# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

# AI-Powered Conversational Agent for Ground water Monitoring and Knowledge Sharing

Dr. B. Rajakumar[1], Jagadeesh.A[2], Balaji.N[3], Diwakar.R[4], Santhosh.S[5]
*Department of Artificial Intelligence and Data Science, J.N.N institute of Engineering, Kannigaipair, Tiruvallur, Tamilnadu, India*

*Abstract: Groundwaterisone ofthemostcriticalnaturalresources,supplying a significant portion of drinking water and irrigation needs world wide.However,rapiddepletionduetoover-extraction,climate change,andpollutionhasled tosevere water crisesinmany regions. Effective groundwater monitoring and management require advanced technological solutions to ensure sustainability. This research introduces an AI-powered chatbot that functions as an intelligent systemforcollating, analyzing,and disseminating real- time groundwater information. The proposed chatbot leverages Natural Language Processing (NLP) and Machine Learning (ML) techniques to interpret user queries, retrieve relevant groundwater data, and provide insightful responses. Utilizing deep learning models such as Sentence Transformers for NLP-based query handling and Convolutional Neural Networks (CNNs) for image- baseddataanalysis,thechatbotensuresaccuracyinunderstanding groundwater patterns and trends.*
*Keywords: Groundwater Monitoring, AI Chatbot, Natural Language Processing, Machine Learning, Water Conservation, IoT-based Groundwater Analysis.*

## I. INTRODUCTION

Groundwaterisoneofthemostessentialnaturalresources,providing nearly 30% of the world's freshwater supply and supporting agricultural, industrial, and domestic needs. However, rapid depletion duetoover-extraction,pollution,andinefficientmanagementhasled to severe groundwater crises globally.Many regions suffer from declining watertables, reduced aquiferrecharge, and contamination, posingsignificantriskstobothhumanpopulationsandecosystems. Addressing these challenges requires effective groundwater monitoring and data-driven decision- making.Traditional groundwater monitoring methods rely on manual data collection, periodic government reports, and observational studies. These approaches, while informative, often lack real-time data accessibility, leading to delays in decision-making and inefficient water resource management. Moreover, most conventional systems are unable to integratelarge-scalegroundwaterdatasetsorprovideinstantinsights to diverse stakeholders, including policymakers, researchers, farmers,andurbanplanners.ArtificialIntelligence(AI)andMachine Learning(ML)haveemergedastransformativetechnologiescapable of enhancing data accessibility and analysis. AI-powered chatbots, equipped with Natural Language Processing (NLP), offer an innovative solution to groundwater monitoring challenges by providing real-time insights through an interactive platform.

This study introduces an AI chatbot designed to collate and disseminategroundwater dataefficiently. By leveraging NLP, ML, and web-based information retrieval, the chatbot enables users to access groundwater information seamlessly, thereby promoting informed decision-making and sustainable water management practices. This paper aims to provide a thorough analysis of the utilization ofAIchatbots inmanagingundergroundwaterdata.Itwill cover the challenges in this field, the capabilities of AI chatbots, frameworks for implementation, and examples of successful applications.Challenges in Underground Water Resource Management One of the significant challenges in managing underground water resources is the lack of adequate data. In many regions, data on groundwater levels, quality, and usage is either scarce or poorly organized. This scarcity leads to fragmented knowledge, making it difficult for stakeholders to make informed decisions.Groundwater systemsareinherentlycomplex duetotheir geological and hydrological variability.

Understanding the dynamics of aquifers, recharge rates, and contamination requires expertise in various fields, including geology, hydrology, and environmental science. This complexity poses challenges in data interpretation and communication among stakeholders.

## II. LITERATURE SURVEYS

Aliteraturesurveyisacomprehensivereviewof existing research, studies,andscholarly articlesrelated to aspecific topic.Itprovidesanoverviewofpreviousfindings,identifiesresearch gaps,andestablishesafoundationfornewstudies.Inthefieldofweb scraping and AI-driven knowledge retrieval, various studies have explored techniques for efficient data extraction and processing.

Researchers have developed web scraping methods using tools like BeautifulSoup, Scrapy, and Selenium, enabling automated data collection from structured and unstructured web sources. Similarly, advancements in natural language processing (NLP) have led to the adoption of Sentence Transformers for generating semantic embeddings,improvinginformationretrievalaccuracy.Priorresearch onsemanticsearchmodels,includingBERTandSBERT,highlights the effectiveness of transformer-based embeddings in capturing contextualrelationships.Studiesalsoemphasizethechallengesofweb scraping, suchasethicalconcerns,websiterestrictions, and theneed for real-time data validation. AI-based chatbots and search engines have integrated precomputed embeddings with cosine similarity methods to enhance response generation. In the domain of groundwaterknowledgeretrieval,limitedworkhas beenconducted, withmostresearchfocusingontraditionaldatabase-drivenapproaches rather than AI-enhanced search models. The existing literature underscores the importance of hybrid systems that combine static knowledge bases with real-time web search to provide up-to-date, relevant responses. Future research should focus on scalable architectures, multilingual support, and adaptive learning models to further refine AI-driven groundwater knowledge systems, ensuring bothefficiencyandaccuracyinretrievingscientificandpolicy-related information.

Google Cloud API provides powerful tools for integrating search functionalityintoapplications,enablingefficientdataretrievalfrom theweb.Oneofthemostcommonlyusedservicesforthispurposeis theGoogleCustomSearchJSONAPI,whichallows developers to createcustomsearchenginestailoredtospecificdomainsortopics. This API helps automate query processing, indexing, and result filtering, making it useful for AI-driven search applications.

TousetheGoogleCustomSearchAPI,developersmustfirstcreatea Custom Search Engine (CSE) in the Google Programmable Search Engine console.Thisinvolvesspecifyingsearchparameters suchas domainrestrictions,rankingpreferences,andresultfiltering.Oncethe CSE is created, it generates a unique Search Engine ID (CX ID), which is required for API requests. The API key, obtained from Google Cloud Console, is also necessary for authentication and request processing.The API returns search results in JSON format, including titles, snippets, and URLs of relevant web pages. These results can be further processed using AI models, NLP techniques, or embedding-based similarity matching for enhanced information retrieval.Google Cloud API combined with Custom Search Engine is crucial for AI-driven applications requiring real-time web data, such as knowledge-based chatbots, research assistants, and environmental monitoring systems. Proper usage of API rate limits and adherence to Google's policies are essential for maintaining ethical and efficient data access.

## III. METHODOLOGIES

To develop an AI-driven groundwater knowledge retrieval system, a multi-layered methodology is employed, integrating web scraping, natural language processing (NLP), semantic search, and cloud-based APIs.

Eachmethodologicalcomponentismeticulouslyengineeredto ensure high accuracy, efficiency, and scalability.

1) *Data Acquisition and Preprocessing:* Thefirststageinvolves automateddataacquisitionthroughwebscrapingandstructured dataintegration.BeautifulSoup(BS4)andScrapyareleveraged to extract relevant groundwater information from scientific articles, government reports, and environmental databases. WebsitesimplementingJavaScript-renderedcontentnecessitate Selenium-based dynamic scraping, ensuring comprehensive data collection.

   Theextractedrawdataundergoestextnormalization,including tokenization,lemmatization,andstopwordremoval,torefineits structure. Data integrity is preserved using checksum algorithms, while Named Entity Recognition (NER) is employed to extract domain-specific terminology such as aquifers, pollutants, and conservation policies.A critical aspect is data deduplication and filtration through Jaccard similarity and TF-IDF weighting, ensuring redundancy elimination and context retention. Preprocessing pipelines utilize spacy and NLTK, reinforcing syntactic coherence and semantic alignment. The final dataset is stored in structured formats— JSON for static QA pairs and Pickle (PKL) for high-speed embeddings, ensuring efficient retrieval.

2) *Embedding Generation and Semantic Representation:* To facilitate semantic retrieval, the "sentence-transformers/all- mpnet-base-v2"modelisemployedforhigh-dimensionalvector encoding. Unlike traditional TF-IDF and BM25 models, Sentence Transformers leverage bi-directional self-attention mechanisms, encapsulating contextual embeddings with 768- dimensional latent representations.Each groundwater-related question undergoes transformer-based encoding, generating densevectorrepresentationsstoredinPyTorchtensorsforhigh- speed similarity computations. The embeddings are optimized using contrastive learning and triplet loss functions, enhancing discriminative representation in semantic space.To further compress embedding dimensions while retaining high variance capture, Principal Component Analysis (PCA) and Singular Value Decomposition (SVD) are applied.

This step reduces computational overhead without compromising retrieval efficacy.Stored embeddings are indexed using FAISS (Facebook AI Similarity Search), which enables approximate nearest neighbor (ANN) retrieval, ensuring sub-millisecond search latency. This indexing methodology significantly outperforms brute-force cosine similarity computations, making the system scalable for large- scale datasets.

3) Data Overview: The chatbot system for Collating and Dissemination of Underground Water relies on two primary data sources: predefined structured knowledge stored in JSON format and real-time data extracted via web scraping. This dual-source approach ensures that users receive accurate, structured, and dynamically updated information regarding groundwater conditions, aquifer health, pollution risks, and conservation strategies.

4) PredefinedKnowledgeBase (JSONData): ThestructuredJSON repository serves as the primary data source, containing well-defined groundwater-related concepts. This dataset is manually curated based on government reports, scientific publications, and hydrologicaldatabases,ensuringhighaccuracyandreliability.The JSON structure is categorized into multiple domains, including: Aquifer Types (Confined, Unconfined, Perched, Artesian) , GroundwaterContamination(Industrialwaste,Agriculturalrunoff, Heavymetalintrusion),RechargeTechniques(Artificialrecharge, Rainwater harvesting, Injection wells) , Water Conservation Methods (Greywater reuse, Sustainable irrigation, Desalination efforts) , Regulatory Frameworks (Water governance policies, Legal compliance, International treaties)

5) Real-Time WebScraping Data: Fordynamically evolvingsuch as groundwater levels, contamination alerts, policy updates, and environmentalreports,thechatbotutilizesautomatedwebscraping techniques to extract the latest data from: Government Portals (USGS, CGWB, UNEP-Water) , Scientific Repositories (Google Scholar, ScienceDirect, IEEE Xplore) , Environmental News Sources(NASAEarthObservatory,WorldWaterCouncil).BeautifulSoup (BS4)isa Python library used forwebscrapingby extracting data from HTML and XML documents.



WEB SCRAPING

HTML WEBSITES → WEB SCRAPING → DATA

## IV. AI EMBEDDING SUSINGS ENTENCE TRANSFORMERS

AI embeddings play a crucial role in natural language processing (NLP), enabling the transformation of textual data into meaningful numerical representations. Sentence Transformers (ST) have emerged as a powerful framework for generating sentence-level embeddings that capture semantic similarity and contextual meaning.Thispaperexploresthearchitecture,implementation,and applications of Sentence Transformers in various NLP tasks.

Traditional word embeddings, such as Word2Vec, GloVe, and FastText, represent words in a fixed-dimensional space based on theirco-occurrencepatterns.However,thesemodelsfailtocapture contextual variations and sentence-level meaning.

Transformer-based embeddings, such as BERT (Bidirectional EncoderRepresentationsfromTransformers)anditsderivatives, address this limitation by leveraging deep contextualized representations.

## V. SENTENCE TRANSFORMERS: OVERVIEW AND ARCHITECTURE

Sentence Transformers, introduced as an extension of BERT and other transformer models, are designed to generate high- quality sentenceembeddingsefficiently.The frame wor kprimarilyutilizes theSiameseandtripletnetworkstructurestocomputesemantically meaningful representations.

International Journal for Research in Applied Science & Engineering Technology (IJRASET)
*ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538*
*Volume 13 Issue IV Apr 2025- Available at www.ijraset.com*

1) *Base Models:* Sentence Transformers are built on pre- trained transformer models, such as BERT, RoBERTa, andDistilBERT. PoolingMechanisms:Sincetransformer outputs are token-level embeddings, Sentence Transformers employ pooling strategies such as mean pooling,maxpooling,orCLS tokenextractiontoderive sentence-level embeddings.

2) *Training Strategies:* Fine-tuning is performed using contrastivelearning,cosinesimilarityloss,andtripletloss, optimizingthemodelforsemantictextualsimilarity(STS) tasks.



3) *Computational Efficiency :* Unlike BERT-based models that requireextensivecomputationalresources,SentenceTransformers optimize inference speed by leveraging:Distillation techniques to create lightweight models.Quantization methods to reduce memory footprint.ONNX conversion for faster inference in production settings.Sentence Embeddings

a) *Information Retrieval*
Sentenceembeddingsenableefficientdocumentretrieval by encoding queries and documents inasharedvectorspace.ModelslikeSBERT improve searchrelevanceby replacing traditional BM25-based ranking algorithms.

b) *Semantic Search*
ApplicationssuchasGoogleSearchanddomain-specific knowledge bases leverage Sentence Transformers for improved semantic understanding and result ranking.

c) *Text Classification*
Sentenceembeddingsserveasrobustfeaturerepresentations for tasks like sentiment analysis, topic classification, and spam detection.

