴

Abstract— This study introduces the Retrieval Augmented Generation (RAG) method to improve Question-Answering (QA) systems by addressing document processing in Natural Language Processing problems. It represents the latest breakthrough in applying RAG to document question and answer applications, overcoming previous QA system obstacles. RAG combines search techniques in vectorstore and text generation mechanism by Large Language Models, offering a time-efficient alternative to manual reading limitations. The research evaluates RAG's that using GPT-3.5-turbo from ChatGPT model impact on document data processing, comparing it with other applications. This research also provides datasets to test the capability of the QA document system. The proposed dataset and SQuAD dataset are used for performance testing. The study contributes theoretically by advancing methodologies and knowledge representation, supporting benchmarking in research communities. Results highlight RAG's superiority: achieving a precision of 0.74 in ROUGE, outperforming others at 0.5; obtaining an F1 score of 0.88 in BERTScore, surpassing other QA apps at 0.81; attaining a precision of 0.28 in BLEU, surpassing others with a precision of 0.09; and scoring 0.33 in Jaccard Similarity, outshining others at 0.04. These findings underscore RAG's efficiency and competitiveness, promising a positive impact on various industrial sectors through advanced AI technology.

Keywords—Natural Language Processing, Large Language Model, Retrieval Augmented Generation, Question Answering, GPT

I. INTRODUCTION

This research proposes a new approach to the increasing reliance on articles and journal documents by introducing a Question-Answering (QA) document processing system [1]. The identification of several critical problems motivates this research. Manual processing for document understanding leads to time-consuming efforts, susceptibility to human error, and inefficient analysis processes. Based on previous methods, the use of modified Large Language Models (LLM) for document processing requires significant resources and pose challenges for widespread implementation. Also, the underutilisation of the recently discovered Retrieval Augmented Generation (RAG) method, particularly in document processing within Question-Answering (QA) systems, provides an opportunity for further exploration. The motivation stems from the challenges associated with manual document processing, resource-intensive Large Language Model (LLM) modifications, and the underutilisation of the Retrieval-Augmented Generation (RAG) method in the document-based question-answering domain [2],

[3]. In addition, there is a tendency to produce hallucinatory responses that lack accuracy and precision in models that rely solely on LLM capabilities for QA systems without modifications. Finally, the implementation of RAG in QA systems for document processing offers the untapped potential to improve the ability of the system to produce accurate and non-hallucinatory responses.

Building on this line of research, this paper proposes the implementation of the Retrieval Augmented Generation model for document question answering tasks, specifically using the ChatGPT model. RAG, introduced in 2021 [4], addresses the limitations of previous methods by merging parametric and non-parametric memory. This hybrid model seamlessly integrates generative capabilities with data retrieval mechanisms, linking language models to external knowledge sources. The distinct advantage of RAG lies in its ability to adapt to dynamic data, its flexibility in working with external data sources, and its ability to mitigate hallucinatory responses [5]. These characteristics make RAG particularly suitable for QA tasks on internal organisational documents, by leveraging external knowledge to reduce response hallucination [6]. The current research aims to exploit the innovative approach of RAG to construct an application capable of automatically processing external text documents. Users can ask questions related to the document content and the system responds based on the processed document. This comprehensive solution not only overcomes the limitations of previous methods, but also promises to significantly speed up research and study exploration in various domains.

II. RELATED WORKS

This study examines the applicability of RAG, its impact on the document processing task, and compares it to the previous methods. This research also investigates the capability of the large language model within the ChatGPT systems, gpt-3.5-turbo within the framework of RAG. This work also highlights the development of Artificial Intelligence (AI) and Natural Language Processing (NLP), so this research focuses on the improvement of intelligence and capabilities of applications [7], [8], [9]. Machine Learning and Deep Learning algorithms, which include BERT Base, and Text-to-Text Transfer Transformer (T5) models, and RAG method, have made significant advances in QA tasks [4], [10], [11], [12]. This research motivated the implementation of RAG for processing documents, integrated in an interactive QA system.

Between 2015 and the present, the evolution of question-answering (QA) systems shows a trajectory characterised by diverse methodologies. Starting with semantic parsing-based systems in 2015, Wen-tau Yih et al. focused on transforming natural language queries into structured logical forms, achieving a performance of 52.5% in F1-score [2]. Subsequent knowledge-based paradigms (KB-QA) by Yanchao Hao et al. in 2017 reformulated questions as predicates, achieving a performance of 42.9% [3]. Progress has been made in integrating of AI technologies. Caiming Xiong's exploration of dynamic memory networks (DMN) in 2016, which achieved an accuracy of 28.79% [7]. In the same year, Minjoon Seo et al.'s Bi-Directional Attention Flow (BiDAF) framework demonstrated significant performance with 68% exact match and 77.3% F1 score, albeit with a computational time of 20 hours [8]. Adams Wei Yu et al. introduced the QANet model in 2018, with a performance of 76.2% exact match and 84.6% F1-Score, within a shorter computational time of 3 hours [9]. As QA systems evolve, in 2019 Wei Yang et al. applied fine-tuning methods with data augmentation techniques, achieving remarkable results with a modified BERT-Base model of 49.2% for exact match and 65.4% for F1-Score [10]. Colin Raffel et al. introduced the Text-to-Text Transfer Transformer (T5), with impressive performance of 63.3% for exact match, 94.1% for F1 score, and a peak accuracy of 93.8%, albeit with an increased number of parameters of 11 billion [11]. In 2020, the focus was on fine-tuning pre-trained models, with Adam Roberts et al. achieving a recall performance of 34.6% using the T5 model [12]. The Retrieval-Augmented Generation (RAG) method, which combines parametric and non-parametric methods, was introduced by Patrick Lewis et al. in 2021. RAG has demonstrated its capabilities in open domain QA tasks, overcoming previous limitations to deliver more efficient and comprehensive QA systems [4].

Large language model called GPT, or Generative Pre-Trained Transformer is developed by OpenAI. Previous research has compared the performance of ChatGPT with other large language models like PaLM and LLaMA in open-domain QA tasks indicates that ChatGPT consistently achieves the highest scores across various open-domain QA datasets [13]. Table 1 presents performance comparisons among LLMs.

TABLE I. LLM PERFORMANCE ON OPEN DOMAIN QA DATASET

Model	TriviaQA	WebQuestion	NQ-Open
PaLM-540B (few-shot)	81.4	43.5	39.6
PaLM-540B (zero-shot)	76.9	10.6	21.2
LLaMA-65B (zero-shot)	68.2	-	23.8
ChatGPT (zero-shot)	85.9	50.5	48.1

The PolyQuery Synthesis test, which identifies multiple queries within a single-query prompt and extracts the answers to all of the questions from the model's latent representation, also shows that ChatGPT outperforms other GPT models from OpenAI (ada-001, babbage-001, curie, and davinci) in terms of accuracy [13]. According to the evaluations, the gpt-3.5-turbo model has been selected for implementation in this research.

III. RESEARCH METHOD

This research undergoes a development phase, starting with designing the application system and integrating the APIs of ChatGPT, LangChain and FAISS. Subsequent stages include extensive system modelling, interface testing and data preparation using the proposed dataset and the SQuAD dataset. The testing phase, which includes a performance comparison with other applications using ground truth metrics (ROUGE, BERTScore, BLEU and Jaccard Similarity), guides the exploration of the capabilities of the proposed system, as shown in the Figure 1.

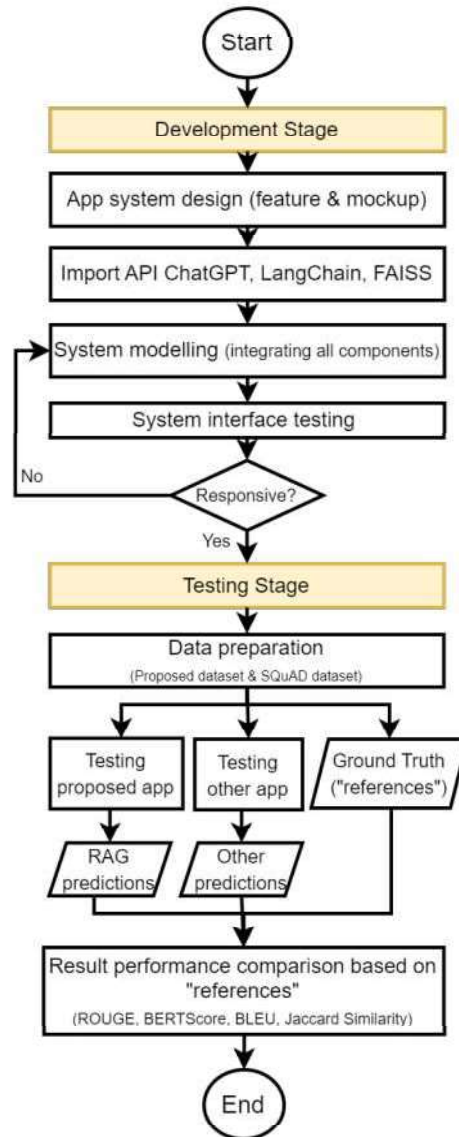


Fig. 1. Research flow diagram.

A. RAG Integration

Retrieval Augmented Generation (RAG) combines retrieval and generation models. It uses a Large Language Model (LLM) to generate text based on commands and integrates information from a separate retrieval system to improve output quality and contextual relevance [14]. The mechanism involves retrieving factual content from a knowledge base via retrieval models and using generative processes to provide additional context for more accurate output [15]. External data sources are used, and the numerical representation is facilitated by embedding methods to ensure compatibility. Based on Figure 2, user queries converted into embeddings are compared with vectors from the knowledge library. Relevant context is added to the queries before they are fed into the base language model.

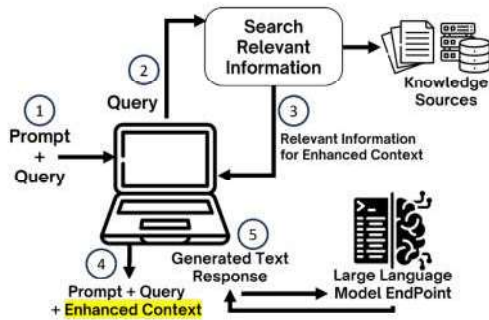


Fig. 2. RAG mechanism with LLM.

OpenAI, the creator of the Large Language Model GPT, conducted comprehensive a few RAG experiments, exploring various implementations such as cosine similarity retrieval, chunk/embedding experiments, reranking, classification steps, and prompt engineering, as depicted in Figure 3. OpenAI's findings, presented in Figure 3, revealed that RAG implementation with prompt engineering achieved the highest accuracy, positioning it as the most effective RAG technique to date [16]. This discovery serves as a catalyst for the integration of RAG with prompt engineering using the LangChain module.

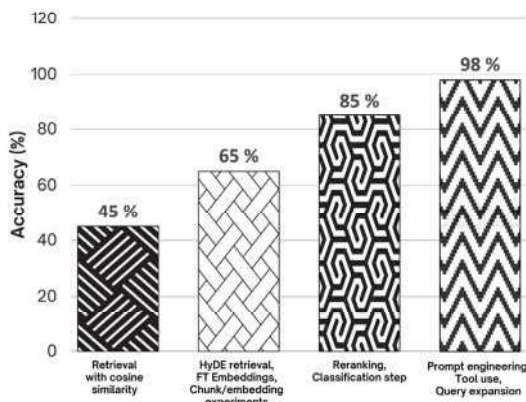


Fig. 3. Accuracy of the RAG method by OpenAI

LangChain provides a robust data processing pipeline that utilises FAISS to perform an efficient retrieval operation in the VectorDB. The query phase transforms inputs into vectors for database searches, and prompt engineering enhances the reusability of retrievals. Output parsers interpret LLM outputs, ensuring consistency [17]. A highly efficient similarity search and vector clustering library, Facebook AI Similarity Search or FAISS [18]. It optimises the trade-off between memory, speed and accuracy, allowing developers to effectively navigate multimedia documents. The mechanism involves the construction of an index for efficient storage, with vector searches retrieving the most similar vectors using cosine similarity scores [19].

B. Proposed Model

This research employs a modified Large Language Model (LLM), ChatGPT, augmented with additional libraries to function as a Question-Answering (QA) system capable of processing external documents for supplementary information. The chosen methodology for QA system development is the Retrieval Augmented Generation (RAG) mechanism. Unlike previous approaches such as semantic parsing-based, knowledge-based, and fine-tuning using LSTM or other DL algorithms, RAG addresses shortcomings like difficulty expanding or revising model memory, inability to provide direct insight into generated predictions, and a tendency to produce hallucinative answers [12]. The solution involves the creation of a hybrid model, merging generative and retrieval models, forming the basis for the RAG method. RAG offers advantages such as adaptive responses to dynamic data, flexibility with external data sources, and minimization of hallucinative responses [5]. Thus, RAG is chosen to construct a text document-based QA system interacting with users through a chatbot interface. The system's workflow, implemented using RAG and supporting libraries like LangChain and FAISS, is illustrated in Figure 4.

The integration of the LangChain framework into the QA document system includes document loading, memory management, and prompting to connect to the LLM model. The process starts with document loading, followed by document splitting into text chunks. These text chunks undergo word embedding, converting them into vectors stored in the vector database. Simultaneously, user-inputted text questions are embedded and converted into word vectors. The system connects these vectors to the vector database, performing semantic search and ranking the relevance between vectors. The semantic search results in relevant context between questions and answers. The system retrieves pertinent answers based on user queries and sends them to the LLM (using the ChatGPT model). The final outcome involves the system receiving LLM-generated answers and delivering them to the user. The application system interacts with users, requiring an interface connecting the user and the system. Mockups, design layouts, and elements for the web application are created using Streamlit framework, facilitating rapid development and sharing of the AI model web application. The mockup for the application system and user interaction within the system is depicted in Figure 5.

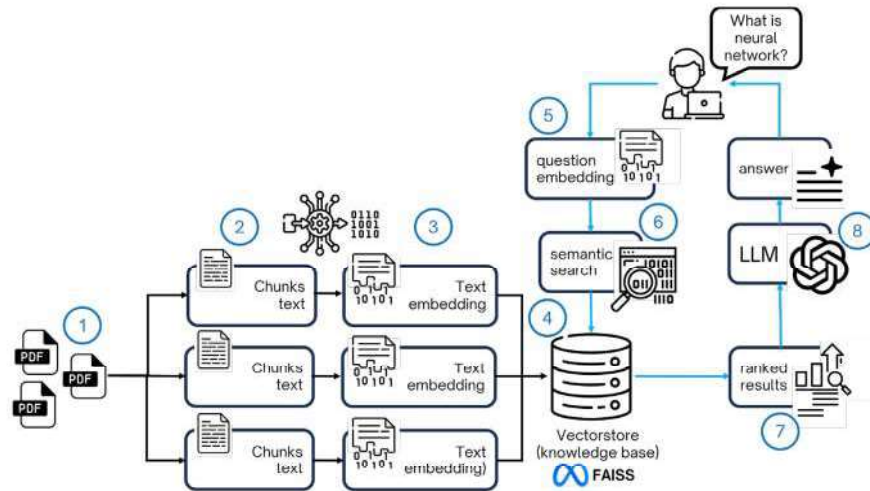


Fig. 4. Integration of LangChain framework in RAG for proposed document QA system

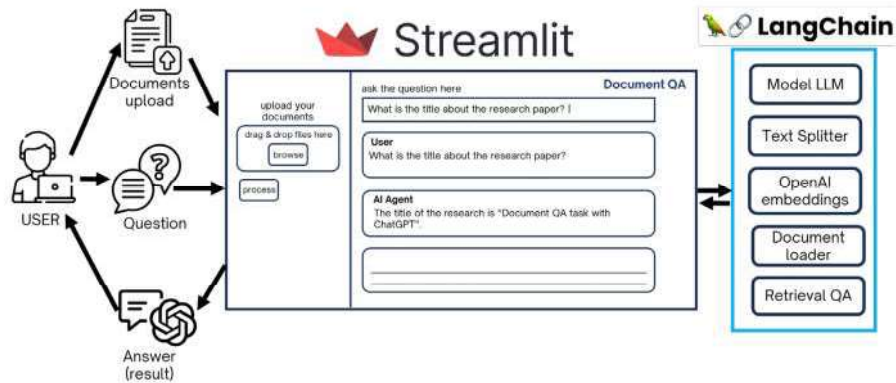


Fig. 5. Mockup of the application system and user interaction for the app

C. Proposed Dataset : DocuQA

The proposed dataset, DocuQA, designed for application-based question-answering systems that process document inputs, consists of 20 diverse documents, encompassing journal articles, news reports, financial documents, and tutorials. Each document file includes 5 questions with corresponding ground truth answers, enabling a thorough evaluation of QA system capabilities, with total 100 questions in dataset. DocuQA consist of journal documents with calculations and formulas, news documents with specific titles, report/financial report/news documents with numbers and currency data, and tutorial documents with step-by-step instructions. Accuracy can be calculated based on correct answers out of 100, providing a metric for information extraction accuracy. The dataset aims to challenge QA systems in understanding context, identifying keywords, and efficiently extracting specific information, offering a robust evaluation tool for developers and researchers across various document and question types. The dataset can be accessed publicly [20].

Proper citation of the dataset is encouraged for research or projects using DocuQA to ensure appropriate credit is given. The preview of the DocuQA dataset can be shown in Figure 6.

Files	Question	Ground Truth
R4	Can you inform the key numbers of fourth-quarter vehicle production and deliveries report for 2023 from Tesla?	Total deliveries Q4 2023 is 484.507 Total production Q4 2023 is 494.989 Total annual deliveries 2023 is 1,808.581 Total annual production 2023 is 1,845.965
R4	How many electric vehicle deliveries and production based on Tesla's report in 2022?	1.31 million deliveries and 1.37 million production
R4	How many units of Chinese automaker BYD's new energy vehicles were sold in 2023?	3.02 million
R5	what is the title of the news report?	Copper could skyrocket over 75% to record highs by 2025 — brace for deficits analysts say
R6	when the news published?	January 2 2024

Fig. 6. Preview of DocuQA test dataset

D. Testing and Evaluation

The tests were performed on two types of test datasets, with DocuQA [20] and SQuAD 1.1 [21]. DocuQA as a dataset originally created by this research, consisting of 100 questions with the ground truth and a total of 20 test documents for document-based QA systems. In addition, the SQuAD dataset was used in the form of modified pdf documents that can be used to test the QA system's ability to process documents and retrieve information based on the questions and related ground truth in the SQuAD dataset. Both types of test datasets will be tested on the QA system developed in this research, and also on other commercial QA systems that process pdf documents, such as typeset.io. The results of these tests will give an idea of QA system performance built by this research, whether it is superior to other document-based QA applications.

The proposed QA document processing system is evaluated through rigorous testing using established metrics such as ROUGE or Recall-Oriented Understudy for Gisting Evaluation, BERTscore, BLEU or Bilingual Evaluation Understudy, and Jaccard Similarity. These metrics provide reliable benchmarks for assessing the system's performance across various dimensions. The testing process involves two key variables. "Predictions RAG" and "Prediction Others" represent the test results from the developed application and comparable commercial applications, respectively. Both sets of predictions are compared to the ground truth data, which is encapsulated in the "references" variable. Different aspects of language models and question answering systems are evaluated using different metrics. ROUGE measures the overlap in summarization [22]. BERTscore assesses semantic similarity using contextual embeddings [23]. BLEU evaluates n-gram precision [24], and Jaccard Similarity compares text similarity based on word or n-gram overlap [25]. Precision in question answering systems is commonly assessed through accuracy, F1 score, and precision metrics, providing insights into their effectiveness. The metrics are used to quantitatively evaluate system performance and establish its superiority over existing commercial applications in document processing and information retrieval tasks.

1) Accuracy: Accuracy defined as the proportion of correct responses from the total number of responses. Accuracy can be calculated by calculating the percentage of correct predictions over the total number of references [26]. In essence, accuracy represents the ability of the system to provide correct answers, which is expressed as a percentage using the following formula (see Eq.1).

$$Accuracy = \frac{\text{correct predictions}}{\text{all predictions}} \times 100\% \quad (1)$$

This metric serves as a valuable indicator of the overall correctness of the model in the response it generates.

2) ROUGE: Recall-Oriented Understudy for Gisting Evaluation can be used to evaluate the text generation models, which is based on the measurement of the overlap between candidate text and reference text [27]. ROUGE has several measurement variants, each depending on the number of overlapping n-grams. The ROUGE-L variant is the most widely

used, because it uses the longest sequence or longest common subsequence or LCS with the longest word sequence that both sentences have. Precision refers to the proportion of n-grams in the candidate that are also in the reference (see Eq. 2.). Recall, on the other hand, refers to the proportion of n-grams that are in the reference text that exactly match in the predicted candidate text (see Eq. 3). The F1-score can be calculated from the precision and recall scores (see Eq. 4).

$$ROUGE-L_{recall} = \frac{LCS(\text{candidate}, \text{reference})}{\# \text{words in reference}} \quad (2)$$

$$ROUGE-L_{precision} = \frac{LCS(\text{candidate}, \text{reference})}{\# \text{words in candidate}} \quad (3)$$

$$ROUGE-L_{F1-Score} = 2 \times \frac{\text{recall} \cdot \text{precision}}{\text{recall} + \text{precision}} \quad (4)$$

Where the reference based on the ground truth in test dataset, and the candidate is from the system predictions. The score generated by the ROUGE measure is between 0 and 1. A score of 1 indicates total agreement between reference and candidate text.

3) BERTScore: BERTScore is an automatic evaluation metric in text generation tasks that evaluates the similarity of each candidate sentence token to each reference sentence token by means of contextual embeddings [23]. The embeddings in BERTScore are contextual, changing depending on the sentence context. The context awareness allows BERTScore to score semantically similar sentences despite their different sentence order. For the recall calculation, each token in x is matched with the most similar token in \hat{x} , as for the precision calculation. Greedy matching is used to maximise the similarity score. The values of precision (see Eq. 5.), recall (see Eq. 6.) and F1 score (see Eq. 7.) for reference x and candidate \hat{x} can be calculated using the following equations.

$$R_{BERT} = \frac{1}{|x|} \sum_{x_i \in x} \max_{\hat{x}_j \in \hat{x}} x_i^T \hat{x}_j \quad (5)$$

Where R_{BERT} is the Recall BERTScore, x is the reference token, \hat{x} is the candidate token, x_i is the sequence vector x , \hat{x}_j is the sequence vector \hat{x} , where $\sum_{x_i \in x}$ is the number of x_i present in x , and also $\max_{\hat{x}_j \in \hat{x}}$ is the maximum value of \hat{x}_j present in \hat{x} , and $x_i^T \hat{x}_j$ is the cosine similarity of x and \hat{x} .

$$P_{BERT} = \frac{1}{|\hat{x}|} \sum_{\hat{x}_j \in \hat{x}} \max_{x_i \in x} x_i^T \hat{x}_j \quad (6)$$

Given P_{BERT} as Precision BERTScore, x as reference token, \hat{x} as candidate token, x_i as sequence vector x , \hat{x}_j as sequence vector \hat{x} , where $\sum_{\hat{x}_j \in \hat{x}}$ is the number of \hat{x}_j present in \hat{x} , and also $\max_{x_i \in x}$ is the maximum value of x_i present in x , and $x_i^T \hat{x}_j$ is the cosine similarity of x and \hat{x} .

$$F_{BERT} = 2 \times \frac{P_{BERT} \cdot R_{BERT}}{P_{BERT} + R_{BERT}} \quad (7)$$

Where F_{BERT} is the F1-score of BERTScore, then P_{BERT} is the precision and R_{BERT} is the recall from BERTScore results. Although the cosine similarity value is theoretically in the interval $[-1,1]$, in practice the value is rescaled so that it is between 0 and 1 in the result of the BERTScore calculation.

4) BLEU: Bilingual Evaluation Understudy is a metric that computes a modification of precision for n-grams, combines it with weights, and applies a brevity penalty to obtain the final BLEU score [28]. The scores range of BLEU is from 0 to 1. The greater the BLEU score, the better the system's performance are considered to be compared to the references. The formula for calculating BLEU can be seen in Eq. 8.

$$BLEU = BP \cdot \exp\left(\sum_{n=1}^N w_n \log p_n\right) \quad (8)$$

BP represents the brevity penalty, adjusting the score to penalize translations shorter than the reference. N denotes the maximum number of considered n-grams. The precision for n-grams, denoted as p_n signifies the n-grams ratios by the candidate text that appearing in any reference translation to the total of n-grams in the candidate text. w_n represents the weight assigned to each n-gram precision score.

5) The Jaccard similarity quantifies the similarities percentage between two sets of data by identifying the common and the different members. This can be calculated by divide the number of observations shared by the sum of the observations in each of the two sets. Jaccard similarity can be expressed as the ratio of the intersection ($A \cap B$) to the union ($A \cup B$) of two sets (see Eq. 9).

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|} \quad (9)$$

$|A \cap B|$ indicates the size of the intersection of the sets A and B, and $|A \cup B|$ indicates the size of the union of the sets A and B. The Jaccard similarity is bounded in the range from 0 to 1. A Jaccard similarity of 1 indicates complete identity between the sets, while a similarity of 0 implies that the sets have no common elements.

IV. RESULT AND DISCUSSION

The interface of the proposed document QA system can accept multiple PDF format documents. If the user clicks the submit button, the system will process the PDF document to convert it to vector form with embedding (as described in the RAG mechanism in Figure 2). Once the document submission process is complete, the user can ask questions related to the submitted document, and the QA system will provide answers based on the source documents provided. The set of questions and answers generated from the user's interaction with the QA system will be in the form of a chatbot, so that it stores the

communication history. The interface of the proposed QA system can be seen in Fig 7.



Fig. 7. Document QA System Interface

A. Accuracy

Accuracy in our system model is expressed as the percentage of correct answers within the entire answer key dataset. To assess accuracy, we calculate the ratio of the number of correct predictions to the total number of predictions [26]. The visualization of this accuracy result can be figured in Figure 8.

The accuracy comparison between the proposed QA document system and other applications reveals the superiority of our method. The proposed system achieved accuracy rates of 96% (our dataset) and 95.5% (SQuAD dev dataset), surpassing the other application's rates of 55% (our dataset) and 85.7% (SQuAD dataset). This underscores the consistently higher accuracy of our proposed method.

B. ROUGE

ROUGE-L score evaluation compares the results of our proposed QA method outperforming other QA applications in terms of precision, recall, and F1-Score. Specifically, on our dataset, our proposed method demonstrated precision, recall, and F1-Score of 73.7%, 23.9%, and 33.7%, respectively. In comparison, other QA applications achieved lower performance metrics with precision, recall, and F1-Score of 50.0%, 10.5%, and 15.2%, respectively. Similarly, on the SQuAD dev dataset, our proposed method excelled with precision, recall, and F1-Score reaching 85.5%, 16.2%, and 26.1%, while other QA applications reported lower scores of

77.2%, 10.4%, and 17.1%, respectively. These results underscore the superior performance of our proposed method across both datasets that can be visualized in Figure 9.

C. BERTScore

BERTScore evaluation compares the results of our proposed QA method outperforming other QA applications in terms of precision, recall, and F1-Score. Specifically, on our dataset, our proposed method demonstrated precision, recall, and F1-Score of 85.2%, 90.1%, and 87.6%, respectively. In comparison, other QA applications achieved lower performance metrics with precision, recall, and F1-Score of 81.6%, 86.3%, and 83.8%, respectively. Similarly, on the SQuAD dev dataset, our proposed method excelled with precision, recall, and F1-Score reaching 82.8%, 87.0%, and 84.8%, while other QA applications reported lower scores of 80.4%, 86.3%, and 83.2%, respectively. These results underscore the superior performance of our proposed method across both datasets that can be visualized in Figure 10.

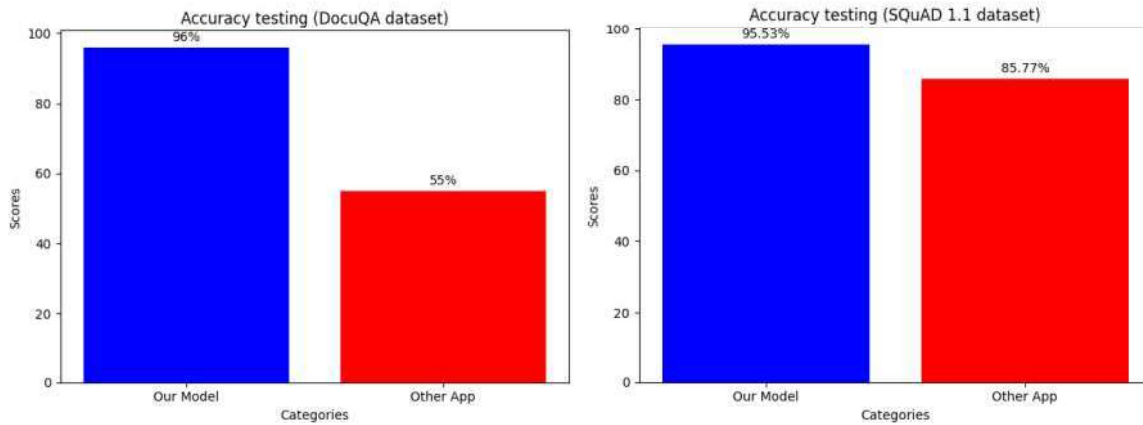


Fig. 8. Accuracy result of proposed method using RAG and other document QA application

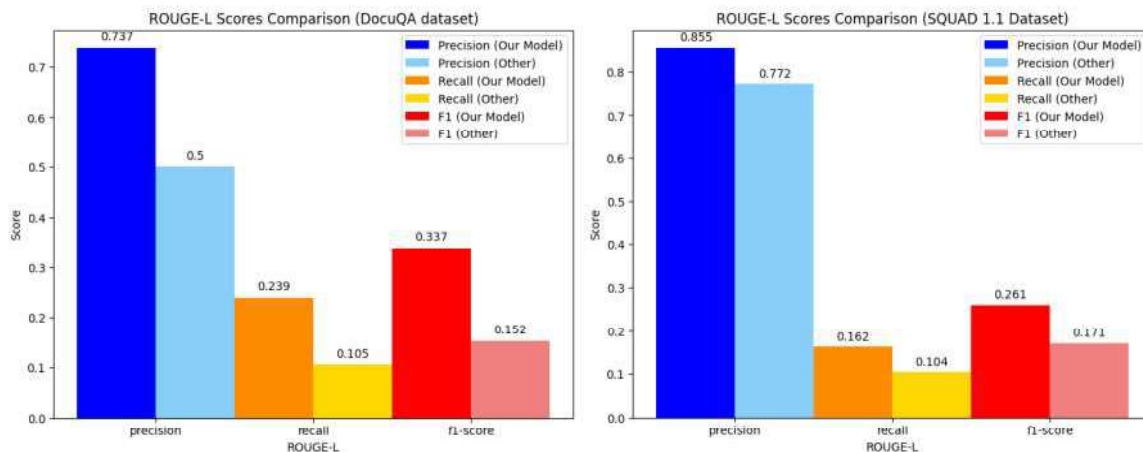


Fig. 9. ROUGE-L result of proposed method using RAG and other document QA application

D. BLEU Accuracy

The BLEU metric score taken is the precision value, to capture the ability of each model to extract keyword answers that match the ground truth. Specifically, on our dataset, our proposed method demonstrated precision of 28.2%. In comparison, other QA applications achieved lower performance precision 9.7%. Similarly, on the SQuAD dev dataset, our proposed method excelled with precision 17.7%, while other QA applications reported lower scores of precision 5.6%. These results underscore the superior performance of our

proposed method across both datasets that can be visualized in Figure 11.

E. Jaccard Similarity

The performance of our QA system, as evaluated through Jaccard Similarity, is outstanding. Our method achieved 33.3% on our dataset and 11.1% on SQuAD dev using RAG method. In comparison, other QA applications scored lower with 4.1% on our dataset and 9.1% on SQuAD dev. These results highlight our method's superiority in Jaccard Similarity on both datasets that can be visualized in Figure 12.

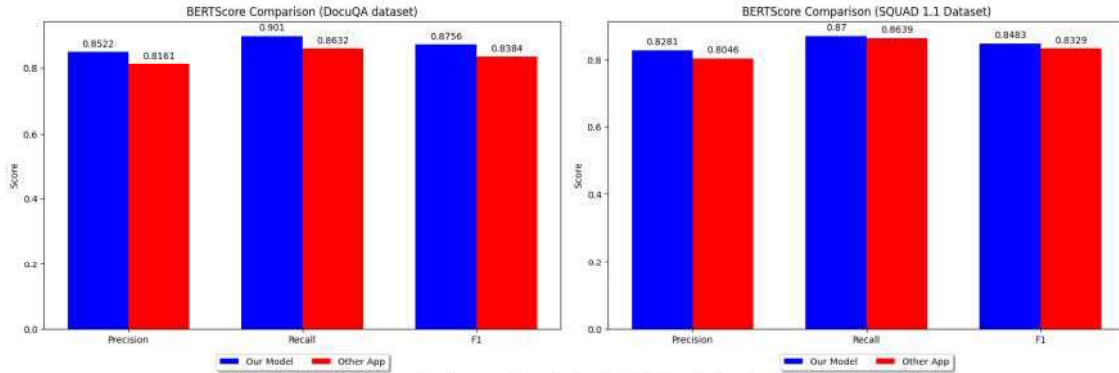


Fig. 10. BERTScore result of proposed method using RAG and other document QA application

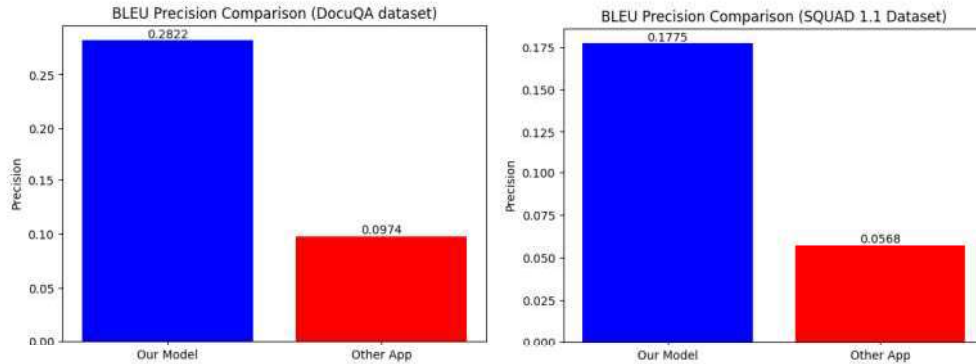


Fig. 11. BLEU precision result of proposed method using RAG and other document QA application

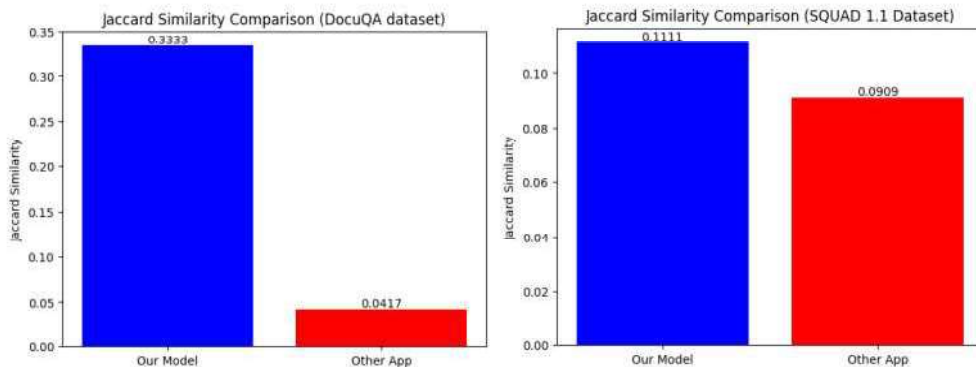


Fig. 12. Jaccard Similarity result of proposed method using RAG and other document QA application

V. CONCLUSION

Our proposed model for Question-Answering (QA) document processing integrates the Retrieval-Augmented Generation (RAG) model. The evaluation of our proposed QA system demonstrates its superiority over existing commercial applications in terms of Accuracy, ROUGE-L scores, BERTScore metrics, BLEU precision, and Jaccard Similarity. The proposed method achieved high accuracy rates of 96% and 95.5% on our dataset and the SQuAD dev dataset, respectively, outperforming other applications tested on the same datasets. Our system's precision, recall, and F1-Score metrics were superior to those of other QA applications on both datasets, as highlighted by the ROUGE-L evaluation. Additionally, the BERTScore metrics consistently showed higher precision, recall, and F1-Score for our proposed method compared to other applications. In addition, our QA system has demonstrated superior performance in keyword extraction and text similarity compared to other applications, as assessed by BLEU precision and Jaccard Similarity.

The accuracy result of 95.5% outperforms other research with 61.5% [30] and 71.4% accuracy [31]. Our system's precision, recall, and F1-Score are 82.8%, 87%, and 84.8%, respectively, which surpass the precision of 62%, recall of 87%, and F1-Score of 67% reported in other research [32]. The proposed QA system's effectiveness is affirmed by the fact that it surpasses the recall result of other research with 42.70% [33] and outperforms other research [31], [34], [35] in terms of F1-Score, which is 42.6% [31], 49% [34], and 70.8% [35]. This positions it as a superior solution for automatic document processing and information retrieval tasks across diverse domains.

In the future, studies could be conducted to refine the architecture of the system, explore additional ways of using external data, and improve the scalability of the model for broader applications. Integration of user feedback mechanisms and continuous learning modules could contribute to the adaptability of the system and further improve its accuracy over time. In addition, exploring ways of processing documents in real time and extending the system's compatibility with different document formats could open up new opportunities for research and study.

REFERENCES

- [1] F. Ganier and R. Querec, "TIP-EXE: A Software Tool for Studying the Use and Understanding of Procedural Documents," *IEEE Trans Prof Commun*, vol. 55, no. 2, pp. 106–121, Jun. 2012, doi: 10.1109/TPC.2012.2194600.
- [2] W. Yih, M.-W. Chang, X. He, and I. Gao, "Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Stroudsburg, PA, USA: Association for Computational Linguistics, 2015, pp. 1321–1331. doi: doi.org/10.3115/v1/P15-1128.
- [3] Y. Hao *et al.*, "An End-to-End Model for Question Answering over Knowledge Base with Cross-Attention Combining Global Knowledge," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Stroudsburg, PA, USA: Association for Computational Linguistics, 2017, pp. 221–231. doi: 10.18653/v1/P17-1021.
- [4] P. Lewis *et al.*, "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," *NIPS'20: Proceedings of the 34th International Conference on Neural Information Processing Systems*, vol. abs/2005.11401, pp. 9459–9474, May 2020, doi: 10.48550/arXiv.2005.11401.
- [5] S. Sriwardhana, R. Weerasekera, E. Wen, T. Kaluarachchi, R. Rana, and S. Nanayakkara, "Improving the Domain Adaptation of Retrieval Augmented Generation (RAG) Models for Open Domain Question Answering," *Trans Assoc Comput Linguist*, vol. 11, pp. 1–17, 2023, doi: 10.1162/tacl_a_00530.
- [6] Y. Ahn, S.-G. Lee, J. Shim, and J. Park, "Retrieval-Augmented Response Generation for Knowledge-Grounded Conversation in the Wild," *IEEE Access*, vol. 10, pp. 131374–131385, 2022, doi: 10.1109/ACCESS.2022.3228964.
- [7] C. Xiong, S. Merity, and R. Socher, "Dynamic Memory Networks for Visual and Textual Question Answering," *Proceedings of The 33rd International Conference on Machine Learning*, pp. 2397–2406, Mar. 2016, doi: 10.48550/arXiv.1603.01417.
- [8] M. Seo, A. Kembhavi, A. Farhadi, and H. Hajishirzi, "Bidirectional Attention Flow for Machine Comprehension," *International Conference on Learning Representations*, Nov. 2016, doi: 10.48550/arXiv.1611.01603.
- [9] A. W. Yu *et al.*, "QANet: Combining Local Convolution with Global Self-Attention for Reading Comprehension," *International Conference on Learning Representations*, vol. abs/1804.09541, Apr. 2018, doi: 10.48550/arXiv.1804.09541.
- [10] W. Yang, Y. Xie, L. Tan, K. Xiong, M. Li, and J. Lin, "Data Augmentation for BERT Fine-Tuning in Open-Domain Question Answering," *ArXiv*, vol. abs/1904.06652, Apr. 2019, doi: 10.48550/arXiv.1904.06652.
- [11] C. Raffel *et al.*, "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer," *Journal of Machine Learning Research*, vol. 21, pp. 140:1-140:67, 2019, doi: 10.48550/arXiv.1910.10683.
- [12] A. Roberts, C. Raffel, and N. Shazeer, "How Much Knowledge Can You Pack Into the Parameters of a Language Model?," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 5418–5426. doi: 10.18653/v1/2020.emnlp-main.437.
- [13] M. T. R. Laskar, M. S. Bari, M. Rahman, M. A. H. Bhuiyan, S. R. Joty, and J. Huang, "A Systematic Study and Comprehensive Evaluation of ChatGPT on Benchmark Datasets," in *Annual Meeting of the Association for Computational Linguistics*, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:258967462>
- [14] W. Yu, "Retrieval-augmented Generation across Heterogeneous Knowledge," in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Student Research Workshop*, Seattle, Washington: Association for Computational Linguistics, Jul. 2022, pp. 52–58. doi: 10.18653/v1/2022.naacl-srw.7.
- [15] D. Thulke, N. Daheim, C. Dugast, and H. Ney, "Efficient Retrieval Augmented Generation from Unstructured Knowledge for Task-Oriented Dialog," *Conference of Association for the Advancement of Artificial Intelligence (AAAI)*, Feb. 2021, doi: 10.48550/arXiv.2102.04643.
- [16] OpenAI, "A Survey of Techniques for Maximizing LLM Performance." Nov. 2023.
- [17] Jacob Lee, "Building LLM-Powered Web Apps with Client-Side Technology." Accessed: Dec. 01, 2023. [Online]. Available: <https://ollama.ai/blog/building-llm-powered-web-apps>
- [18] J. Johnson, M. Douze, and H. Jégou, "Billion-Scale Similarity Search with GPUs," *IEEE Trans Big Data*, vol. 7, no. 3, pp. 535–547, 2021, doi: 10.1109/TBDATA.2019.2921572.
- [19] J. Zhu, J. Jang-Jaccard, I. Welch, H. Al-Sahaf, and S. Camtepe, *A Ransomware Triage Approach using a Task Memory based on Meta-Transfer Learning Framework*. 2022. doi: 10.48550/arXiv.2207.10242.
- [20] K. M. Fitria, "DocuQA: Document Question Answering Dataset." Feb. 2024. doi: 10.6084/m9.figshare.25223990.v1.
- [21] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "SQuAD: 100,000+ Questions for Machine Comprehension of Text," in

- Conference on Empirical Methods in Natural Language Processing, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:11816014>
- [22] A. Chen, G. Stanovsky, S. Singh, and M. Gardner, "Evaluating Question Answering Evaluation," in *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, A. Fisch, A. Talmor, R. Jia, M. Seo, E. Choi, and D. Chen, Eds., Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 119–124. doi: 10.18653/v1/D19-5817.
- [23] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi, "BERTScore: Evaluating Text Generation with BERT," *International Conference on Learning Representations*, vol. abs/1904.09675, Apr. 2019, doi: 10.48550/arXiv.1904.09675.
- [24] B. Ojokoh and E. Adebisi, "A Review of Question Answering Systems," *Journal of Web Engineering*, vol. 17, no. 8, pp. 717–758, 2019, doi: 10.13052/jwe1540-9589.1785.
- [25] J. Soni, N. Prabakar, and H. Upadhyay, "Behavioral Analysis of System Call Sequences Using LSTM Seq-Seq, Cosine Similarity and Jaccard Similarity for Real-Time Anomaly Detection," in *2019 International Conference on Computational Science and Computational Intelligence (CSCI)*, IEEE, Dec. 2019, pp. 214–219. doi: 10.1109/CSCI49370.2019.00043.
- [26] J. F. BELL and A. H. FIELDING, "A review of methods for the assessment of prediction errors in conservation presence/absence models," *Environ Conserv*, vol. 24, no. 1, pp. 38–49, 1997, doi: DOI: 10.1017/S0376892997000088.
- [27] C.-Y. Lin, "ROUGE: A Package for Automatic Evaluation of Summaries," *Association for Computational Linguistics*, vol. Text Summa, no. 12, pp. 74–81, 2004, [Online]. Available: <https://aclanthology.org/W04-1013/>
- [28] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a Method for Automatic Evaluation of Machine Translation," in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, P. Isabelle, E. Charniak, and D. Lin, Eds., Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, Jul. 2002, pp. 311–318. doi: 10.3115/1073083.1073135.
- [29] N. C. Chung, B. Miasojedow, M. Startek, and A. Gambin, "Jaccard/Tanimoto similarity test and estimation methods for biological presence-absence data," *BMC Bioinformatics*, vol. 20, no. 15, p. 644, 2019, doi: 10.1186/s12859-019-3118-5.
- [30] A. Stricker, "Question answering in Natural Language: the Special Case of Temporal Expressions," in *Proceedings of the Student Research Workshop Associated with RANLP 2021*, S. Djabri, D. Gimadi, T. Mihaylova, and I. Nikolova-Koleva, Eds., Online: INCOMA Ltd., Sep. 2021, pp. 184–192. [Online]. Available: <https://aclanthology.org/2021.ranlp-srw.26>
- [31] S. Min, V. Zhong, R. Socher, and C. Xiong, "Efficient and Robust Question Answering from Minimal Context over Documents," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, I. Gurevych and Y. Miyao, Eds., Melbourne, Australia: Association for Computational Linguistics, Jul. 2018, pp. 1725–1735. doi: 10.18653/v1/P18-1160.
- [32] H. Bahak, F. Taheri, Z. Zojaji, and A. Kazemi, "Evaluating ChatGPT as a Question Answering System: A Comprehensive Analysis and Comparison with Existing Models," *ArXiv*, vol. abs/2312.07592, Dec. 2023, doi: 10.48550/arXiv.2312.07592.
- [33] T. Cakaloglu, C. Szegedy, and X. Xu, "Text Embeddings for Retrieval From a Large Knowledge Base," *Research Challenges in Information Science*, vol. abs/1810.10176, Oct. 2018, doi: 10.48550/arXiv.1810.10176.
- [34] S. Gholami and M. Noori, "Zero-Shot Open-Book Question Answering," *ArXiv*, vol. abs/2111.11520, Nov. 2021, doi: doi.org/10.48550/arXiv.2111.11520.
- [35] G. Nur Ahmad and A. Romadhony, "End-to-End Question Answering System for Indonesian Documents Using TF-IDF and IndoBERT," in *2023 10th International Conference on Advanced Informatics: Concept, Theory and Application (ICAICTA)*, 2023, pp. 1–6. doi: 10.1109/ICAICTA59291.2023.10390111.

ORIGINALITY REPORT

10%

SIMILARITY INDEX

7%

INTERNET SOURCES

4%

PUBLICATIONS

4%

STUDENT PAPERS

PRIMARY SOURCES

1	Submitted to Stefan cel Mare University of Suceava Student Paper	2%
2	export.arxiv.org Internet Source	2%
3	Jiaxin Zhang, Zhilin Yu, Yunqin Li, Xueqiang Wang. "Uncovering Bias in Objective Mapping and Subjective Perception of Urban Building Functionality: A Machine Learning Approach to Urban Spatial Perception", Land, 2023 Publication	1%
4	"Pattern Recognition and Image Analysis", Springer Science and Business Media LLC, 2019 Publication	<1%
5	oaji.net Internet Source	<1%
6	ipfs.io Internet Source	<1%
7	www.iieta.org Internet Source	

<1 %

8

Submitted to University of Bristol

Student Paper

<1 %

9

Yongrui Chen, Huiying Li. "DAM: Transformer-based relation detection for Question Answering over Knowledge Base", Knowledge-Based Systems, 2020

Publication

<1 %

10

www.se.cuhk.edu.hk

Internet Source

<1 %

11

ts2.space

Internet Source

<1 %

12

Budhi, Gregorius Satia, and Rudy Adipranata. "Comparison of bidirectional associative memory, counterpropagation and evolutionary neural network for Java characters recognition", 2014 International Conference of Advanced Informatics Concept Theory and Application (ICAICTA), 2014.

Publication

<1 %

13

aircconline.com

Internet Source

<1 %

14

artemis.cslab.ece.ntua.gr:8080

Internet Source

<1 %

mdpi-res.com

15

Internet Source

<1 %

16

Ahmed, Toufique. "Learning Program Embedding From Unlabeled Source Code", University of California, Davis, 2023

Publication

<1 %

17

Mohammed Yousif Zeain, M. Abu, Z. Zakaria, Ahmed Jamal Abdullah Al-Gburi, R. Syahputri, A. Toding, Sriyanto Sriyanto. "Design of a wideband strip helical antenna for 5G applications", Bulletin of Electrical Engineering and Informatics, 2020

Publication

<1 %

18

Submitted to University College London

Student Paper

<1 %

19

Yasmina Al Ghadban, Huiqi (Yvonne) Lu, Uday Adavi, Ankita Sharma et al. "Transforming Healthcare Education: Harnessing Large Language Models for Frontline Health Worker Capacity Building using Retrieval-Augmented Generation", Cold Spring Harbor Laboratory, 2023

Publication

<1 %

20

ebin.pub

Internet Source

<1 %

21

www.mdpi.com

Internet Source

<1 %

22

d-nb.info

Internet Source

<1 %

23

dspace.ut.ee

Internet Source

<1 %

24

www.intechopen.com

Internet Source

<1 %

25

arxiv.org

Internet Source

<1 %

Exclude quotes Off

Exclude matches Off

Exclude bibliography On



Kaira Milani Fitria <kairaamilanii@gmail.com>

IJACSA March 2024: Paper Submission Received

1 message

Editor IJACSA <editorijacsa@thesai.org>

Tue, Feb 20, 2024 at 12:56 PM

To: kurnia@darmajaya.ac.id, kairaamilanii@gmail.com, joko.triloka@darmajaya.ac.id, sutedi@darmajaya.ac.id

Dear Corresponding Author,

Thank you for submitting your paper entitled:

1. "Retrieval-Augmented Generation Approach: Document Question Answering using Large Language Model"

for publication with International Journal of Advanced Computer Science and Applications (IJACSA) March 2024 Edition (Volume 15 No 3).

Your paper will be reviewed by the IJACSA technical committee and the evaluation outcome will be communicated up to 15 March 2024.

Regards,
Editor
IJACSA
The Science and Information (SAI) Organization

P.S. You can now rewatch the keynote talks from previous conferences available on our [Youtube channel](#). Press play and get inspired!



Kaira Milani Fitria <kairaamilanii@gmail.com>

IJACSA March 2024 : Reviewers Feedback

1 message

Editor IJACSA <editorijacsa@thesai.org>

Sat, Mar 23, 2024 at 10:51 PM

To: kurnia@darmajaya.ac.id, kairaamilanii@gmail.com, joko.triloka@darmajaya.ac.id, sutedi@darmajaya.ac.id

Dear Author,

Please find the attached Reviewer Feedback of your manuscript "Retrieval-Augmented Generation Approach: Document Question Answering using Large Language Model".

Kindly revise your paper as per the feedback attached here with and send us an updated version following the SAI Paper format (attached). Please submit your camera ready paper (both .docx and .pdf format) on or before March 26, 2024 for publication in IJACSA March 2024.

Tentative Publication Date - 30 March 2024

If you have prepared your paper in Latex, there is no need to submit a .docx file (Submit Latex sources with .pdf file). You may download the Latex Paper Format from <http://thesai.org/Home/Downloads>

Our publication team is experienced in handling most of the formatting issues in the manuscripts. While there are instances when an issue cannot be resolved, only in those cases the manuscript may be shifted to the next issue. There will be no other extra charges nor there will be any liabilities. We are fully committed to the satisfaction of the authors and are always there to assist you in the best possible manner.

Thank you for considering IJACSA as a medium for publication of your work.

Regards,
Editor
IJACSA
The Science and Information (SAI) Organization

P.S. You can now rewatch the keynote talks from previous conferences available on our [Youtube channel](#). Press play and get inspired!

3 attachments**SAI_PAPER_FORMAT.docx**
38K**Reviewer Feedback Form_2.pdf**
24K**Reviewer Feedback Form_1.pdf**
25K

Reviewer Feedback Form

Date March 7, 2024

International Journal of Advanced Computer Science and Applications (IJACSA)

Paper Title

Retrieval-Augmented Generation Approach: Document Question Answering using Large Language Model

Reviewer Recommendation

Accept with modifications

Reviewer Ratings

The authors contribution to the paper	Fair
Potential interest to research community	Good
Originality of the work	Fair
Use of examples and illustrations	Good
Quality of questions or problems raised by the Author	Fair
Reader's confidence in Author's knowledge	Fair
Formatting and Presentation	Fair
Awareness of related work	Good
Scientific Impact or Practical Utility	Good
Citations and References	Fair
Paper Organization	Good

Detailed Comments:

The 'Introduction' section needs to be elaborated with more details.

The previous work is adequately reported.

Details of the dataset will improve the quality of this paper. Based on the evaluation done for this work, the results are sufficient but more discussion on results and evaluation of different dataset is required to prove the scalability of this proposed model.

It is advisable to separate the 'Results and Discussion' sections as their combination may cause ambiguity in discerning the boundary between the description of the results and the start of the discussion.

The essential sections of the paper have been included, but it would be beneficial to add a separate section for the 'Conclusion' and 'Future Work'.

Reference [29] is not cited in the text.

A round of proof reading is recommended.

Grammar, punctuation, or spelling errors: