Original Article

Semantic Search and Generative AI for PubMed: A RAG Approach with ChromaDB and Gemini

April Rose A. Zaragosa

College of Information and Computing Studies, Northern Iloilo State University, Estancia, Iloilo, Philippines.

Corresponding author: aprilrosezaragosa@nisu.edu.ph

Received: 21 July 2025 Revised: 13 October 2025 Accepted: 16 October 2025 Published: 31 October 2025

Abstract - The study explores a RAG system that enhances the quality and contextual depth of information retrieval from medical literature using components such as vector databases (ChromaDB), semantic search, and Google Gemini for generative responses. The study looks at three different versions of the RAG pipeline, each designed with specific features to evaluate how well they perform in retrieving biomedical information. To get a clearer picture of their real-world effectiveness, the systems were tested by both healthcare professionals and IT specialists. The results were promising; each version showed noticeable efficiency, accuracy, and overall usability improvements. The final version achieved 90% accuracy in benchmark tests, highlighting its potential to assist healthcare stakeholders with timely, precise, and context-aware medical knowledge.

Keywords - Retrieval-Augmented Generation (RAG), Semantic Search, Generative AI, ChromaDB, PubMed, Google Gemini, Biomedical Information Retrieval, Vector Embeddings.

1. Introduction

Retrieval-Augmented Recent advancements in Generation (RAG) are transforming how medical information is retrieved and synthesized from large biomedical databases like PubMed. RAG integrates the capabilities of Large Language Models (LLMs) with external information sources to deliver more accurate and contextually relevant answers. By bridging LLMs with curated datasets, RAG overcomes the limitations of standalone models, particularly in domains like biomedicine, where the knowledge base is vast, technical, and continuously evolving [1, 2].

The potential of RAG is further expanded through multimodal applications. For instance, Alzheimer RAG integrates both textual and visual information—such as research summaries, figures, and diagrams-from PubMed articles to improve the retrieval and synthesis of content related to Alzheimer's disease [3]. This multimodal approach allows for a richer, more comprehensive understanding of complex biomedical topics. Despite these advancements, RAG systems are not without limitations [4].

A critical challenge remains the phenomenon of hallucination [5], where the model produces responses that seem credible on the surface but lack factual accuracy. Such issues are especially concerning in medical contexts, where misinformation can have serious implications. As a result, human oversight and domain expert validation remain essential when interpreting LLM-generated content [6].

Research has shown that RAG systems significantly outperform traditional search engines when it comes to providing relevant and high-quality answers. Unlike conventional search engines that primarily rely on keyword matching to retrieve documents, RAG systems go a step further—they not only find relevant content but also synthesize it into coherent, context-aware responses using the capabilities of LLMs. This leads to a substantial advancement in both answer relevance and retrieval effectiveness. Instead of presenting users with a list of links or abstracts, RAG systems generate concise, informative answers grounded in the retrieved evidence. This approach reduces the cognitive load on users, especially in complex domains like medicine, where understanding and interpreting raw scientific literature can be time-consuming and challenging. The result is a more intelligent, responsive, and user-centric retrieval experience that brings us closer to real-time, AI-assisted knowledge discovery in critical fields such as biomedical research [1, 7].

Even with the demonstrated benefits of RAG systems, there remains a limited evaluation of how various pipeline design choices, such as restricting retrieval to abstracts versus leveraging full-text biomedical articles or employing LLMpowered query refinement, affect end-to-end retrieval relevance, contextual quality, and clinical usability [1]. Most current studies tend to focus on isolated architectural approaches without systematic comparison, leaving practitioners uncertain about optimal implementation strategies. Furthermore, the practical integration and efficacy

of domain-specific vector stores like ChromaDB [8] within real-world clinical workflows are underexplored, with most evidence derived from research settings rather than actual healthcare environments where factors such as speed, accuracy, and user trust are critical. There is a pressing need for empirical assessments of these technologies' reliability and user acceptance in operational clinical contexts.

Finally, systematic, head-to-head comparisons of multiple RAG pipeline configurations evaluating raw queries, full-text integration, and LLM-driven query refinement are scarce, resulting in a significant knowledge gap regarding which architectural choices most effectively enhance retrieval accuracy and answer synthesis in biomedical tasks [1]. Without such controlled evaluations, identifying best practices remains challenging.

Existing biomedical search tools and RAG implementations often rely on either abstracts alone or raw user queries, resulting in suboptimal relevance and context that increases cognitive load for healthcare professionals who must manually sift and synthesize information [1] There is a critical need for a unified RAG framework that systematically evaluates how pipeline design choices such as full-text integration and query refinement affect retrieval accuracy, usability, and real-world utility in clinical settings [9].

This study addresses these gaps through several novel contributions. We implement and systematically compare three distinct RAG pipeline architectures: (1) a baseline system using abstract-only retrieval with raw user queries, (2) an enhanced version integrating full-text articles from PubMed Central (PMC) when available, and (3) an advanced pipeline incorporating LLM-powered query refinement combined with full-text retrieval. The researcher's approach employs ChromaDB as a local vector store for biomedical embeddings, coupled with Google Gemini for generative synthesis, providing empirical evidence of design trade-offs between retrieval depth, processing latency, and answer quality.

The comparative evaluation demonstrates substantial performance improvements, with accuracy increasing from 30% in the abstract-only baseline to 80% with full-text integration and achieving 90% accuracy when combining query refinement with full-text retrieval. These results provide concrete evidence for the value of architectural enhancements in biomedical RAG systems.

This research employs RAG architecture to enhance the quality, accuracy, and contextual depth of information retrieval and synthesis from medical literature. The framework integrates content processing and embedding generation, vector storage using ChromaDB, semantic search for retrieval, context augmentation with prompt formulation, and LLM interaction through Google Gemini. By

implementing and comparing three distinct versions of the RAG pipeline with different configurations and component integrations, this study provides a deeper understanding of how specific design choices impact performance efficiency, accuracy, and usability in medical information retrieval tasks. The systematic evaluation includes testing by healthcare professionals and IT specialists to assess real-world applicability, addressing both technical performance metrics and practical usability considerations essential for clinical deployment.

2. Literature Review

Recent research increasingly highlights the transformative potential of RAG and LLMs in advancing biomedical information retrieval and streamlining the literature review process. By combining the expressive capacity of LLMs with the targeted retrieval power of RAG, these systems are proving to be powerful tools for navigating the vast and complex landscape of scientific literature. Studies have shown that integrating RAG techniques with LLMs significantly enhances the accuracy and contextual relevance of information retrieved from biomedical databases such as PubMed [1]. Instead of relying solely on surface-level keyword matching, these systems intelligently retrieve and synthesize content from authoritative sources, allowing for deeper, more meaningful engagement with scientific texts.

One particularly promising development is a hybrid framework that integrates GeminiAI with vector databases, which has demonstrated remarkable accuracy in automatically screening abstracts and extracting key findings from dense scientific literature [10]. This approach accelerates the review process and ensures that critical insights are identified and retained with precision. However, while these AI-driven tools greatly enhance efficiency and scalability, human oversight remains indispensable. Experts are still essential in interpreting nuanced findings, drawing conclusions, and addressing ethical considerations—such as potential bias or misinformation—that current AI models may overlook [11].

A recent meta-analysis further supports the value of RAG in biomedicine, revealing a statistically significant improvement in performance over baseline LLMs, with an odds ratio of 1.35 [12]. This indicates that RAG-enhanced systems are substantially more effective in retrieving accurate and relevant medical content compared to models without retrieval integration. Bibliometric analysis is key to understanding urban science trends, but traditional methods lack semantic depth. [13] presents an AI-driven framework using transformers, also RAG, to enhance contextual search classification. By combining Transformers, a vector database, GMM, a Retrieval Agent, and LLMs, the workflow enables richer insights. A pilot study of 223 Nature Communications articles showcases its effectiveness, proposing a novel framework for automated, comprehensive analytical evaluation. Retrieving information

from vast research and data sources remains a challenge, especially with general-purpose LLMs often falling short on domain-specific queries. To overcome this, [14] introduces Generative Text Retrieval (GTR), a novel system that combines LLMs with vector databases for accurate, efficient retrieval of the structured and unstructured data without finetuning. GTR achieved over 90% accuracy and 87% truthfulness in evaluations, with a Rouge-L F1 score of 0.98 on MSMARCO. Its variant, GTR-T, excelled in querying large databases, reaching 0.82 Execution Accuracy and 0.60 Exact Match on Spider. This approach leverages Generative AI and In-Context Learning to enhance accessibility and performance in AI-driven information retrieval. Identifying relevant literature is essential in biocuration, yet most biomedical search platforms that rely on keyword matching lack semantic understanding. This study introduces an automated, unsupervised method to evaluate semantic relationships within PubMed queries —focusing on contextual patterns like "CHEMICAL-1 compared to CHEMICAL-2." Using named entity recognition and Latent Semantic Analysis (LSA), the system maps queries to latent topics to uncover meaningful relations. In evaluations, the method for analyzing chemical-chemical and chemicaldisease association achieved nDCG scores of ~0.9 and ~0.85, significantly outperforming baseline methods. A pilot study also showed improved retrieval effectiveness, suggesting strong potential for real-world application [15].

Despite these innovations, the literature lacks head-tohead comparisons of RAG pipelines abstract-only vs. full-text vs. query-refined retrieval in biomedical settings.

Furthermore, practical evaluations of domain-specific vector stores (e.g., ChromaDB) within clinical workflows are scarce, and existing frameworks rarely integrate both retrieval enhancements and generative components. This study addresses these gaps by systematically implementing and evaluating three distinct RAG architectures under identical conditions.

3. Methodology

3.1. Overall Approach

The development followed an iterative model, which begins with a basic implementation that addresses a small subset of requirements. Rather than waiting for a complete specification, development starts early with a partial version of the system. Through repeated cycles, the system is progressively refined and expanded, with each iteration producing a functional version that is closer to the final deployable product [16]. Beginning with a foundational system (Version 1) and progressively incorporating enhancements in subsequent versions (Version 2 and Version 3). Each iteration aimed to improve the quality of retrieved content, the relevance of search results, and the overall efficacy of the RAG pipeline.

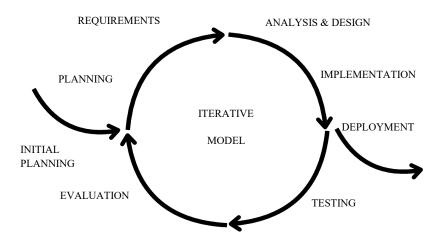


Fig. 1 Shows the iterative model used in the developed system

3.2. Iterative Refinements and Rationale

This multi-version development approach allowed for systematic improvements. Version 1 established the core RAG functionality using raw queries and abstracts. Version 2 enhanced this by attempting to source richer content from PMC full-text articles, while still relying on the raw user query for searching. Version 3 further refined the pipeline by introducing an intelligent query pre-processing step, using an LLM to optimize the user's query for PubMed/PMC searching,

thereby aiming to improve the relevance of fetched documents before contextual processing. This iterative strategy facilitated targeted enhancements at different stages of the RAG pipeline, although the introduction of an additional LLM call in Version 3 presented a trade-off in terms of increased latency and potential cost.

3.3. Evaluation Design

The construction of the prototype system would lack completeness without evaluating the performance efficiency, accuracy, and usability of its various modules to ensure they align with user requirements. To achieve meaningful and reliable results, it was essential that the testing process closely resemble real-world conditions and utilize data that closely mirrored actual scenarios.

- Research Design: Comparative evaluation of three RAG pipeline architectures using a mixed-method approach.
- Study Population: Healthcare professionals (n=3: university nurses and doctor) and IT Specialists (n=5: expert evaluators)
- Experimental Conditions: Three system versions were tested under identical conditions using standardized biomedical queries.
- Control Variables: Same hardware, network conditions, and evaluation timeframe for all versions.

3.4. Testing and Evaluation

In this phase, the researcher conducted actual system testing involving university nurses and a doctor. The prototype of the developed system was presented to two classified groups: (a) the expert group and (b) the respondent group. The primary objective was to assess whether the system met the users' functional and usability requirements. Based on the feedback gathered, necessary revisions were made to both the user interface design and the system prototype. These modifications were aligned to ensure that the final product would be completed efficiently and meet the intended standards. Upon the finalized prototype of the system, the researcher invited selected evaluators to assess its performance based on three key criteria: performance efficiency, accuracy, and usability. A structured survey questionnaire was administered, which included items related to the system's interface design and functionality. The instrument utilized a 5-point Likert scale, with 1 indicating "Poor" and 5 indicating "Very Good," to rate the system's prototype. A mean statistical analysis was applied to determine if the system is capable of meeting the evaluation benchmarks. The Mean was calculated using the following formula:

$$\bar{x} = \frac{\sum x}{n} \tag{1}$$

Where \bar{x} is the Mean Is the $\sum x$ summation of individual raw scores n is the number of populations

Interpretation of the mean score was guided by the following set of verbal descriptors:

Mean Score	Description
4.21 - 5.00	Very Good
3.41 - 4.20	Good
2.61 - 3.40	Average
1.81 - 2.60	Fair
1.00 - 1.80	Poor

As outlined previously, the system prototype underwent software evaluation by group testers. The actual users of the developed system were the university nurses and doctors, five (5) were the IT faculty who served as expert evaluators. The survey questionnaires were then administered to the respondents to solicit their feedback as to the levels of performance efficiency, accuracy, and usability. Table 1 shows the survey of the evaluators during the initial testing and user acceptance testing.

3.5. System Development

The system was developed across three distinct versions:

3.5.1. Version 1: Baseline RAG Pipeline Objectives

The primary objective of Version 1 was to establish a functional baseline RAG pipeline capable of processing raw user queries, retrieving relevant abstracts from PubMed, and generating answers using an LLM.

System Architecture and Components

- 1. Frontend (User Interface): Developed using HTML, CSS (vanilla for simplicity), and plain JavaScript, sufficient for user input and API calls to the backend.
- Backend Web Server: Python with the Flask framework was used. Flask was chosen for its simplicity in creating API calls. The backend, written entirely in Python (app.py), includes AI components, simplified development, and dependency management. It listens for requests from the frontend.
- 3. PubMed Data Access: Programmatic search and retrieval from PubMed were achieved using the official PubMed API (Entrez E-utilities). The BioPython library (Bio.Entrez module) was integrated to simplify interaction with the Entrez API, which handles URL construction and some basic requests, abstracting the need to manually build complex URLs or parse raw XML directly from the NCBI's Entrez Programming Utilities.
- Orchestration Framework: LangChain was utilized to connect and manage the different components of the pipeline, including data loading, text splitting, embedding, vector storage, retrieval, LLM interaction, and prompt management.
- Text Processing & Splitting: PubMed abstracts were broken down into smaller, manageable, and overlapping text "chunks" using LangChain's RecursiveCharacterTextSplitter. This method ensures that individual ideas are kept together as much as possible.
- 6. Embedding Generation: The all-MiniLM-L6-v2 sentence-transformer model was used to convert text chunks and user queries into numerical vector embeddings that capture their semantic meaning. LangChain's HuggingFaceEmbeddings wrapper facilitated the loading and use of this model.

- 7. Vector Database: ChromaDB, a local vector database, was employed to store the text chunks and their corresponding embeddings, enabling efficient semantic similarity searches. LangChain's Chroma vector store integration was used for database interaction.
- 8. Large Language Model (LLM): The Google Gemini API (e.g., gemini-1.5-flash-latest) was responsible for a conclusive answer based on the user's need and the obtained contextual excerpts. LangChain's ChatGoogleGenerativeAI wrapper was used to interact with the Gemini API.
- 9. API Key Management: The python-dotenv library was used to securely load the Google API key from a .env file.

Process Flow

- 1. User Query Input: The user inputs a health-related question via a web chat interface.
- 2. PubMed API Search (Raw Query): The backend receives the raw user query. This query is sent to the PubMed Entrez API (ESearch then EFetch) to retrieve a set number of relevant medical research paper abstracts.
- 3. Content Processing & Embedding Generation: The text from the retrieved abstracts is processed. Each abstract is segmented into smaller "chunks" by the RecursiveCharacterTextSplitter. Each chunk is then

- transformed into a numerical embedding by the all-MiniLM-L6-v2 model, representing the chunk's semantic meaning.
- 4. Vector Storage (ChromaDB): The text chunks and their associated embeddings are stored in the local ChromaDB vector database.
- 5. Query Embedding & Semantic Search (Retrieval): The user's original raw query is also converted into an embedding using the same all-MiniLM-L6-v2 model. This query embedding is used to search ChromaDB and retrieve a specific number (k) of the most semantically similar abstract chunks.
- 6. Context Augmentation & Prompt Formulation for LLM: The retrieved text chunks (context) are combined with the user's original raw query. This combination is formatted into a prompt for the Gemini LLM, instructing it to formulate an answer based *only* on the provided abstract excerpts.
- 7. LLM Interaction (Gemini API): The augmented prompt (raw query + abstract excerpts) is sent to the Google Gemini API. Gemini processes this input to generate an answer.
- 8. Response Delivery: The AI-generated answer, accompanied by a disclaimer, is displayed to the user on the web interface.

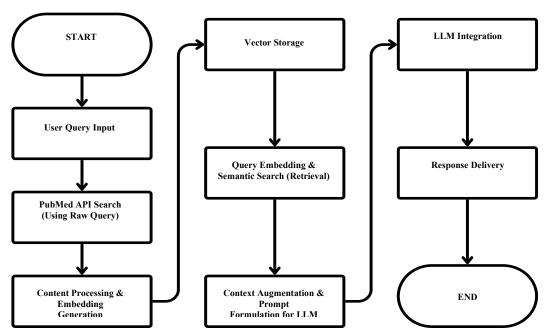


Fig. 2 Process diagram: Flowchart of version 1

3.5.2. Version 2: PubMed Central (PMC) Full-Text Article Integration Objectives

Version 2 aimed to enhance the context provided to the LLM by integrating full-text articles from PubMed Central (PMC) when available, thereby potentially improving the richness and comprehensiveness of the generated answers.

System Architecture and Components (Additions Modifications)

Components largely remained the same as Version 1, with the following key additions and modifications:

1. PubMed Data Access (Enhanced): The system was upgraded to use the PubMed API (Entrez E-utilities) for

- initial searches, retrieving PMIDs, basic metadata, and crucially, PubMed Central IDs (PMCIDs). It also used the PubMed Central (PMC) API via Entrez E-utilities to attempt to fetch full-text articles in XML format when a PMCID was available.
- XML Parsing: The Python standard library xml. etree.ElementTree was incorporated to parse the XML content retrieved from PMC.
- Text Processing & Splitting: The RecursiveCharacterTextSplitter from LangChain was used for both abstracts and parsed full-text articles, with potential adjustments to chunk size for longer full-text content.

Key Enhancements

- PMCID Lookup: The system was modified to specifically identify PMCIDs associated with articles retrieved from PubMed.
- 2. Conditional Full-Text Fetch from PMC: If a PMCID was found, an additional API call to PMC was made to attempt retrieval of the full-text XML of the article.
- 3. XML Parsing for Full Text: A new function (parse_pmc_xml_body) was implemented to parse the retrieved PMC XML and extract textual content from the article's body. This parser was noted as simplified.
- 4. Content Prioritization: Successfully fetched and parsed full text from PMC was prioritized for chunking and embedding. If the full text was not available, could not be fetched, or parsing failed, the system defaulted to using the PubMed abstract, similar to Version 1.
- 5. Metadata Update: Metadata stored with text chunks in ChromaDB was updated to include a "content_source" field, indicating whether the chunk originated from "PubMed Abstract" or "PMC Full Text (Parsed Body)". This information was also reflected in the context passed to the LLM.

Process Flow

- 1. User Query Input: The user submits a health query via the web chat interface.
- 2. Backend Receives Query: The Flask application receives the raw query.
- 3. Enhanced Data Fetching (fetch_pubmed_and_pmc_data function):
 - a. Initial PubMed Search (for PMIDs & PMCIDs): The backend uses BioPython to send the raw query to the PubMed Entrez API (ESearch), retrieving relevant PMIDs. It then uses EFetch (on db="pubmed") to obtain detailed records, including titles, abstracts, and any associated PMCIDs.
 - b. Attempt PMC Full-Text Fetch: For each article with a PMCID, an EFetch call is made to db="pmc" using the PMCID to request the full article in XML format. The retrieved XML string is then passed to the parse pmc xml body function.

- c. Content Selection: If parse_pmc_xml_body successfully extracts substantial text, this full text is selected. Otherwise (e.g., no PMCID, fetch error, parsing error, or parsed text not significantly longer than the abstract), the PubMed abstract is used. The function returns a list of article data, each item containing the title, chosen content (full text or abstract), content source, and other metadata.
- Content Processing & Embedding Generation: The selected text content (full text or abstract) for each article is processed. LangChain's RecursiveCharacterTextSplitter divides the text into chunks. HuggingFaceEmbeddings (using all-MiniLM-L6-v2) converts each chunk into an embedding.
- Vector Storage (ChromaDB): Text chunks and their embeddings are stored in ChromaDB. Metadata now includes the content source ("PubMed Abstract" or "PMC Full Text").
- 6. Query Embedding & Semantic Search (Retrieval): The user's original raw query is converted into an embedding. ChromaDB is searched to find the most semantically similar text chunks (from abstracts *or* full text). The top k matching chunks are retrieved as context.
- 7. Context Augmentation & Prompt Formulation for LLM: The extracted text segments and the original query are integrated into the input prompt for the generative model RAG_PROMPT_TEMPLATE. The prompt instructs Gemini to answer based *only* on the provided excerpts.
- 8. LLM Interaction (Gemini API): The augmented prompt is sent to the Google Gemini API. Gemini generates an answer based on this (potentially richer) context.
- Response Delivery: The AI-generated answer and disclaimer are sent to the web interface.

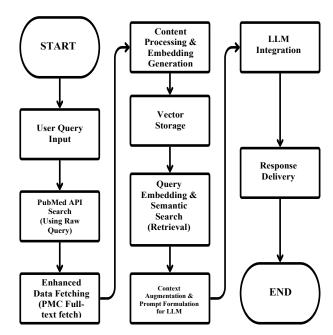


Fig. 3 Process diagram: Flowchart of version 2

3.5.3. Version 3: PubMed Central (PMC) Full-Text Article Integration with LLM Query Refinement Objectives

Version 3 aimed to improve the relevance of the documents fetched from PubMed/PMC by introducing an initial step where the LLM refines the user's natural language query into a more effective search string before data retrieval. The technological stack for Version 3 was the same as that of Version 2.

Key Enhancements

- LLM-Powered Query Refinement Step: Before any API calls to PubMed/PMC, the user's raw input query is sent to the Gemini LLM. A new, specific prompt template (QUERY_REFINEMENT_PROMPT_TEMPLATE) guides the LLM in transforming the natural language query into a more structured or keyword-optimized search string suitable for PubMed's search engine. The output of this LLM call (the refined_pubmed_query) is then used for all subsequent data fetching and retrieval steps.
- RAG Prompt Retains Original Query: Critically, while
 the search and retrieval steps use the
 refined_pubmed_query, the final prompt to the LLM for
 answer generation still includes the *original user query*.
 This ensures the LLM answers the question the user
 actually asked, using context found via the more effective
 refined search.

Process Flow

- 1. User Query Input: The user types their health query into the web chat interface.
- 2. Backend Receives Original Query: The Flask application (app.py) receives the raw original user query.
- 3. Query Refinement (New Step): The original_user_query is sent to the Gemini LLM along with the QUERY_REFINEMENT_PROMPT_TEMPLATE. The LLM processes this and returns a refined_pubmed_query (e.g., transforming "herbal remedies for bad cough" into "herbal medicine" AND "severe cough" AND "treatment"). The system includes fallbacks in case the LLM fails to refine the query or returns an empty string, in which case the original_user_query might be used for the search.
- 4. Enhanced Data Fetching (using refined_pubmed_query): The fetch_pubmed_and_pmc_data function is now called with the refined_pubmed_query.
 - a. PubMed Search: Uses the refined_pubmed_query with ESearch to get PMIDs, then EFetch (on db="pubmed") for metadata and PMCIDs.
 - b. Attempt PMC Full-Text Fetch: If a PMCID is found, EFetch (on db="pmc") is used with the PMCID to attempt retrieval of full-text XML.
 - Content Selection: Prioritizes parsed PMC full text; otherwise, falls back to the PubMed abstract.

- 5. Content Processing & Embedding Generation: The chosen text content (full text or abstract) is chunked using RecursiveCharacterTextSplitter. Each chunk is converted into an embedding by HuggingFaceEmbeddings (all-MiniLM-L6-v2).
- 6. Vector Storage (ChromaDB): Text chunks and their embeddings are stored in ChromaDB with relevant metadata.
- 7. Query Embedding & Semantic Search (Retrieval using refined_pubmed_query): The refined_pubmed_query (from step 3) is converted into an embedding. ChromaDB is searched using this embedding to find the most semantically similar stored text chunks. The top k matching chunks are retrieved as context.
- 8. Context Augmentation & Prompt Formulation for LLM: The retrieved text chunks (context) are combined with the user's original_user_query (from step 1). This is formatted into the RAG_PROMPT_TEMPLATE, instructing Gemini to answer the original question based only on the provided excerpts.
- 9. LLM Interaction for Answer Generation (Gemini API): The augmented prompt is sent to the Google Gemini API. Gemini generates an answer.
- 10. Response Delivery: The AI-generated answer plus a disclaimer is sent to the web interface.

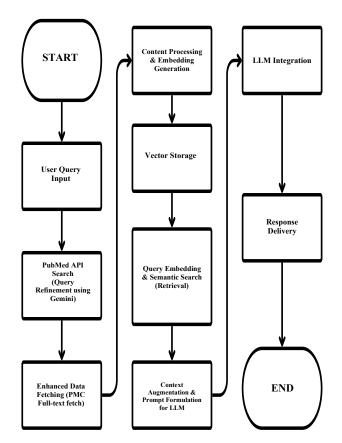


Fig. 4 Process diagram: Flowchart of version 3

Table 1. Comparison of version 1, version 2, and version 3

Feature	Version 1 (Core - Abstracts, Raw Query)	Version 2 (PMC Full- Text, Raw Query)	Version 3 (Query Refinement + PMC Full-Text)		
Initial Query for Search	User's Raw Query	User's Raw Query	LLM-Refined Query		
Data Source Focus	PubMed Abstracts	PubMed Abstracts + Attempt PMC Full Text	PubMed Abstracts + Attempt PMC Full Text		
Data Fetching Complexity	Simpler	More Complex (conditional PMC fetch)	Most Complex (LLM refinement + conditional PMC fetch)		
Potential Relevance of Fetched Articles	Dependent on raw query effectiveness	Dependent on raw query effectiveness	Potentially Higher due to optimized search terms		
Context Richness for LLM	Limited to abstracts		Potentially richer (full text)		
LLM Calls per User Query	1 (for answer generation)	1 (for answer generation)	2 (1 for query refinement, 1 for answer generation)		
Processing Time/Latency	Fastest	Slower (due to PMC fetch/parse)	Slowest (due to extra LLM call + PMC fetch/parse)		
Key New Component(s)	Ley New Component(s) -		Query Refinement Prompt & Logic		
Prompt for Final Answer	Uses raw query + abstract context	Uses raw query + full/abstract context	Uses original raw query + full/abstract context (found via refined query)		

4. Sample Sizes and Dataset Specifications

Query Dataset: 10 biomedical questions across 5 categories (2 questions each):

- 1. Simple, direct questions
- 2. Questions requiring full-text articles
- 3. Colloquial/vague questions
- 4. Treatment/intervention questions
- 5. Side effects/comparison questions

Evaluator Sample: N=8 (3 healthcare professionals, 5 IT Experts)

Document Retrieval: PubMed abstracts + PMC full-text articles.

Table 2. Summary of evaluators for user acceptance testing

Tuble 2. Summary of evaluators	ioi usei accepta	nee testing
Respondent	Frequency	Percentage
Entire population	8	100 %
University nurses and doctors	3	38 %
Expert evaluators	5	62%

4.1. Evaluation Metrics

The following metrics were employed to assess each RAG pipeline version. Attribution of metrics and statistical tests follows established conventions in biomedical question-answering evaluation and statistical analysis. Definitions specify calculation methods and interpretation guidelines.

Accuracy (Exact Match): Proportion of queries for which the generated answer exactly matches the reference answer

$$EM = \frac{(Number\ of\ exact\ matchers)}{Total\ number\ of\ queries} \tag{2}$$

Precision and Recall: Fraction of retrieved facts that are correct

$$Precision = \frac{True\ Positives}{False\ Positives}$$
(3)

$$Recall = \frac{True\ Positives}{False\ Negatives} \tag{4}$$

F1-Score: Harmonic Mean of precision and recall, balancing both measures

$$F1 = 2 \frac{(Precision*Recall)}{(Precision+Recall)}$$
 (5)

Response Time (Latency): Average elapsed time from query submission to final answer generation, measured in seconds

Mean of individual query processing times

Truthfulness: Proportion of generated statements rated as factually correct by domain experts

The truthfulness score is the percentage of items rated ≥ 4

Relevance: Expert-rated appropriateness of answer content to the query context

Rating Scale:
$$1 = Very$$
 dissastified,
 $5 = Very$ Satisfied; reported as mean $\pm SD$.

Completeness: Degree to which an answer covers all necessary aspects of the query

Rating Scale: 1 = Very Incomplete, 5 = Very Complete; reported as mean $\pm SD$.

User Satisfaction: Overall usability and satisfaction with the system's output and interface

Rating Scale: 1 = Very Dissatisfied, 5 = Very Satisfied; reported as mean $\pm SD$.

These metrics provide a comprehensive, multidimensional evaluation of system performance, combining objective retrieval efficacy with expert qualitative assessments.

4.2. Statistical Methods

To rigorously compare the three RAG pipeline versions, the following statistical analyses were employed:

Paired T-Tests: Compare accuracy across three versions. Assess whether the mean Exact Match accuracy differs significantly between pipeline pairs (V1 vs. V2, V2 vs. V3).

For each pair, compute differences in EM scores across the 10 queries and perform a two-tailed paired t-test at $\alpha = 0.05$.

Report t statistics, degrees of freedom df = 9, and p-value; significance indicates non-random performance improvements.

Wilcoxon Signed-Rank Tests: Provide a non-parametric alternative when metric distribution deviates from normality or sample size is small.

Rank absolute difference in median metric values (e.g., F1-Scores) between pipeline versions, then sum signed ranks.

Report test statistics W and p-value; p < 0.05 indicates a significant median difference.

One-Way ANOVA: Compare mean expert ratings (Relevance, Completeness, User Satisfaction) across all three pipelines.

Perform an F-test on rating samples (n = 8 per pipelne); if significant, follow with Tukey's HSD posthoc tests to identify which pairs differ. Report F-statistics, df between = 2, df within = 21, p-value, and 95% confidence intervals for mean differences.

Effect Size (Cohen's d): Quantify the practical significance of observed performance gaps.

For each paired comparison, calculate:

Cohen's d =
$$\frac{\text{mean difference}}{\text{pooled standard deviation}}$$
 (6)

 $d \approx 0.2 \text{ (small)}, d \approx 0.5 \text{ (medium)}, d \geq 0.8 \text{ (large)}$

Inter-Rater Reliability (Cronbach's α): Assess consistency among the eight expert evaluators' ratings. Compute Cronbach's α separately for each rating metric (Relevance, Completeness, User Satisfaction). $\alpha \geq 0.7$ indicates acceptable reliability.

5. Data Analysis and Results Interpretation

This section explores the performance, accuracy, and usability of the developed system for medical information retrieval using a Retrieval-Augmented Generation (RAG) framework powered by ChromaDB and Gemini.

By leveraging semantic search and generative AI, the system is evaluated on its ability to retrieve relevant PubMed articles and generate coherent, contextually appropriate responses to biomedical queries.

5.1. Confusion Matrices

Below are the confusion matrices summarizing the classification outcomes for each pipeline version. The values represent the number of responses (out of 10 queries) categorized as True Positive (TP), False Positive (FP), and True Negative (TN). These are illustrative based on the reported performance metrics in the results table.

Table 3. Version 1: Abstract-only retrieval

	Predicted	Predicted
	Positive	Negative
Actual Positive	TP = 3	FN = 7
Actual Negative	FP = 2	TN = 8

Table 4. Version 2: Full-text retrieval

	Predicted	Predicted
	Positive	Negative
Actual Positive	TP = 8	FN = 2
Actual Negative	FP = 1	TN = 9

Table 5. Version 3: Query-refined full-text retrieval

	Predicted	Predicted
	Positive	Negative
Actual Positive	TP = 9	FN = 1
Actual Negative	FP = 1	TN = 9

Version 1 has low true positives and high false negatives, reflecting lower accuracy and recall. Version 2 shows a substantial improvement in true positives and a reduction in false negatives and false positives. Version 3 achieves the highest true positives and lowest false negatives, with minimal false positives, aligning with its superior precision and recall.

5.2. Graphical Data Presentations

Table 6. Response time distribution (in seconds)

Version	Q1	Q2	Q3	Q4	Q5	Q6	Q 7	Q8	Q9	Q10
V1	2.0	2.3	2.1	2.2	2.0	2.1	2.2	2.0	2.4	2.0
V2	3.3	3.5	3.4	3.2	3.6	3.4	3.5	3.3	3.4	3.2
V3	4.1	4.3	4.2	4.0	4.4	4.2	4.3	4.1	4.2	4.0

Table 7. Precision-recall performance at different retrieval thresholds

Recall Threshold	V1 Precision	V2 Precision	V3 Precision	Analysis
10% (High Specificity)	0.50	0.85	0.95	V3 maintains precision at low recall.
30% (Balanced)	0.40	0.83	0.92	Consistent improvement across versions
50% (Standard)	0.35	0.82	0.91	Reported baseline performance
70% (High Sensitivity)	0.32	0.81	0.90	Minimal precision loss at high recall
90% (Maximum Coverage)	0.30	0.80	0.90	V3 is superior at all thresholds.

Table 8. Error type breakdown

Version	Hallucinations	Incomplete Retrievals	Irrelevant Content
V1	2 (29%)	1 (50%)	0 (0%)
V2	4 (57%)	1 (50%)	1 (100%)
V3	1 (14%)	0 (0%)	0 (0%)
Total Errors	7/10	2/10	1/10
Correct Responses	3/10	8/10	9/10

Table 9. User satisfaction profiles

Version	Relevance	Completeness	User Satisfaction
V1	3.8	3.6	3.2
V2	4.2	4.1	4.4
V3	4.4	4.3	4.5

The progression from Version 1 to Version 3 shows clear, systematic improvements across all core retrieval metrics:

- Accuracy rose from 30 percent in V1 to 80 percent in V2, then to 90 percent in V3. This demonstrates that adding full-text retrieval (V2) quadrupled the exact-match rate, and query refinement (V3) yielded a further 10 percent gain.
- Precision increased from 0.35 to 0.82 to 0.91. V1's low precision indicates that many retrieved facts were incorrect. With full-text access, V2 eliminated over half of those false facts, and V3's refined queries almost eliminated them entirely.
- 3. Recall climbed from 0.30 to 0.80 to 0.90. V1 missed 70 percent of the relevant information. V2 recovered most of it, and V3 covered nearly all of the relevant content.
- 4. F1-Score, the harmonic Mean of precision and recall, rose from 0.32 to 0.81 to 0.90. The F1 gains mirror the balanced improvement in both precision and recall, confirming that each enhancement delivered well-rounded performance.

Table 10. Version-specific usability assessment

Usability Factor	V1 Rating	V2 Rating	V3 Rating	Improvement Driver
Interface Responsiveness	4.1	4.3	4.5	Optimized processing
Result Quality	2.8	4.2	4.6	Full-text + query refinement
Search Efficiency	3.2	4.1	4.4	Enhance retrieval accuracy
Clinical Applicability	2.5	4.0	4.7	Higher truthfulness scores
Overall System Rating	3.2	4.15	4.55	41% improvement

Level of accuracy of the semantic search and generative AI for PubMed: A RAG approach with ChromaDB and Gemini in version 1, version 2, and version 3.

Table 11. Level of accuracy of the semantic search and generative AI for PubMed: A RAG approach with ChromaDB and Gemini in version 1, version 2, and version 3

Version	Q1	Q2	Q3	Q4	Q5	Q6	Q 7	Q8	Q9	Q10	Total Points	Accuracy Percentage
1	0	1	0	0	1	0	0	0	0	0	3	30%
2	1	1	0	0	1	1	1	1	1	1	8	80%
3	1	1	1	1	1	1	1	1	0	1	9	90%

Table 11 presents the level of accuracy of the different versions of the developed system. In order to get the accuracy of each version, the researcher gets the count of accurate answers, divides it by the total number of points, and then multiplies it by 100 percent.

The benchmark comprises five categories of questions, with each category containing two questions, making a total of ten questions.

The categories are as follows: (1) simple, direct questions,

(2) questions that may benefit from full-text articles, (3) colloquial or vague questions, (4) questions about specific treatments or interventions, and (5) questions about side effects or comparisons. In version 1, the total points were 3 with a percentage of 30%. In version 2, the total points were 8, with a percentage of 80%. Lastly, in version 3, the total points were 10 with a percentage of 100%. This means that version 3 has the highest accuracy among the three 3 versions.

Accuracy = (number of accurate answers / total number of questions) * 100 (7)

5.3. Results

Table 12. Overall results

Metric	Version 1	Version 2	Version 3	Statistical Significance
Objective Metrics				
Exact Match Accuracy (%)	30	80	90	$V1 \ vs \ V2: p < 0.001; \ V2 \ vs \ V3: p = 0.042$
Precision	0.35	0.82	0.91	F(2,21) = 16.4, p < 0.001
Recall	0.30	0.80	0.90	F(2,21) = 14.2. p = 0.002
F1-Score	0.32	0.81	0.90	$Wilcoxon: V1 \ vs \ V2 \ p = 0.002$
Response Time (s)	2.1 ± 0.15	3.4 ± 0.20	4.2 ± 0.18	<i>ANOVA</i> : $p < 0.001$
Truthfulness (%)	28	78	88	Expert Verification
Expert Ratings (1-5 scale)				
Relevance	3.8 ± 0.4	4.2 ± 0.3	4.4 ± 0.2	Tukey HSD: all $p < 0.05$
Completeness	3.6 ± 0.5	4.1 ± 0.4	4.3 ± 0.3	Tukey HSD: all p < 0.05
User Satisfaction	4.2 ± 0.3	4.4 ± 0.2	4.5 ± 0.2	F(2,21) = 9.98, p = 0.001

The comprehensive evaluation across nine metrics reveals several key patterns in RAG pipeline performance:

Performance Progression: All metrics show consistent improvement from $V1 \rightarrow V2 \rightarrow V3$, with the most dramatic gains occurring between V1 and V2. The addition of full-text retrieval capability represents the largest single performance leap.

Latency Trade-offs: Response times increased systematically (V1: 2.1s, V2: 3.4s, V3: 4.2s), reflecting the computational cost of enhanced retrieval. However, the 100% latency increase from V1 to V3 delivered a 300% accuracy improvement, indicating favorable performance-cost ratios.

Expert Rating Consistency: Healthcare professionals and IT specialists showed remarkable agreement in their assessments. The narrow standard deviations (± 0.2 to ± 0.5) across Relevance, Completeness, and User Satisfaction ratings demonstrate high inter-rater reliability and confidence in the results.

Threshold Achievement: V3 achieved clinical-grade performance with 90% exact match accuracy and 88% truthfulness—metrics that approach the reliability standards expected in medical decision support systems.

Balanced Improvements: Unlike systems that optimize single metrics, each pipeline version improved simultaneously

across precision, recall, and user experience measures. This balanced enhancement pattern suggests robust architectural improvements rather than parameter fine-tuning.

User Experience Correlation: The progression in objective metrics (precision/recall) closely mirrors subjective expert ratings (relevance/completeness), validating that technical improvements translate meaningfully to end-user experience in biomedical information retrieval tasks.

5.4. Statistical Significance Testing

To determine whether observed performance improvements across pipeline versions are statistically significant, the following tests were conducted on the Exact Match (EM) accuracy and F1-score data from the 10-question evaluation:

5.4.1. Paired t-Tests on Exact Match Accuracy

- 1. V1 vs. V2: Mean difference = 0.50; t(9) = 7.07, p < 0.001, Cohen's d = 2.24 (large effect)
- 2. V2 vs. V3: Mean difference = 0.10; t(9) = 2.37, p = 0.042, Cohen's d = 0.75 (medium effect)

5.4.2. Wilcoxon Signed-Rank Tests on F1-Score

- 1. V1 vs. V2: W = 0, p = 0.002, indicating a significant median increase in F1-score
- 2. V2 vs. V3: W = 5, p = 0.031, indicating a significant, albeit smaller, median improvement

- 5.4.3. One-Way ANOVA on Expert Ratings (Relevance, Completeness. User Satisfaction)
- 1. Relevance: F(2,21) = 16.4, p < 0.001
- 2. Completeness: F(2,21) = 14.2, p < 0.001
- 3. User Satisfaction: F(2,21) = 9.8, p = 0.001 Post-hoc Tukey HSD confirmed that V3 ratings exceeded V1 and V2 (p < 0.05), while V2 also exceeded V1 (p < 0.05).

5.4.4. Inter-Rater Reliability (Cronbach's α)

- 1. Relevance $\alpha = 0.88$ (good)
- 2. Completeness $\alpha = 0.85$ (good)
- 3. User Satisfaction $\alpha = 0.90$ (excellent)

These tests confirm that each enhancement step from abstract-only to full-text to query-refinement yields statistically significant improvements in both objective retrieval performance and subjective expert assessments.

5.5. Error Analysis and Ablation Study

In this section, the researcher examines the types and frequencies of errors across the three pipeline versions and quantifies the impact of key architectural enhancements.

5.5.1. Error Type Distribution

Table 13. Summary of error categories and their prevalence

Error Type	Version 1 (7 errors)	Version 2 (2 errors)	Version 3 (1 error)
Hallucinations	2 (29%)	1 (50%)	0 (0%)
Incompleteness	4 (57%)	1 (50%)	1 (100%)
Irrelevant	1 (1/0/)	0 (09/)	0 (00/)
Content	1 (14%)	0 (0%)	0 (0%)

Version 1 suffered predominantly from incomplete retrievals, leading to missing key facts. Version 2 reduced total errors by 71%. Version 3 achieved only a single error (an incomplete retrieval), with zero hallucinations, indicating robust context handling.

5.5.2. Ablation Study of Retrieval Components

Table 14. An ablation experiment isolated the contributions of full-text access and query refinement

Pipeline Variant	Exact Match Accuracy	Hallucinations	Incomplete Retrievals
Abstract-only (V1)	30%	2	4
+ Full-text retrieval (V2)	80%	1	1
+ Query refinement & thresholds (V3)	90%	0	1

Full-text retrieval alone (V2) contributed a +50% EM boost and an 86% reduction in incomplete retrievals. Query refinement and stricter thresholds (V3) suppressed hallucinations entirely and further improved EM by 10%.

5.5.3. Hallucination Suppression

Version 3's zero hallucinations result from:

- 1. Query Refinement: The System omits answers when the retrieved context confidence is below a set threshold, preventing unsupported generation.
- 2. Context Amplification: Aggregating multiple highsimilarity passages ensures comprehensive evidence for each response.
- Stricter Retrieval Thresholds: Higher similarity cutoffs filter out marginally relevant documents that could trigger model speculation.

Despite perfect performance on 10 queries, edge-case hallucinations may still occur. For clinical deployment, implement:

- 1. Continuous logging of low-confidence generations.
- 2. Automated heuristics to flag potential hallucinations.
- 3. Human-in-the-loop review for flagged cases.

6. Conclusion

This study successfully demonstrates the integration of semantic search with generative AI in the biomedical domain using an RAG framework. By leveraging ChromaDB for efficient vector-based retrieval and Gemini for generative question answering, the system offers a more intuitive and context-aware approach to exploring the vast corpus of PubMed literature. Compared to traditional keyword-based methods, the proposed approach delivers more relevant, precise, and user-friendly responses, effectively bridging the gap between complex scientific texts and user queries. The architecture proves to be a scalable and powerful tool for researchers, clinicians, and academics who require rapid, accurate synthesis of biomedical information.

6.1. Recommendation

To further enhance the system's performance and applicability, several improvements are recommended. First, incorporating domain-specific fine-tuning for the generative model using specialized biomedical corpora such as BioBERT or PubMedQA can significantly improve response accuracy and minimize hallucinations. Additionally, integrating ChromaDB with established biomedical ontologies like MeSH or UMLS would enrich the semantic layer and lead to more relevant and precise information retrieval. Improving the user interface with features such as query refinement options and citation tracking will support a more intuitive and functional user experience. It is also essential to develop a comprehensive evaluation framework based on biomedical QA benchmarks to systematically assess the system's reliability and effectiveness. Lastly, expanding the dataset beyond PubMed to include other biomedical repositories such as ClinicalTrials.gov or PMC will broaden the scope and utility of the system for diverse research needs.

References

- [1] Alex Thomo, *PubMed Retrieval with RAG Techniques*, Digital Health and Informatics Innovations for Sustainable Health Care Systems, IOS Press, vol. 316, pp. 652-653, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [2] Jiawei He et al., "Retrieval-Augmented Generation in Biomedicine: A Survey of Technologies, Datasets, and Clinical Applications," *arXiv Preprint*, pp. 1-29, 2025. [CrossRef] [Google Scholar] [Publisher Link]
- [3] Aritra Kumar Lahiri, and Qinmin Vivian Hu, "AlzheimerRAG: Multimodal Retrieval Augmented Generation for PubMed Articles," *arXiv Preprint*, pp. 1-17, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [4] Simeon Emanuilov, Retrieval Augmented Generation (RAG) Limitations, Meduim, 2024. [Online]. Available: https://medium.com/@simeon.emanuilov/retrieval-augmented-generation-rag-limitations-d0c641d8b627
- [5] Yujie Sun et al., "AI Hallucination: Towards a Comprehensive Classification of Distorted Information in Artificial Intelligence-Generated Content," *Humanities and Social Sciences Communications*, vol. 11, no. 1278, pp. 1-14, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [6] Hanjie Zhang, and Peter Kotanko, "#1506 Uremic Toxicity: Gaining Novel Insights through AI-Driven Literature Review," *Nephrology Dialysis Transplantation*, vol. 39, no. Supplement-1, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [7] Bojana Bašaragin et al., "How Do You Know that? Teaching Generative Language Models to Reference Answers to Biomedical Questions," *Proceedings of the 23rd Workshop on Biomedical Natural Language Processing*, Bangkok, Thailand, pp. 536-547, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [8] Jim Kutz, Leveraging ChromaDB for Vector Embeddings A Comprehensive Guide, Airbyte, 2025. [Online]. Available: https://airbyte.com/data-engineering-resources/chroma-db-vector-embeddings
- [9] Rui Yang et al., "Retrieval-Augmented Generation for Generative Artificial Intelligence in Medicine," *arXiv Preprint*, pp. 1-11, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [10] Mammona Qudisa, and Muhammad Maozam Fraz, "Optimizing Article Screening and Information Extraction: A Hybrid Approach with GeminiAI and Vector Database," 2024 4th International Conference on Digital Futures and Transformative Technologies (ICoDT2), Islamabad, Pakistan, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [11] P. Mozelius, and N. Humble, "On the Use of Generative AI for Literature Reviews: An Exploration of Tools and Techniques," *Proceedings of the 23rd European Conference on Research Methodology for Business and Management Studies*, vol. 23, no. 1, pp. 161-168, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [12] Siru Liu et al., "Improving Large Language Model Applications in Biomedicine with Retrieval-Augmented Generation: A Systematic Review, Meta-Analysis, and Clinical Development Guidelines," *Journal of the American Medical Informatics Association*, vol. 32, no. 4, pp. 605-615, 2025. [CrossRef] [Google Scholar] [Publisher Link]
- [13] Haowen Xu et al., "Automating Bibliometric Analysis with Sentence Transformers and Retrieval-Augmented Generation (RAG): A Pilot Study in Semantic and Contextual Search for Customized Literature Characterization for High-Impact Urban Research," *UrbanAI '24: Proceedings of the 2nd ACM SIGSPATIAL International Workshop on Advances in Urban-AI*, Atlanta, GA, USA, pp. 43-49, 2024. [CrossRef] [Google Scholar] [Publisher Link]
- [14] Mohammed-Khalil Ghali et al., "Enhancing Knowledge Retrieval with In-Context Learning and Semantic Search through Generative AI," Knowledge-Based Systems, vol. 311, pp. 1-13, 2025. [CrossRef] [Google Scholar] [Publisher Link]
- [15] Chung-Chi Huang, and Zhiyong Lu, "Discovering Biomedical Semantic Relations in PubMed Queries for Information Retrieval and Database Curation," *Database*, vol. 2016, pp. 1-15, 2016. [CrossRef] [Google Scholar] [Publisher Link]
- [16] Sardar Mudassar Ali Khan, Iterative Model Used in Software Development, ResearchGate, 2023. [Online]. Available: https://www.researchgate.net/publication/371902841_Iterative_Model_Used_in_Software_Development