*d) Text Cluster in gand Summarization*

Bycapturingcontextualsimilarity,embeddingsfacilitate document clustering and extractive summarization techniques, improving NLP-driven analytics.
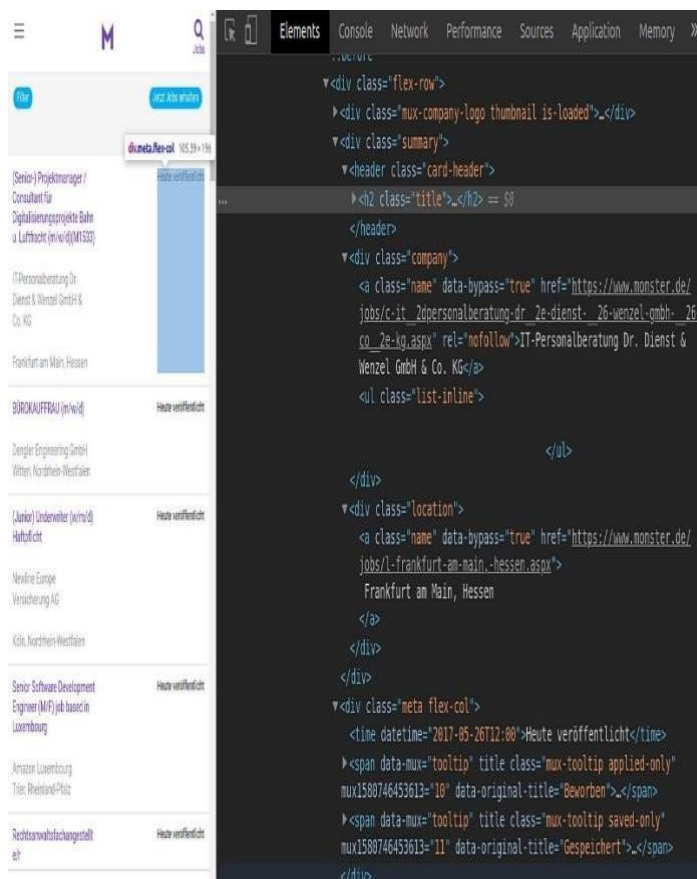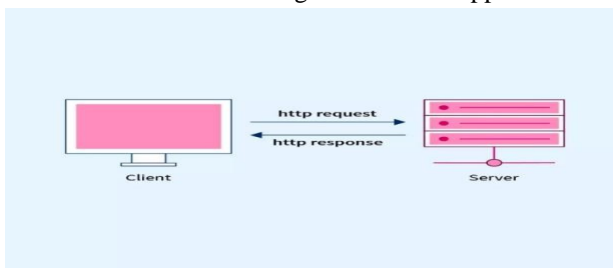
SentenceTransformersgenerateembeddingsbyencodingtext into high-dimensional vector spaces that preserve semantic meaning. Unlike traditional word embeddings such as Word2Vec or GloVe, which provide static representations, Sentence Transformers produce context-aware embeddings usingtransformer-basedarchitectureslikeBERT,RoBERTa, and DistilBERT.The embedding process involves tokenizing input text and passing it through a transformer model to obtain contextualized token-level representations. Since transformers output embeddings for each token, a pooling mechanism (e.g., mean pooling, max pooling, or CLS token extraction) aggregates these representations into a single sentence-level vector.Fine-tuning is performed using supervised contrastive learning approaches, including cosine similarity loss and triplet loss, optimizing embeddings for semantic similarity tasks. The resulting embeddings enable applicationssuchasinformationretrieval,textclustering,and semantic search by efficiently mapping similar sentences closer in the vector space.To enhance efficiency, Sentence Transformers leverage model distillation, quantization, and ONNX conversion, reducing inference time and computationalcosts.Despitetheiradvantages,challengeslike domain adaptation and handling out-of-vocabulary words persist, driving ongoing research in self-supervised training and hybrid embedding models.



SentenceTransformersrepresentasignificantadvancementin NLP,offeringefficient,high-qualitysentenceembeddingsfor awiderangeof applications. Byleveraging transformer-based architectures, pooling mechanisms, and contrastive learning techniques,they outperformtraditionalembeddingapproaches in capturing semantic similarity.

Webscrapingistheprocessofextractingdatafromwebsitesand storingitforfurtheruse,ofteninstructuredformatslikeJSONor databases. BeautifulSoup (BS4) is a widely used Python library for parsing HTML and XML, making it easier to navigate and extract specific elements from a webpage. The general web scrapingprocessinvolvessendinganHTTPrequesttoawebsite, parsing the retrieved HTML content, extracting relevant data, and storing it in a structured format such as JSON ora database. To begin with, a request is made to a target website using the requestslibrary, which fetches the page's HTML content.

Once the response is received, BeautifulSoup parses the HTML usingparserslike"html.parser"or"lxml",allowingeasy navigation and data extraction through tag-based searching and CSS selectors. The extracted data is then formatted into Question-Answer (QA) pairs, which can be useful for AI models, FAQs,

and research purposes. For example, if a website contains multiple headings **as** questions (e.g., <h2>tags) and correspondingparagraphsasanswers(<p>tags),wecanscrape and pair them together. The extracted data is often stored in JSON format using Python's built-in jsonmodule, which ensures easy accessibility and compatibility with various applications.

This script performs three main steps: fetching the webpage content using requests, parsing and extracting relevant QA datausingBeautifulSoup**,**andsavingtheextracteddatainJSON format**.** However, JSON is not the only storage option. In many applications, scraped data is stored in a database for efficient querying, retrieval, and further processing. Databases such as SQLite, MySQL, and PostgreSQL offer structured storage, making it easy to manage and analyze large-scale scraped data. BelowisanexampleofhowtostoretheextractedQApairsinan SQLite database using Python's sqlite3 module. usefulformachinelearningmodels,chatbots,andanalyticalappli- cations.



Another important aspect of web scraping is handling dynamic content. Many websites use JavaScript to load data dynamically, which cannot be scraped using BeautifulSoup alone. For such cases, tools like Selenium and Playwright are used to automate browserinteractionsandextractdynamicallyloadedcontent.Sele- nium launches a headless browser instance, executes JavaScript, andretrievesthefullyrenderedHTMLcontent,whichcanthenbe processedusingBeautifulSoup.Here's asimple example ofusing Selenium to scrape dynamically loaded data.

Selenium can handle interactions such as clicking buttons, scrollingpages,andsubmittingforms,makingitidealforscraping websites that require user interactions before revealing content.

However, it is slower than requests + BeautifulSoup because it simulatesarealbrowser.WebscrapingusingBeautifulSoupisa powerful technique for extracting structured data from websites and storing it for various applications .

## VI. DATASET DESCRIPTION FOR GROUNDWATER KNOWLEDGE RETRIEVAL SYSTEM

The datasetusedinthis researchis astructuredcollection ofgroundwater-relatedquestion-answer(QA)pairs designedforan AI-driven knowledge retrievalsystem. The data is embedded using Sentence Transformers to facilitate efficient semantic search. The primaryobjectiveistoprovideaccurate,context-awareresponsesto groundwater-related queries. The dataset is stored in multiple formats, including JSON for structured QA pairs and Pickle (PKL) for embedded vectors, ensuring optimized storage and retrieval.



### A. Dataset Components
Thedatasetcomprisesthreeprimarycomponents:
1) Question-AnswerPairs(qa_db.json)
2) PrecomputedEmbeddings(groundwater_qa.pkl)
3) MetadataandAuxiliaryInformation

Question-Answer Pairs (qa_db.json) : The qa_db.json file contains structured question-answer pairs related to groundwater. These pairs are carefully curated from multiple sources, including academicresearch,governmentreports,andexpertknowledge.The datasetincludesover1000QApairs,categorizedintothefollowing themes:

### B. BasicGroundwaterConcepts
1) Definitionsofgroundwater,aquifers,andrecharge processes
2) Differencesbetweenconfinedandunconfinedaquifers
3) Importanceofgroundwaterinwatersupplyand ecosystems

### C. GroundwaterQualityandPollution
1) Common groundwater contaminants (arsenic, nitrates, heavy metals)
2) Causesandconsequencesofgroundwaterpollution
3) Methodstopreventcontaminationandensuresafe drinking water.

### D. GroundwaterDepletionandManagement
1) Causesofgroundwaterdepletion(over-extraction, urbanization, climate change)
2) Conservationmethods,includingrainwaterharvesting and recharge wells
3) Governmentinitiativesandpoliciesforgroundwater management

### E. MetadataandAuxiliaryInformation
InadditiontoQApairsandembeddings,thedatasetincludes metadata such as:
4) Question Categories: Each question is labeled under themeslike"Pollution,""Depletion,"or"Management."
5) Source Information: Some answers are sourced from governmentreports,scientificliterature,ornewsarticles.
6) Confidence Scores: Responses retrieved through similaritymatchingareassignedaconfidencescore based on cosine similarity.

*F. ImplementationintheAISystem*

ThedatasetisintegratedintoanAI-drivenchatbotandsearch engine through the following steps:

1) LoadingPrecomputedEmbeddings:
- Theembeddingsfromgroundwater_qa.pklare loaded using torch for fast similarity comparison.
- UserQueryProcessing:
2) Theuser'squeryisembeddedusingthesame SentenceTransformer model.
- Cosine similarity is com+-puted with the storedembeddingstofindthebest-matching question.
- ResponseRetrieval:
3) Ifthesimilarityscoreisabove0.7,thematched answer from qa_db.json is returned.
4) If the score is between 0.4-0.7, a hybrid responseisgenerated,combiningstoredanswers andlivewebsearchresults(viaGoogleCustom Search API).
5) Ifthesimilarityisbelow0.4,themodelinforms the user that the knowledge base lacks the required information.

This dataset is a specialized, structured knowledge base tailored for groundwater-related inquiries. By leveraging QA pairs,precomputedembeddings,andmetadata,itenablesreal- time, AI-driven responses while ensuring accuracy and scalability.Itshybriddesign—integratingstaticQAdatawith dynamic web search—makes it a powerful tool for groundwater awareness and policy research. Future enhancements may include expanding the dataset,integrating userfeedbackloops,andincorporatingmultilingualsupportto make groundwater knowledge more accessibl

## VII.QUERY PROCESSING AND SIMILARITY COMPUTATION

The query processing pipeline initiates with user input normalization,leveraging characterembeddingmodels tohandle typographical variations. The input query is transformed into sentence-level embeddings via Sentence Transformers, ensuring compatibility with the precomputed vector space.Cosine similarity metrics are employed to determine query proximity within the stored embedding space. A similarity threshold of 0.7 isdefinedforhigh-confidencematches,whileascorerangeof0.4 to 0.7 triggers hybrid retrieval, integrating precomputed embeddings with real-time web search results.

For ambiguous queries, Latent Semantic Indexing (LSI) and Word Mover's Distance (WMD) are incorporated to infer underlying contextual semantics. If a query lacks a direct match within the knowledge base, Google Custom Search API is invoked,leveragingprogrammaticwebsearchtoretrieveexternal corroborative sources.The response ranking mechanism utilizes weighted fusion scoring, where retrieval confidence, semantic coherence, and source credibility contribute to ranking the most relevant results. The response synthesizer aggregates high- ranking results using extractive summarization techniques (e.g., BART and T5 Transformers), ensuring concise, high-fidelity responses.



Deployment and Scalability Considerations: The final deployment architecture integrates containerized microservices, ensuring modularity and scalability. The core NLP model and embedding store are deployed within a Kubernetes cluster, enabling auto-scaling based on query demand metrics.The application interface is powered by Streamlit, providing a user-friendlyconversationalUI.ThebackendincorporatesFastAPIfor asynchronous request handling, minimizing latency in query processing. This methodology ensures high retrieval precision, low latency, and scalable cloud-based deployment, making the groundwater knowledge retrieval system robust, efficient, and adaptable for future research and enhancements.

## VIII. CONCLUSION

This research presents a hybrid AI-driven groundwater knowledge retrieval system, integrating semantic search, web scraping, and real-time information retrieval. By leveraging Sentence Transformers, the system generates context-aware embeddings that enhance the accuracy and efficiency of knowledge retrieval. The BeautifulSoup-based web scraping pipeline ensures continuous data acquisition, while Google CustomSearchAPIsupplementsknowledgegapswithreal-time external sources.A key contribution is the embedding-based semantic search, optimized through FAISS indexing and cosine similarity,achievinghigh-speed,high-precisionquerymatching. Additionally, the integration of hybrid retrieval mechanisms— combining static QA pairs, deep learning embeddings, and web search—improves responseaccuracy.Thesystem's deployment onGPU-accelerated cloudinfrastructure,alongwith FastAPIand Streamlit, ensures scalability, real-time interaction, and low- latency responses.Despite its advancements, challenges such as domainadaptation,computationalcosts,andevolvingdataneeds remain. Future research can explore self-supervised learning, multilingual adaptation, and federated AI models to improve contextualgeneralization and real-world applicability.This study demonstrates the efficacy of AI-enhanced groundwater knowledgeretrieval,offeringascalable,efficient,andintelligent solution for environmental research, policy-making, and public awareness.

## REFERENCES

[1] Mitchell, R. (2018). Web Scraping with Python: CollectingMoreDatafromtheModernWeb.O'Reilly Media.
[2] Kougia,V.,Kalogiros,C.,&Daras,P.(2021)."Asurvey on web crawling and data scraping for open-source intelligence (OSINT)." IEEE Access, 9, 29513-29537.
[3] Boehmke,B.C.,&Greenwell,B.(2020).Hands-On Machine Learning with R. CRC Press.
[4] Singh,H.,&Singh,A.(2020)."Acomparativeanalysis of web scraping techniques for data extraction." International Journal of Computer Science and Information Security (IJCSIS), 18(4), 76-82.
[5] Reimers, N., & Gurevych, I. (2019). "Sentence-BERT: SentenceembeddingsusingSiameseBERT-networks." Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP), 3982-3992.
[6] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). "BERT: Pre-training of deep bidirectional transformersforlanguageunderstanding."Proceedingsof NAACL-HLT 2019, 4171-4186.
[7] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,L.,Gomez,A.N.,Kaiser,Ł.,&Polosukhin,I. (2017)."Attentionisallyouneed."AdvancesinNeural InformationProcessingSystems(NeurIPS),30,5998-6008.
[8] Johnson, J., Douze, M., & Jégou, H. (2019). "Billion- scalesimilaritysearchwithGPUs."IEEETransactions on Big Data, 7(3), 535-547.
[9] Guo,J.,Fan,Y.,Ai,Q.,&Croft,W.B.(2016)."Adeep GoogleDevelopers.(2024).CustomSearchJSONAPI documentation. Retrieved from https://developers.google.com/custom- search
[10] Dean, J., Ghemawat, S., & Sanjay, G. (2008). "MapReduce: Simplified data processing on large clusters."CommunicationsoftheACM,51(1),107-113.
[11] Mikolov,T.,Sutskever,I.,Chen,K.,Corrado,G.,&Dean, J. (2013). "Distributed representations of words and phrasesandtheircompositionality."AdvancesinNeural Information Processing Systems (NeurIPS), 26, 3111-3119.
[12] Nogueira, R., Cho, K., & Lin, J. (2019). "Passage re- rankingwithBERT."arXivpreprintarXiv:1901.04085.
[13] Famiglietti,J.S.(2014)."Theglobalgroundwater crisis."Nature ClimateChange,4(11), 945-948.
[14] Scanlon,B.R.,Ruddell,B.L.,Reed,P.M.,Hook,S.J., & Longuevergne, L. (2017). "Drought risk mitigation: Water management and hydrologic infrastructure." Water Resources Research, 53(7), 5468-5476.
[15] Gleeson,T.,Wada,Y.,Bierkens,M.F.,& VanBeek,L. P.(2012)."Waterbalanceofglobalaquifersrevealedby groundwater footprint." Nature, 488(7410), 197-200.
[16] Bierkens, M. F., & Wada, Y. (2019). "Non-renewable groundwateruseandgroundwaterdepletion:Areview." Environmental Research Letters, 14(6), 063002.
[17] Karmakar,S.,Simonovic,S.P.,Peck,A.,& Blackport, R. (2010). "Flood forecasting using artificial neural networks: Methodological issues and applications." EnvironmentalModelling&Software,25(5),805-818.
[18] Trilles,S.,Luján,A.,Díaz,L.,&Huerta,J.(2020)."An artificial intelligence approach for modeling groundwater resources using machine learning techniques." Hydrology, 7(3), 56.
[19] Moriasi,D.N.,Arnold,J.G.,VanLiew,M.W., Bingner, R. L., Harmel, R. D., & Veith, T. L. (2007). "Model evaluation guidelines for systematic quantification of accuracyinwatershedsimulations."Transactionsofthe ASABE,50(3),34.

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY