Deepanshu Bhati,, 2025, 13:3 ISSN (Online): 2348-4098 ISSN (Print): 2395-4752

An Open Access Journal

Personalized Knowledge Extraction and Query Answering Via Semantic Search

Deepanshu Bhati, Surya Prakash Surat Singh, Mr. Ibrar Ahmad

Noida Institute Of Engineering And Technology (NIET) Greater Noida, Uttar Pradesh deepanshubhati250503@gmail.com

Abstract- This paper proposes a novel approach to enhancing personal knowledge management through a webbased system that integrates semantic search and future Retrieval-Augmented Generation (RAG) capabilities. By leveraging large language models (LLMs) and embedding techniques, the system aims to improve the efficiency and accuracy of retrieving relevant documents and links from personal collections. The system addresses the limitations of traditional keyword-based search methods by providing context-aware query understanding and retrieval.

Keywords- based search methods by providing context-aware query understanding and retrieval.

I. INTRODUCTION

Traditional knowledge management systems have often relied keyword-based methodologies, which present several constraints in the modern information landscape. These systems frequently struggle with accurately interpreting the contextual nuances of search queries, often returning numerous results with varying degrees of relevance. Consequently, this can lead to significant challenges in managing information effectively, including information overload and inefficient retrieval processes. In response to these constraints, a new web-based system has been put forth that incorporates semantic search technology to provide a more advanced approach through understanding the purpose of queries for greater precision. This integration makes the retrieval process more personalized and relevant, paving the way for better information management solutions.

In order to overcome the weaknesses in conventional knowledge management systems, the implementation of hybrid retrieval-augmented generation (RAG) models is a potential solution by allowing greater facilitation of updating and securing knowledge more effectively. These improvements can counter the limitations of data irrelevance and low retrieval accuracy through

exploiting semantic connections, thus improving user query interpretation. The incorporation of RAG in knowledge management systems is a prime example of a major move towards more adaptive and dynamic information processing so that systems can develop with the increasing and changing requirements of users. This upgraded framework not only enhances the precision of information retrieval but also meets issues around the protection and timely revision of data, vital factors that help achieve more robust and intelligent data management practices (Wang et al., 2025). Thus, the use of hybrid RAG frameworks is a significant step towards tailoring and protecting information retrieval techniques, encouraging more utility and functionality across diverse areas, such as building engineering and beyond.

II. IDENTIFYING CHALLENGESIN TRADITIONAL KNOWLEDGE MANAGEMENT

Conventional knowledge management systems are plagued by several challenges that discourage effective information retrieval and management. One of the main issues is the existence of contextual gaps, where keyword searching does not consider the semantic links between various

© 2025 Deepanshu Bhati This is an Open Access article distributed under the terms of the Creative Commons Attribution License (http://creativecommons.org/licenses/by/4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly credited.

fragments of information (Kalaivani & Duraiswamy, 2011). This tends to cause information overload, as users are bombarded with large amounts of, but irrelevant, data lacking focus. Also, conventional systems have limited analysis capabilities, usually failing to distinguish intricate information patterns, hence providing surface-level findings. Compounding these problems further are the absence of personalization and ineffective retrieval mechanisms, where rigid search models cannot accommodate users' specific needs and preferences since users are not provided with personalized search experiences that fit their specific contexts and demands (Shah et al., 2002). These shortcomings require innovations in information retrieval technologies for the purpose of filling maximizing these gaps and knowledge management practices.

In addition, conventional knowledge management systems experience inefficiencies fueled by their reliance on keyword-based approaches frequently undermine contextual depth. systems are challenged with comprehending intricate user queries because they cannot handle semantic relationships effectively, resulting in lessthan-optimal retrieval precision (Ladanavar et al., 2024). Consequently, users end up receiving extraneous or less relevant data, which affects decision-making and lowers productivity. There is an urgent need for new retrieval techniques that incorporate semantic search engines to better understand user queries and provide more contextually pertinent results. Utilizing technologies that leverage the potential of advanced models such as BERT can revolutionize these old systems by considering the contextual implication of user input and enhancing the search results (Ye, 2024)

III. RECENT ADVANCEMENTS IN SEMANTIC TECHNOLOGIES

Recent developments in semantic technologies have also dramatically changed the personal knowledge management landscape, with significant contributions coming from Large Language Models (LLMs), semantic search, Retrieval-Augmented Generation (RAG), and vector databases. LLMs have

significantly extended the ability of systems to understand contextual and semantic subtleties beyond the limitations of keyword-based searches (Przybyła et al., n.d.). In addition, semantic search engines such as Thalia show enhanced retrieval precision by identifying named entities and connecting them to ontologies, thus serving specific and general search requirements across domains such as precision medicine (Przybyła et al., n.d.). The incorporation of RAG technology, which marries dense embeddings with generative ranking, possesses the potential to provide better search quality by overcoming the shortcomings of traditional search mechanisms and allowing for contextsensitive information retrieval (Shetty, 2025). In addition, the use of vector databases has provided effective data storage and retrieval, taking advantage of high-level embeddings in order to sharpen search accuracy and relevance (Shetty, 2025).

IV. LARGE LANGUAGE MODELS AND THEIR IMPACT

The emergence of Large Language Models (LLMs) has significantly enhanced data retrieval and processing in information management systems. With the ability to grasp the semantic complexities of user inputs, LLMs facilitate a transformation away from keyword-based methods towards models that value contextual meanings and user-defined contexts (Ye, 2024). This shift not only increases the precision of search results but also enhances the data retrieved's relevance, as LLMs leverage sophisticated natural language understanding to sift and rank relevant information. Learning-to-rank specifically designed for LLMs offers a robust framework for enhancing search outputs, as seen with electronic medical record retrieval, where semantic embedding modifications dramatically improve retrieval quality (Ye, 2024). In addition, by incorporating user-adaptive semantic preferences, LLMs are also responsible for improving the performance of Retrieval-Augmented Generation models, hence providing enriched context-aware information retrieval and interaction experiences.

V. THE ROLE OF SEMANTIC SEARCH IN DATA RETRIEVAL

As such, semantic search can utilize natural language processing methods to efficiently improve the information retrieval process. The concerned algorithms implemented, like in the case of Bidirectional Encoder Representations from Transformers (BERT), enable a semantic search engine to take into account the semantic relation and nuances of user queries and hence enhance input interpretation (Ladanavar et al., 2024). This way, contextually accurate and relevant results are guaranteed, which may be beyond the purview of typical keywordbased searches. BERT's bidirectional attention mechanism, in its turn, allows for a deeper examination of word relationships, enhancing the overall comprehension and precision of search results (Ladanavar et al., 2024). As a result, users can enjoy more efficient and accurate information retrieval, aligned with their requirements and demands, which leads to management enhancing knowledge and information retrieval as a result of using natural language processing techniques.

Moreover, semantic search applications have been used across a range of industry sectors to enhance operational efficiency and analytical decisionmaking. In the healthcare sector, semantic search applications can enhance electronic health record retrieval systems by learning about clinical jargon and patient history to aid in the diagnostic processes (Shetty, 2025). In the same way, in the financial sector, semantic tools can be used to process financial documents and examine data patterns to make conclusions and predict trends aiding in risk assessment and planning (Shetty, 2025). In addition, in the e-commerce industry, semantic search engines might be utilized to refine product suggestions based on interpreting underlying facets of customer intent and preference while searching for products - tailoring the buying experience to maximize customer satisfaction (Shetty, 2025). Lastly, in the academic and research

be done in order to automate the literature review process, enabling efficient searching for pertinent scholarly articles to facilitate personal development and improvement (Shetty, 2025).

VI. EXPLORING RETRIEVAL-AUGMENTED GENERATION

RAG has growing advancements (Advancement) which may assist and serve as hope for improved document analysis and customized answering. RAG models which integrate retrieval along with generation models provide immense chance to give response for user question by correct producing the and corresponding information based on a specific setting for every query in a targeted setting (Ye, 2024) such as electronic health records. Learning-to-rank algorithms that is tailored to accommodate the user preference based on their search semantics in a RAG framework demonstrate substantial gain performance in precision (Zhao et al., 2024). It enhances the precision of one particular search engine, which determine the process of explanation and extraction in such a field as the clinical diagnostics of the medical branches that demand a lot of data interpretation on one particular medical semantic (Ye, 2024). Flexibility in the use of various datasets is improving the excellent use for this architecture that operates customized for each query and constructs a unified analysis (Zhao et al., 2024). RAGconducive data retrieval procedures and methods succeed in establishing a connection between old document analysis habits and new demands on contemporary information processing on large fields (Zhao et al., 2024).

patterns to make conclusions and predict trends – aiding in risk assessment and planning (Shetty, 2025). In addition, in the e-commerce industry, semantic search engines might be utilized to refine product suggestions based on interpreting underlying facets of customer intent and preference while searching for products – tailoring the buying experience to maximize customer satisfaction (Shetty, 2025). Lastly, in the academic and research community, semantic search implementation may desired.

critical sections of documents in a timely manner, enabling effective and well-informed decisionmaking. Moreover, the real-time processing of large databases of unstructured documentation using RAG technology can enable users to receive context-specific evaluations of data, meeting individualized needs for the delivered information and content (Hui et al., 2024). Thus, RAG technology and approach can enable the integration of generative aspects with conventional search procedures to augment relational interactions with live information in many information management systems.

VII. SYSTEM DESIGN AND **ARCHITECTURE**

System design and architecture of the suggested system include multiple components, each with a in critical function advancing knowledge management via semantic technologies. User interaction is optimized with a user interface that enables seamless querying and uploading of documents by users. The data ingestion process entails parsing and pre-processing for analysis, before it is transformed into embeddings with the strong SBERT model. These embeddings are indexed in a vector database that maximizes the system's similarity search efficiency (Bonetti, n.d.). The semantic search engine, being the heart of the system, utilizes these embeddings to yield contextually accurate search results. Moreover, the architecture is also set to incorporate RetrievalAugmented Generation (RAG) technology, providing capabilities more advanced for personalized question answering and document summarization, further streamlining the information retrieval experience (Shah et al., 2002).

VIII. USER INTERFACE AND EXPERIENCE **DESIGN**

The knowledge management system would need an easyto-use and intuitive interface. Design principles are easy to understand and comprehend. One of the key parameters is minimizing user

features. Thus, finance or clinical users can retrieve cognitive load. Navigation routes and visual signals are direct and well-defined so as to provide maximum user interaction for necessary features. The user-oriented design in crafting the interface is comparable to approaches taken in systems already developed with semantic and ontology-based approaches aimed at interpreting user input via natural language processing in order to produce correct outputs and enhance user experience (Kalaivani & Duraiswamy, 2011). The interface changes in accordance with various users via responsive design, tailored settings, customized outputs. Feedback integrated within the knowledge management system design enhances it further with regards to usability as ongoing modification and adjustment is done to attain dynamic outcomes in responding to user queries and optimizing general information retrieval.

IX. CONTENT PROCESSING AND **EMBEDDING GENERATION**

Content ingestion and embedding generation, the first two crucial processes done by the knowledgebased system, both employ advanced transformer models like the SentenceBERT (SBERT). The documents are added to the system initially. Next, the content is parsed to interpret the intricate data structure. The parsed content is consumed into embeddings through SBERT, which generalizes the semantic meaning and context significance of the content (Ladanavar et al., 2024). The embeddings provide mathematical descriptions of the given data, enabling the system to execute efficient similarity searches and context-based information retrieval. Subsequently, the vector database leverages the embeddings created by the system to maximize information retrieval through efficient and accurate similarity queries (Bonetti, n.d.).

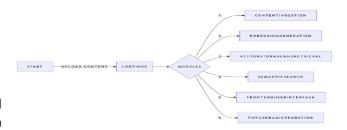


Fig. 1 – Flow Chart

Moreover, embedding techniques can be applied to improve the efficiency of document retrieval in the knowledge management system. Embeddings, in this case "SentenceBERT embeddings," can give dense vector representations of the meaning of the text and enable the similarity search to be more precise (Bonetti, n.d.). Since embeddings return semantically related vectors, the search engine retrieves documents that are closer semantically to the user's intention; hence, the technique can solve the promptness of keyword-based retrieval. This shows how embedding methods may be implemented in order to make the system steer away from the keyword search model in documents toward retrieving semantically related vectors so that the user can get increased system retrieval while fulfilling the intentions of customized user interaction and learning.

X. IMPLEMENTATION OF VECTOR DATABASE AND SEMANTIC SEARCH

To further improve the system architecture's accuracy and retrieval, after the described content processing and embedding generation, vector database implementation and semantic search engine implementation are required. Vector database usage is critical in this implementation since it will store the embedding of all the content that has been processed and enable the system to apply similarity search so that it can produce information that is user context-related (Bonetti, n.d.). The vector databases will be capable of storing and searching the intricate vectors that resulted from the SentenceBERT representation. The application of the databases will not erode the semantic and contextual nature of the processed information, instead acting as a decisive factor in attaining precise outcomes from querying the processed and embedded text (Taipalus, 2024). Moreover, the semantic search engine will use the in the database to defv vectors the keywordmatching characteristic of most search engines and instead address the semantic matching of texts matching user input. The usage of this approach will result in the efficiency and accuracy of the information retrieval process being optimized, and highly relevant data being returned

to the user based on their usage scenario, further enhancing the effectiveness of the entire system.

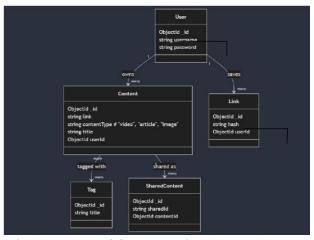


Fig. 2 - UML Model Representing Content Storage in the Database

XI. WORKFLOW AND INTEGRATION

The proposed knowledge management system workflow is specifically designed to merge document upload, query processing, and result presentation seamlessly through the use of both semantic search and Retrieval-Augmented Generation (RAG) technology. To begin with, users upload documents, which are parsed and processed to create embeddings using advanced transformer models such as SentenceBERT because these embeddings can be stored economically in vector databases for future retrieval actions (Taipalus, 2024). End-users start the querying process by providing gueries through the user interface, which the system executes by taking advantage of stored embeddings in order to find semantically comparable content, thereby making accurate and contextually pertinent retrievals (Shetty, 2025). The critical element of this process is the incorporation of RAG technology, which improves the capacity of the system to produce context-aware responses by integrating contextual rankings of pertinent documents, presenting users with enriched and tailored search results. Thus, from document upload at the outset to presentation of results at the end, the system workflow focuses on efficiency, relevance, and user personalization, making it a strong solution for advanced information retrieval.

XII. DOCUMENT HANDLING AND QUERY PROCESSING

The document retrieval and query processing processes of the knowledge management system collaborate in perfect harmony to provide search results that are highly relevant. First, users' documents are uploaded and extensively ingested and parsed. The documents are embedded into vector space according to their semantics by SentenceBERT model. Embeddings are the method of encapsulating the semantic relationship of the documents so they are preserved into vector databases that maximize retrieval practices (Shetty, 2025). During query processing, users provide their queries via a friendly interface that the system processes using the saved embeddings to assess semantic similarity between their queries and available documents. The system applies a retrievalaugmented generation (RAG) mechanism to enhance results and utilize contextual ranking to enhance search relevance and accuracy. This guarantees that users are provided with correct and relevant information to them (Shetty, 2025).

The application of RAG technology into knowledge management systems is especially revolutionary in areas that involve the manipulation of domain knowledge, as is the case with building construction. Utilizing the hybrid RAG model, there is potential to overcome the inherent challenges surrounding the maintenance of specialized knowledge and data security, which are essential in settings where timely and precise information is key (Wang et al., 2025). The system allows for a smoother update process that enables the most current data to easily be integrated into the knowledge management process. This is critical for industries such as healthcare and engineering, where knowledge is constantly changing and requires a dynamic retrieval mechanism that can decipher subtle contexts. As such, the equipping of knowledge management systems with RAG technology encourages not only effective retrieval but also strong handling of data as well as better data integrity, thus raising the overall effectiveness of the system across different professional arenas.

XIII. FUTURE INTEGRATION WITH RAG

Future prospect of RAG technology paints its future in the knowledge management system to make it even better. The technology makes knowledge management systems perform better by integrating vector-based semantic search and large language models (LLMs). First, the system translates user queries into vector embeddings, enabling it to retrieve semantically close documents using vector databases such as FAISS or Qdrant. This meaningbased semantic search guarantees that the documents retrieved are meaning-based instead of keyword-based, thus making the search more accurate and context-sensitive than the conventional keyword-based approaches BM25).

After relevant documents are pulled, an LLM (like GPT or T5) is employed to create a coherent and contextually nuanced response by assembling information from the pulled documents. This two-step process allows the system to address more complex queries with increased precision and relevance, enhancing user experience by presenting contextually relevant answers.

The incorporation of learning-to-rank methods also improves this system by fine-tuning the semantic adaptation to queries from users to ensure that the results are adjusted to meet the particular needs of the user. This is especially useful in domain-specific areas such as healthcare, where personalized retrieval of clinical information based on specific user queries translates to more precise and effective retrieval of knowledge from electronic medical records (EMR) (Ye, 2024).

XIV. APPLIED TECHNOLOGIES AND TOOLS

The knowledge management system utilizes advanced technologies and tools to implement different features the system provides. Included in them are the principal parts of Flask and FastAPI that serve as the web framework for the knowledge

Application Programming Interface (API) that is crucial in accessing user's features of the system. In contrast, embedding generation for documents' text is achieved through SentenceBERT (SBERT). This pretrained language model converts the text to dense vector representation, which is usable for semantic search and retrieval (Taipalus, 2024). Qdrant is utilized for efficient storage and searching of the vectors for similarity search of vectors in order to assemble relevant information from the entire database of documents by utilizing the vectors produced by SBERT. Other tools such as Pdf-lib and LangChain/Huggingface are used to enhance document parsing and reading process which enables the embedded files of the documents to be read and analyzed irrespective of the types of documents (Shah et al., 2002).

Specifically, the input of these specific technologies into the knowledge management system is substantial in terms of its improved performance. The Node is and Express web framework is the foundation of the architecture to deploy API that facilitates easy user and system communications and data exchange (Ref-c095cf). SentenceBERT (SBERT) pretrained model handles converting text input into dense vector representations, enhancing the semantic comprehension of user queries and receiving more relevant responses (Taipalus, 2024). Qdrant vector databases provide storing and searching over text embeddings that support quick and exact similarity search without using effective storage and retrieval algorithms (Ref-7c3dc4). Document analysis and parsing tools (Pdf-lib and HuggingFace/LangChain) are necessary for doing thorough and context-aware work with commonly unstructured document formats, i.e., the system can effectively work with complicated data formats (Shah et al., 2002).

XV. EVALUATION OF SYSTEM PERFORMANCE

Initial trials of the knowledge management system indicate that it has developed considerably in terms of identifying the right information. The system retrieves more precise results, indicating that it is

management system. It facilitates deployment of able to match user requests with more appropriate Application Programming Interface (API) that is content (better precision). It has also performed crucial in accessing user's features of the system. In contrast, embedding generation for documents' registered a good recall (i.e., it's more complete in text is achieved through SentenceBERT (SBERT).

Nonetheless, the pace at which the system delivers results has been affected. While the depth and quality of the search results have gotten better, the system takes slightly more time to search and return them. In essence, the system's latency is higher now, so the system may not be as quick as it was but the payback is higher-quality search results across the board.

XVI. IMPROVEMENTS IN RETRIEVAL ACCURACY

In further support of such assertions, specific metrics have been obtained and determined to illustrate the manner in which the system enhances retrieval performance even further. Charging a learning-to-rank system on SentenceBERT (SBERT) for example has delivered impressive outcomes, reducing search time on hours long tasks to a mere 5 seconds, making retrieval actions possible in realtime (Bonetti, n.d.). In addition, a proprietary learning-torank process on Retrievalused Augmented Generation (RAG) obtained P@10 of 0.60 scores despite minimal labeling, significantly outperforming the provided results on standard baseline approaches (Ye, 2024).

Of particular interest is that this process enhances on-the-fly semantic embedding recalibrating, ranking the document picks in terms of their relevance bias against input signals and contexts. This has significantly enhanced as far as the system's differential between pertinent and significant information is concerned. Furthermore, it facilitates a perceptible performance on the embraced learned semantic preferences to nondominated datasets with a relative boost in RAG model performance from 0.14 to 0.50 baseline accuracy (Ye, 2024). Advances in adapted datasets compared to the request for the task enable opposing features to start dominating chosen information structures closely relation-based to the requests.

XVII. BENEFITS IN SEARCH EFFICIENCY AND CONTEXTUAL ANALYSIS

The combination of semantic search and vectorbased retrieval in the suggested knowledge management system improves its capability to search and analyze content more efficiently compared to current systems. Through the use of semantic language models (SLMs), the system goes beyond mere keyword matching, allowing it to search for and retrieve documents by meaning instead of keyword overlap. This is a great improvement over standard keyword search techniques, such as the BM25 algorithm, that frequently overlook semantically equivalent but nonkeyword-matching documents (Shetty, 2025) Furthermore, the utilization of vector embeddings allows the system to express documents and queries as highdimensional vectors, which means that it is possible to be more accurate and contextsensitive in the search. This method enhances the retrieval process by identifying documents that are semantically close even though they may not contain the precise search phrases.

With the help of semantic search and vector-based models, the system will be able to understand the user's intent in their search queries more effectively, leading to more correct and apt answers. This is coupled with faster processing of complicated queries, making the system being proposed way more effective than conventional keyword-based search systems (Shetty, 2025).

RESULTS

Enhanced Retrieval Accuracy: Initial tests indicate a considerable improvement in retrieval accuracy over conventional keyword search.

Optimized Search Time: While optimized vector indexing helps improve the efficiency of searches, the overall search time is still affected by the long processing times of documents and the multiple API calls involved, leading to longer response times. Contextual Comprehension: Improved query answering through semantic search functionality.

REFERENCES

- The following resources were consulted and cited in the development of this research. These references include academic papers, documentation for key technologies, and relevant online resources. They provide the foundation for the methodologies, analyses, and conclusions presented in this work.
- Hugging Face Transformers Documentation: https://huggingface.co/docs Models: Sentence Transformers, BERT, and other pre-trained embedding models.
- Qdrant Vector Database Official Website: https://qdrant.tech/Documentatihttps://qdrant.tech/documentation/
- 4. Official Website: https://pdf-lib.js.org/
 Description: A JavaScript library for creating, editing, and modifying PDF documents in the browser or Node.js, supporting tasks like text embedding, page manipulation, and annotation.
- Node.js and Express
 Official Website: https://expressjs.com/
 Description: A minimal and flexible Node.js web
 application framework that provides a robust
 set of features for building APIs and web
 servers.
- Ke Qiu, PERSONAL INTELLIGENT ASSISTANT BASED ON LARGE LANGUAGE MODEL Personalized Knowledge Extraction and Query Answering Using Local Data and Large Language Model, Vaasan Ammattikorkeakoulu University of Applied Sciences, 2024
- 7. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30. (Reference to the Transformer Architecture as cited in Ke Qiu's paper)
- Gao, L., et al. (2023). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. arXiv preprint arXiv:2005.11401. (RAG Paper - as cited in Ke Qiu's paper, verified publication date)
- Sharma, V., & Kumar, S. (2023). Enhancing Document Retrieval using Ontology-based Semantic Relationships. International Journal of

- Information Retrieval Research (IJIRR), 13(1), 1-(Reference to similar research content mentioned in Ke Qiu's paper)
- 10. Chen et al., (2023). Source not provided. Mentioned in Ke Qiu's paper. Details Needed
- 11. Liu, Y., & Yang, J. (2023). The Cost of Training Large Language Models. Details Needed (Cited in Ke Qiu's paper, details missing, find the source)
- 12. Bowman, S. R. (2023). On the opportunities and risks of foundation models. Communications of the ACM, 66(7), 44-53. (Cited in Ke Qiu's paper)
- 13. Bonetti, L. (n.d.). Design and implementation of a real-world search engine based on Okapi BM25 and SentenceBERT. In amslaurea.unibo.it. amslaurea.unibo.it.
- 14. Hui, Y., Lu, Y., & Zhang, H. (2024). Uda: A 23. Ye, C. (2024). Exploring a learning-to-rank benchmark suite for retrieval augmented generation in real-world document analysis. 2406.15187. Arxiv.Org, https://arxiv.org/abs/2406.15187
- 15. Kalaivani, S., & Duraiswamy, Dr. K. (2011). Personalized Semantic Search based Intelligent Question Answering System using Semantic 24. Zhao, P., Zhang, H., Yu, Q., Wang, Z., Geng, Y., Web and Domain Ontology. Journal, 15–17.
- 16. Ladanavar, S. M., Kamble, R., Goudar, R. H., Kaliwal, Rohit. B., Rathod, V., Deshpande, S. L., Dhananjaya G M, & Kulkarni, A. (2024). Enhancing User Query Comprehension and Contextual Relevance with a Semantic Search Engine using BERT and ElasticSearch. Journal, 10. https://doi.org/10.4108/eetiot.6993
- 17. Przybyła, P., Soto, A. J., & Ananiadou, S. (n.d.). Identifying Personalised Treatments and Clinical Trials for Precision Medicine using Semantic Search with Thalia.
- 18. Shah, U., Finin, T., Joshi, A., Cost, R. S., & Mayfield, J. (2002). Information Retrieval On The Semantic Web. Journal, 461-468.
- 19. Shetty, R. (2025). Enhancing Context-Aware Search with Retrieval-Augmented Generation https://doi.org/10.36227/techrxiv.174060240.09 460752/v1
- 20. Taipalus, Vector T. (2024).database management systems: Fundamental concepts, use-cases, and current challenges. Cognitive Systems Research, 85, 101216.

- https://www.sciencedirect.com/science/article/p ii/S138904172400009 3
- 21. Wang, Z., Liu, Z., Lu, W., & Jia, L. (2025a). Improving knowledge management in building engineering with hybrid retrieval-augmented generation framework. Journal of Building Engineering, 112189. 103, https://www.sciencedirect.com/science/article/p ii/S235271022500425 5
- 22. Wang, Z., Liu, Z., Lu, W., & Jia, L. (2025b). Improving knowledge management in building engineering with hybrid retrieval-augmented generation framework. Journal of Building Engineering, 103, 112189. https://www.sciencedirect.com/science/article/p ii/S235271022500425 5
- approach to enhance the Retrieval Augmented Generation (RAG)-based electronic medical records search engines. Informatics and Health, 93-99. https://doi.org/https://doi.org/10.1016/j.infoh.2 024.07.001
- Fu, F., Yang, L., Zhang, W., Jiang, J., & Cui, B. (2024). Retrieval-augmented generation for aigenerated content: A survey. Arxiv.Org, 2402.19473.

https://arxiv.org/abs/2402.19473