



**ARTICLE**

# Efficiency-Driven Custom Chatbot Development: Unleashing LangChain, RAG, and Performance-Optimized LLM Fusion

S. Vidivelli\*, Manikandan Ramachandran\* and A. Dharunbalaji

School of Computing, SASTRA Deemed University, Thanjavur, Tamilnadu, 613401, India

\*Corresponding Authors: S. Vidivelli. Email: vidieng@gmail.com; Manikandan Ramachandran. Email: srmanimt75@gmail.com

Received: 26 May 2024 Accepted: 03 July 2024 Published: 15 August 2024

## ABSTRACT

This exploration acquaints a momentous methodology with custom chatbot improvement that focuses on proficiency close by viability. We accomplish this by joining three key innovations: LangChain, Retrieval Augmented Generation (RAG), and enormous language models (LLMs) tweaked with execution proficient strategies like LoRA and QLoRA. LangChain takes into consideration fastidious fitting of chatbots to explicit purposes, guaranteeing engaged and important collaborations with clients. RAG's web scratching capacities engage these chatbots to get to a tremendous store of data, empowering them to give exhaustive and enlightening reactions to requests. This recovered data is then decisively woven into reaction age utilizing LLMs that have been calibrated with an emphasis on execution productivity. This combination approach offers a triple advantage: further developed viability, upgraded client experience, and extended admittance to data. Chatbots become proficient at taking care of client questions precisely and productively, while instructive and logically pertinent reactions make a more regular and drawing in cooperation for clients. At last, web scratching enables chatbots to address a more extensive assortment of requests by conceding them admittance to a more extensive information base. By digging into the complexities of execution proficient LLM calibrating and underlining the basic job of web-scratched information, this examination offers a critical commitment to propelling custom chatbot plan and execution. The subsequent chatbots feature the monstrous capability of these advancements in making enlightening, easy to understand, and effective conversational specialists, eventually changing the manner in which clients cooperate with chatbots.

## KEYWORDS

LangChain; retrieval augumental generation (RAG); fine tuning

## 1 Introduction

The dental business flourishes with clear correspondence and informed patients. Conventional strategies for patient collaboration, such as calls and appointment arrangements, can be tedious for the two patients and staff. This restricted admittance to promptly accessible data can prompt disappointment and botched open doors for safeguard care. While existing dental chatbots offer a brief look into the eventual fate of patient communication, they frequently miss the mark regarding information profundity, functional productivity, and the capacity to enable patients. This exploration



proposes a clever way to address these difficulties by creating custom dental chatbots focusing on viability and execution proficiency.

Envision a patient encountering unexpected tooth responsiveness and looking for data at night. Current chatbots could offer restricted pre-customized reactions, leaving the patient baffled and without replies. Our exploration proposes an earth-shattering methodology that prepares chatbots to be effective and far-reaching wellsprings of dental data. This exploration uses the force of three key innovations: LangChain, Retrieval Augmental Generation (RAG), and performance-efficient fine-tuned Large Language Models (LLMs).

LangChain permits us to tailor chatbots for explicit functionalities inside dental considerations quickly and fastidiously. By zeroing in on the one-of-a-kind requirements of dental practices and patients, we can make chatbots that address a more extensive scope of requests with more prominent exactness. It guarantees engaged and essential communication, enabling patients to adopt a more proactive strategy for their dental well-being.

RAG's web-scratching capacities engage dental chatbots to access a vast store of dental data from respectable internet-based sources. This enables the chatbot to give exhaustive and enlightening reactions to patient requests, going from essential dental cleanliness tips to additional mind-boggling techniques. In the previously mentioned situation with tooth responsiveness, a chatbot fueled by RAG could get to and combine important data, offering the patient primer direction on relief from discomfort choices while suggesting planning an arrangement for a legitimate conclusion. This engages the patient as well as reduces nervousness and cultivates trust in the dental practice.

The third mainstay of our methodology lies in the execution and effective calibrating of Large Language Models (LLMs). LLMs are strong artificial intelligence models prepared on massive measures of text information, permitting them to create human-quality text in light of many prompts and questions. Be that as it may, customary LLM adjusting can be computationally costly. Our exploration dives into procedures like LoRA and QLoRA, which focus on execution productivity while keeping up with the adequacy of LLM calibrating. This guarantees that the chatbot works productively, limiting reaction times and augmenting patient fulfilment.

By combining these elements, our proposed approach offers a three-pronged benefit for dental practices and patients:

- i) **Further Developed Adequacy and Patient Strengthening:** Chatbots become adroit at taking care of patient inquiries precisely and productively, decreasing the weight on dental staff and smoothing out quiet cooperation. Straightforward requests about booking, protection, or essential dental data can be tended to by the chatbot, saving staff time for additional intricate assignments. All the more critically, patients are engaged to get an abundance of data, permitting them to come to informed conclusions about their dental considerations.
- ii) **Upgraded Patient Experience:** Educational and logically significant reactions make a more regular and drawing in collaboration, encouraging patient trust and fulfilment. Envision a patient looking for data about a particular dental technique. The chatbot, furnished with RAG-scratched data and tweaked LLMs, can clarify the method and answer follow-up inquiries in a characteristic and drawing way. This customized approach fabricates trust and makes patients feel more responsible for their dental well-being.
- iii) **Extended Admittance to Data and Potential Finding Help:** Web scratching enables chatbots to address a more extensive assortment of patient requests, permitting them to respond to essential inquiries, plan arrangements, and even give primer dental data daily. Patients can get to data and get fundamental direction whenever the timing is ideal, eventually prompting a

more educated and proactive way to deal with their dental well-being. In certain circumstances, the chatbot, in light of the patient's side effects, can offer possible findings while stressing the significance of looking for proficient assessment for a conclusive determination and treatment plan.

This exploration digs into the complexities of performance-productive LLM fine-tuning and highlights the primary job of web-scratched information in enhancing chatbot reactions. By synergistically consolidating these innovations, this study offers a critical commitment to propelling custom dental chatbot plans and execution. The subsequent chatbots feature the tremendous capability of these advances in making enlightening, easy-to-understand, and productive conversational specialists, eventually changing how patients cooperate with dental practices and assume responsibility for their oral health. The rest of this paper is organized as follows. [Section 2](#) depicts the current works. [Section 3](#) clarifies the proposed approach, including the combination of LangChain, RAG, and fine-tuned LLMs. [Section 4](#) dives into the advancement cycle of the chatbot, displaying its functionalities and capacities. [Section 5](#) presents the consequences of our assessment, examining the chatbot's performance and client experience.

## 2 Related Work

This task proposes a clever way to deal with dental chatbots that focuses on viability, performance productivity, and thorough dental data access, expanding on existing exploration. We dive into critical examinations that enlighten headways and limits in medical care conversational specialists, emphasising dental chatbots.

Alejandrino et al. in 2023 [1] propose a thorough framework configuration custom-fitted for private dental facilities to robotize key business processes, including dental record storage, charging handling, report generation, and reconciliation of man-made consciousness through chatbot innovation for brilliant arrangement frameworks. Availability to different partners, like patients, dental specialists, and staff, is underscored. The concentrate likewise frames the utilization of the V-Model philosophy for programming advancement and testing, with partner interviews illuminating significant work bundles and assessing time and cost prerequisites.

This overview paper [2] investigates the ramifications of incorporating ChatGPT, a computerized reasoning system for creating text reactions, into different areas inside open dental wellbeing, the scholarly community, and clinical practice. Through a methodical survey system, 84 pertinent papers were recognized, from which 21 examinations were chosen for top-to-bottom investigation. Discoveries show that ChatGPT, as an occurrence of large language models (LLMs), has exhibited adequacy in supporting researchers with logical examination and dental examinations, offering appropriate reactions and saving time for trial and error stages. In any case, difficulties like predisposition in preparing information, degrading of human ability, possible logical wrongdoing, and lawful reproducibility concerns were distinguished. Notwithstanding these difficulties, the remarkable benefit of coordinating ChatGPT into general well-being dentistry rehearses was recognized, with an accentuation on the correlative job of computer-based intelligence innovation close by human skill in the nuanced field of dentistry.

Parviainen et al. in 2022 [3] propose a proactive work examining task-oriented chatbots' implications as partially automated consulting systems on clinical practices and expert-client relationships. It highlights the need for novel approaches in professional ethics to accommodate the widespread deployment of artificial intelligence, potentially revolutionizing decision-making and interactions in

healthcare organizations. The article argues that chatbot implementation amplifies rationality and automation in clinical practice, altering traditional decision-making processes based on epistemic probability and prudence. Through this discussion, the paper contributes valuable insights into the ethical challenges of chatbots within healthcare professional ethics.

In 2021, Nguyen et al. presented an AI-based chatbot to provide students with instant updates on curriculum, admissions procedures, tuition fees, IELTS writing task II scores, and more [4]. The chatbot was developed using Deep Learning models integrated into the Rasa framework, with a proposed rational pipeline for Vietnamese chatbots to ensure optimal accuracy and prevent model overfitting. The model can detect over fifty types of user questions with an impressive accuracy rate of 97.1% on the test set. Implemented on the National Economics University's official admission page on the Facebook platform, the chatbot showcases the potential of AI technology to streamline communication processes. This exploration offers point-by-point rules on building a simulated intelligence chatbot without any preparation, joined by methods that can be applied across various languages and settings. This study [5] presents the development of a university enquiry chatbot using the Rasa framework, which leverages deep learning to provide a responsive and intelligent system for handling common queries in an academic setting and these works [6–10] explore various aspects of chatbot technology, from deep learning frameworks and language models to question-answering systems and domain-specific applications, offering valuable insights into the efficiency-driven evolution of chatbot systems. Drawing motivation from the writing surveyed [11–15], this paper proposes a model which leverages the creative advancements of LangChain, Retrieval Augmental Generation (RAG), and Performance-Enhanced Large Language Model (LLM) Combination. Expanding upon the experiences acquired from past exploration of chatbot advancement and its applications across different spaces, this model plans to improve the productivity and adequacy of chatbot frameworks. By coordinating LangChain for careful customization, RAG for thorough data retrieval, and Performance-Streamlined LLM Combination for further developed reaction generation, the proposed model looks to alter the capacities of chatbots in conveying engaged, significant, and high-performing connections. This comprehensive methodology not only addresses the difficulties distinguished in existing writing but additionally sets another norm for the advancement of canny conversational specialists. Below [Table 1](#) compares the various attention mechanisms in terms of computational complexity as reviewed in references [16–18].

**Table 1:** Attention mechanisms comparison

Attention mechanism	Description	Computational complexity	Impact on performance
Dot-product attention	Computes attention weights using the dot product of query and critical vectors.	$O(n^2 \cdot d)$	Efficient for more minor sequences, perform well on tasks like machine translation (e.g., Transformer architecture).

(Continued)

**Table 1 (continued)**

Attention mechanism	Description	Computational complexity	Impact on performance
Self-attention	Each element in the input sequence attends to all other elements, allowing for modelling dependencies regardless of distance.	$O(n^2 \cdot d)$	This is crucial for tasks requiring long-range dependencies like text generation (e.g., GPT models) and BERT for various NLP tasks.
Low-rank attention	Approximates the attention matrix with a lower-rank matrix to reduce computational complexity.	$O(n \cdot r \cdot d)$	Balances between efficiency and performance apply when a full attention matrix is unnecessary.
Performer attention	Uses a kernelized approximation to self-attention, reducing complexity while preserving performance.	$O(n \cdot d \log d)$	Suitable for very long sequences, showing competitive performance in language modelling tasks with significant contexts.

Large language models (LLMs) find diverse applications. In their 2024 works, Zhang et al. [19] and Peng et al. [20] extensively review the pivotal role of LLMs in enhancing human-robot interaction (HRI) across various contexts. Notably, LLMs empower intricate multimodal input handling mechanisms, reshaping how robots perceive and engage with users.

### 3 Proposed Method and Workflow

Our examination presents an original procedure for building a dental chatbot that leverages the LangChain structure and tweaks large language models (LLMs) for efficient and educational communications. The workflow of the proposed model is diagrammatically represented in Fig. 1. LangChain gives a particular establishment, permitting consistent reconciliation of exchange the board and data retrieval. We explicitly pick Peft and Qlora, pre-prepared LLMs, and adjust them for the dental area to guarantee they comprehend dental phrasing and client questions. This calibrating could include directed learning with an immense corpus of dental text information and support in figuring out how to compensate for educational and precise reactions. LangChain's implicit aim order abilities are utilized to arrange client inquiries (e.g., "side effects of toothache"). A discourse on the board framework inside LangChain directs the discussion in light of the plan, deciding activities like data retrieval or posing explaining inquiries. We intend to incorporate Retrieval-Augmented Generation (RAG) with LangChain to convey cutting-edge and dependable data. RAG empowers the chatbot to create reactions and access and incorporate pertinent data from valid web sources through centred web scratching of respectable dental sites and clinical diaries. Finally, client testing and iterative improvement in light of criticism guarantee the chatbot's viability, ease of use, and data

precision. This philosophy, consolidating LangChain's adaptability, tweaked LLMs, and centred data retrieval, endeavours to make a dental chatbot that engages patients with instructive and easy-to-use encounters.

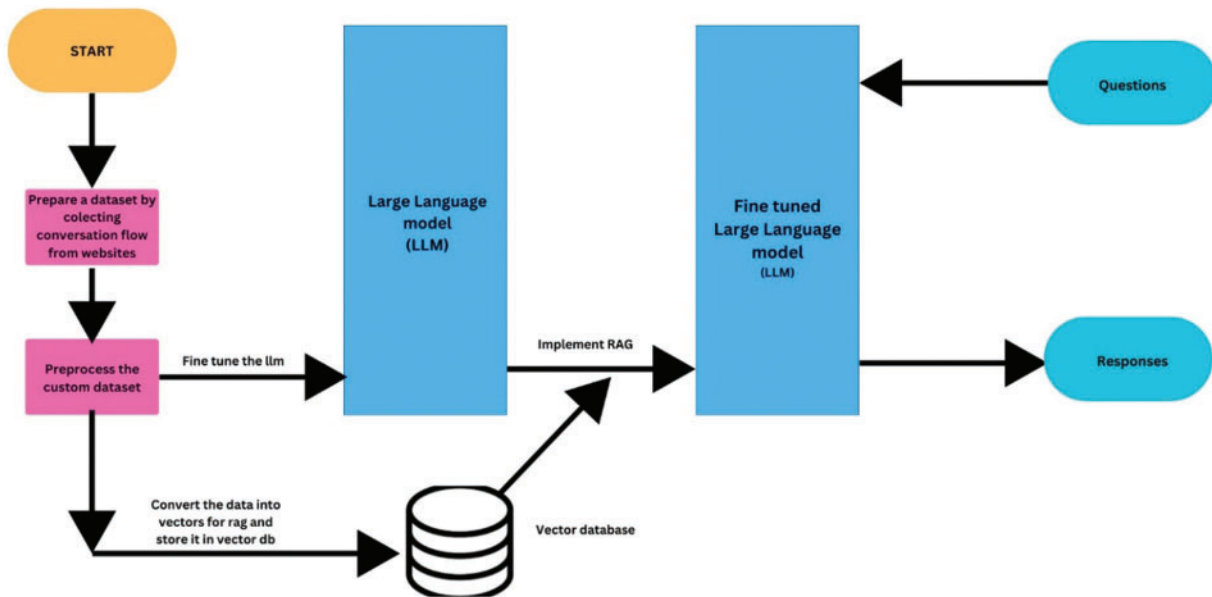


Figure 1: Workflow diagram

### 3.1 Large Language Model

In this project, we use TinyLlama-1.1B-Chat-v1.0 as a prominent language (LLM), a chatbot model built upon the TinyLlama project. The TinyLlama project focuses on creating compact and efficient large language models (LLMs) through pretraining. Here, a 1.1B parameter LLM is fine-tuned for chat conversations. The document details the training process, leveraging synthetic and human-ranked datasets to achieve practical conversational abilities.

Attention mechanisms enable LLMs to selectively focus on relevant parts of the input sequence, capturing long-range dependencies and improving the model's understanding of the context.

The TinyLlama-1.1B-Chat-v1.0 model uses the exact attention mechanism as Llama 2, specifically the scaled dot-product attention. This type of attention is standard in Transformer architectures. It involves calculating attention weights as the dot product of query and key vectors, scaling by square root of the dimension of the key vectors, and applying a softmax function to obtain probabilities.

Scaled dot-product attention has a computational complexity of  $O(n^2 \cdot d)$ , where  $n$  is the sequence length and  $d$  is the dimension of the model. This complexity arises because each element in the sequence must attend to every other element, making it computationally expensive for long sequences.

The TinyLlama model, being a compact version with only 1.1 billion parameters, is optimized for tasks requiring lower computational and memory resources while maintaining good performance on a variety of text generation tasks. Attention mechanisms in TinyLlama-1.1B-Chat-v1.0, particularly self-attention with scaled dot-product attention, play a pivotal role in its ability to handle diverse and complex language tasks effectively. These mechanisms ensure that the model can maintain performance while being computationally efficient, making it suitable for various applications.

### **3.2 Preprocessing Data**

The expanding field of medical care chatbots fueled by large language models (LLMs) presents a groundbreaking chance to upgrade patient training, data access, and, generally speaking, medical care conveyance. Notwithstanding, opening this likely depends on fastidious information arrangement. This exploration investigates the fundamental advances engaged with preprocessing medical care discussion information for LLM preparation, drawing upon the lavita/ChatDoctor-HealthCareMagic-100k dataset on Hugging Face as a source of perspective point. Here, we dig into the thinking behind apparently basic strategies and their effect on the adequacy of the last chatbot model.

Our underlying step includes fragmenting the crude dataset into individual discussion strings. This is usually accomplished by recognizing newline characters or other delimiters that differ in one interview. Separating these discussions permits us to investigate the data stream and recognize the jobs of various members (client and medical services colleagues). Our centre movements pinpoint and extricate explicit text segments in every discussion section. This involves confining the client's feedback (various forms of feedback) and comparing medical services colleague's reactions (answers or guidance). Here, we fastidiously eliminate superfluous names that could mess up the information and obstruct the growing experience.

A fundamental part of information change includes guaranteeing that the model can recognize various speakers inside a discussion. We accomplish this by reformatting the separated human info and right-hand reaction into an organized organization. By utilizing speaker labels and developing the discourse in a particular request, we empower the model to perceive the consecutive idea of the discussion. This organized arrangement works by learning examples of responsive trades essential for viable chatbot collaborations. Genuine discussions are not impeccably organized 100% of the time. The given dataset could contain circumstances where a client's feedback probably does not have a comparing partner reaction. Here, even sections without a total discourse (i.e., just holding back client input) are held. This opens the model to a more practical conversational situation where the client could suggest a conversation starter that might not have a promptly accessible response. By integrating these unpaired fragments, we possibly upgrade the model's flexibility and capacity to deal with unexpected circumstances during certifiable patient collaborations.

### **3.3 Tokenization**

Following preprocessing, the text information goes through an interaction called tokenization. Here, the content separates the whole discussion (counting human information, aide reaction, and speaker labels) into more modest units called tokens. These tokens can be individual words or subword units, contingent upon the particular model engineering. Tokenization permits the LLM to process and comprehend the text information in a manner that works with learning and language generation.

By carefully applying these preprocessing strategies, medical services discussion information can be changed into an organization that engages the LLM to learn successfully. This prepares for improving a vigorous and enlightening medical care chatbot equipped to participate in significant discussions with patients. As analysts in this field, we should ceaselessly assess the adequacy of these methods and investigate imaginative ways to deal with enhanced information groundwork for building the up-and-coming generation of medical care chatbots.

#### **3.3.1 Fine-Tuning LLM**

This exploration project examines the utilization of tweaking methods to upgrade the capacities of a large language model (LLM) in medical care chatbots. Our essential goal is to address the test

of excellent client inquiries that fall outside the domain of the LLM's pre-prepared information base and could prompt mistaken or pointless reactions. By leveraging tweaking, we plan to outfit the LLM with the capacity to deal with these exceptional cases successfully, at last working on the quality and dependability of the chatbot's cooperation with clients.

The pre-prepared LLM utilized in this undertaking fills in as an establishment. After that, we construct space-explicit skills. While pre-prepared LLMs like Tiny Llama 1.1B Visit from Hugging Face offer extraordinary language understanding abilities, they could battle with medical services explicit wording, subtleties, and excellent client inquiries. Tweaking overcomes this issue by giving further preparation on an engaged dataset custom-fitted to the medical services space. This dataset is pivotal in moulding the LLM's capacity to understand medical services-related ideas and precisely answer client requests that veer off of its pre-prepared information.

The dataset picked from Hugging Face is a rich asset for calibrating. Here, the particular subtleties of the dataset are not significant for the more extensive exploration concentrate. However, it is essential to recognize that the information is tokenized. Tokenization alludes to the most common way of separating text into more modest units, frequently individual words or subword units. This pre-handling step is fundamental for the LLM to comprehend and deal with the data inside the dataset during tweaking. The tokenized design permits the LLM to gain proficiency with the particular examples and connections between words ordinarily experienced in medical care discussions.

### *3.3.2 Optimizing Performance: Quantization and Performance-Efficient Techniques*

While fine-Tuning empowers the LLM to handle exceptional healthcare queries, we must also consider the computational demands of this process. We implemented two key strategies to address this: quantization and performance-efficient fine-tuning techniques. Quantization reduces the precision of the LLM's internal parameters, leading to a smaller model footprint and faster inference (response generation) during deployment. This optimization technique enables the fine-tuned LLM to operate more efficiently in resource-constrained environments, making it more practical for real-world chatbot applications.

We specifically employed LoRA (Low-Rank Adapters) and Q-LoRA (Quantized Low-Rank Adapters) within performance-efficient fine-tuning techniques. These techniques focus on supervised fine-tuning, where the LLM learns from a dataset of labelled examples. The dataset likely consists of user queries (inputs) and corresponding informative and accurate healthcare responses (outputs). By incorporating LoRA or Q-LoRA as shown in [Fig. 2](#) during fine-tuning, we can significantly reduce computational cost while maintaining a desirable level of accuracy in the LLM's responses to exceptional healthcare inquiries.

## **3.4 Retrieval Augmented Generation with LangChain**

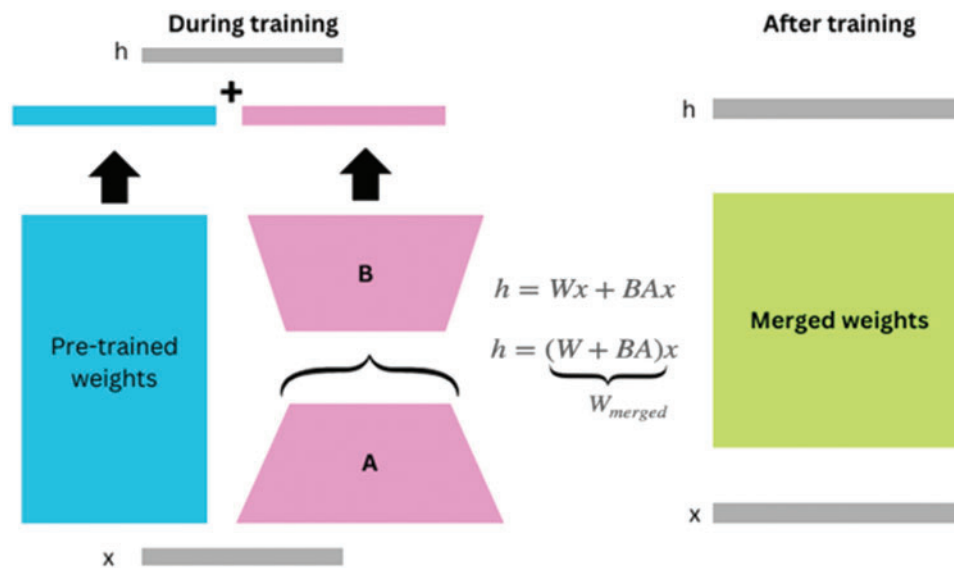
In the previous section, we explored the process of fine-tuning a large language model (LLM) to enhance its capabilities within a specific domain. While fine-tuning equips the LLM with domain-specific knowledge, exceptional user queries inevitably fall outside its knowledge base. This is where Retrieval-Augmented Generation (RAG) emerges as a powerful complementary technique.

- i) Expanding the Knowledge Horizon: Retrieval-Augmented Generation.

RAG operates under a two-stage paradigm that extends the capabilities established through fine-tuning. Here is how RAG builds upon the foundation laid by fine-tuning.



- ii) Leveraging Fine-Tuned Information: The fine-tuned LLM fills in as a center part inside the RAG structure. The space-explicit information gained through fine-tuning permits the LLM to grasp the client's question more and decipher the recovered data from the outer information base in a logically necessary way.



**Figure 2:** LoRA algorithm

Getting to Outside Information: RAG acquaints the vital capacity with access and leverage data from an outer information base. This information base can be fastidiously arranged and customized to the particular area of the chatbot. For example, in our venture zeroed in on medical services chatbots, the information base could envelop clinical archives, research papers, and online assets pertinent to medical care points.

By joining fine-tuning and RAG, we accomplish a synergistic impact that engages the chatbot to convey outstanding performance:

- i) Tending to Remarkable Inquiries: Fine-tuning reinforces the LLM's grasping inside a particular space. RAG permits it to get to outer information for uncommon questions outside its pre-prepared or fine-tuned information base. This consolidated methodology brings a more vigorous and flexible chatbot fit for dealing with a more extensive scope of client requests.
- ii) Upgraded Precision and Factuality: Fine-tuning assists the LLM with learning area explicit wording and accurate data. RAG further improves precision by permitting the chatbot to get to and coordinate modern data from the outside information base. This joined methodology guarantees that the chatbot's reactions are grounded in dependable and exact data, a fundamental viewpoint for spaces like medical services.

Further Developed Client Fulfillment: By giving complete and instructive reactions that address both average and outstanding inquiries, the mix of fine-tuning and RAG prompts a fulfilling client experience. Clients can trust the chatbot to convey precise and enlightening reactions, in any event, for mind-boggling or surprising inquiries. The steps involved in RAG Index Processing are shown in Algorithm 1.

**Algorithm 1:** Algorithm for an LLM application using RAG

## Index Process

```

1: embeddings ← load("embedding_model")
2: doc ← load("file_name")
3: c ← chunker.chunk(doc)
4: ce ← embeddings.embed(c)
5: db ← index(ce)

```

## Query Process

```

1: INITIALIZE sys_prompt ← "You are an AI..."
2: INITIALIZE model ← load("llm_model")
3: while TRUE do
4: q ← get.user_query()
5: qe ← embeddings.embed(q)
6: chunks ← db.search(qe)
7: context ← merge(chunks)
8: prompt ← create_prompt(sys_prompt, q, context)
9: answer ← model.generate(prompt)
10: end while

```

**3.5 Executing RAG with LangChain**

LangChain, a flexible library for building NLP pipelines, works with the consistent reconciliation of RAG inside our fine-tuned LLM structure. We can leverage Langchain's abilities to:

- i) **Load and Preprocess the Outer Information Base:** LangChain gives apparatuses to deal with different information designs, including PDFs. This permits us to efficiently stack the information base of the arranged medical services into the framework and preprocess the substance for robust retrieval.
- ii) **Produce Thick Portrayals:** LangChain coordinates with libraries like Sentence Transformers to make thick vector portrayals for the client's inquiry and the information sources inside the outside information base. These portrayals catch the central semantic importance of the text, empowering efficient retrieval of significant data.
- iii) **Retrieve and Rank Information:** LangChain considers joining with retrieval libraries like Faiss or Elasticsearch. These libraries work with the ID and positioning of the most significant information sources from the outer information base because of the client's question and the thick vector portrayals.
- iv) **Generate Informed Responses:** The positioned rundown of recovered information sources and the client's inquiry are passed to the fine-tuned LLM inside LangChain. The LLM leverages this data to form an extensive reaction to the client's inquiry by integrating pertinent subtleties from the outer information base.

Fine-tuning and RAG address reciprocal methods that work as one to make a more hearty and enlightening chatbot experience. By joining the space, explicit information is obtained through fine-tuning with the capacity to access and leverage outside data utilizing RAG. We engage chatbots to deal with a more extensive scope of client questions with upgraded precision and client fulfilment. This consolidated methodology makes ready for the advancement of savvy chatbots that are fit to participate in significant and enlightening discussions across areas and this sequence of action is shown in [Fig. 3](#).

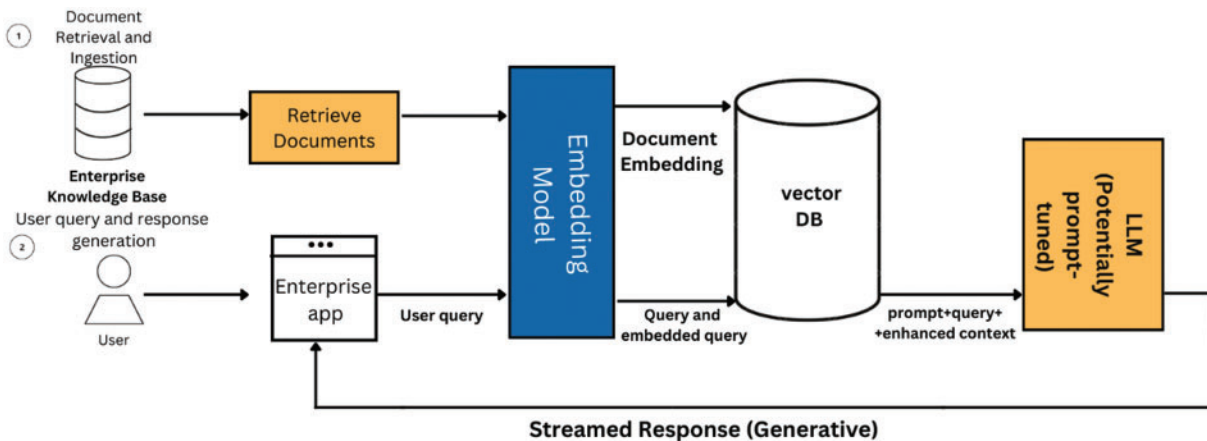


Figure 3: RAG sequence diagram

### 3.6 ChromaDB as a Vector Database

Our venture investigates using ChromaDB, a strong vector information base, inside the RAG setting for chatbot improvement. ChromaDB spends significant time putting away and overseeing thick vector portrayals of text information. These portrayals catch the central semantic significance of text in a compacted design, empowering efficient likeness examinations during retrieval. Strategies like word and sentence embeddings can create these vector portrayals. ChromaDB is intended to deal with large volumes of information, making it appropriate for putting away broad information bases applicable to explicit areas. This versatility guarantees that the framework can efficiently handle the data expected to address various client questions.

By coordinating ChromaDB into the RAG system, we accomplish a few benefits. ChromaDB's efficient vector storage and retrieval systems empower fast ID of the most critical data inside the information base given the client's question. This guarantees that the chatbot recovers the most appropriate information sources to illuminate its reaction generation. ChromaDB succeeds at semantic comparability search, going past straightforward catchphrase coordinating. It leverages the semantic connections caught inside the thick vector portrayals to recognize information sources with comparative importance to the client's inquiry, regardless of whether they contain specific catchphrases. It improves the precision and pertinence of the recovered data. ChromaDB can be flawlessly incorporated with LangChain, a flexible library for building NLP pipelines. LangChain interfaces different NLP parts, considering a smooth joining of ChromaDB inside the RAG system.

The work process with ChromaDB includes preprocessing the outer information base. Message information is changed into thick vector portrayals utilizing strategies like word embeddings or sentence embeddings. ChromaDB then, at that point, stores these portrayals alongside their related metadata. When a client presents an inquiry, it is likewise changed into a thick vector portrayal utilizing similar methods utilized for the information base. LangChain connects with ChromaDB, using the client's question vector portrayal to distinguish the most crucial information sources inside the information base. ChromaDB performs a semantic closeness search, recovering information sources with high vector comparability to the client's question. The recovered information sources are positioned in light of pertinence scores, and this data, alongside the client's question, is passed to the fine-tuned LLM inside LangChain. The LLM leverages this data to figure out a thorough reaction that tends to the client's inquiry.

### ***3.7 Pruning***

Pruning aims to reduce number of parameters in a model by removing redundant or insignificant connections. This technique can significantly decrease model size and improve inference speed while maintaining comparable accuracy.

Pruning can be applied to various layers within LLMs and LMMs. Structured pruning techniques targeting specific network parts (e.g., removing weights with low magnitudes) can be particularly effective.

### ***3.8 Quantization***

Quantization reduces the precision of model weights and activations, typically from 32-bit floating-point numbers to lower precision formats like 8-bit integers. This significantly reduces memory footprint and can accelerate inference on resource-constrained devices.

Quantization has been successfully applied to various LLM architectures, often with minimal performance degradation. Quantization techniques like post-training quantization and adaptive quantization can be effective for LLMs.

### ***3.9 Knowledge Distillation***

Knowledge distillation involves transferring knowledge from a complex, pre-trained teacher to a smaller student model. This allows the student model to achieve performance closer to the teacher model with a smaller footprint.

Knowledge distillation is a promising technique for compressing large pre-trained LLMs into smaller, more efficient models. Techniques like attention distillation and intermediate layer distillation can be used to transfer knowledge effectively.

### ***3.10 Transfer Learning***

This approach leverages knowledge from a pre-trained model on a large, general-purpose dataset (source task) and applies it to a new, specific task (target task). Using transfer learning techniques, pre-trained LLMs like GPT-3 and LMMs like VL-BERT can be fine-tuned on target tasks. This allows them to leverage their learned representations for new domains while adapting to specific task requirements.

### ***3.11 Few-Shot Learning***

This technique aims to train models to perform well on new tasks with only a limited amount of labelled data for the target task. Few-shot learning techniques like prompt-based learning and adaptation algorithms are being explored for LLMs. These techniques allow the model to learn from a few examples in the target domain and adapt its behaviour accordingly.

## **4 Experiment**

Our examination investigates the capability of consolidating Boundary Efficient Fine-Tuning (PEFT) strategies and quantization to make an asset-efficient and exact medical care chatbot inside the limits of a free Google Colab environment. We further explore the utilization of Retrieval-Augmented Generation (RAG) with LangChain to upgrade the chatbot's capacity to address client inquiries by leveraging an outside knowledge base.

#### ***4.1 Fine-Tuning with PEFT for Resource Efficiency***

Since medical services chatbots expect admittance to a tremendous measure of space explicit data, pre-prepared large language models (LLMs) are a natural beginning stage. However, fine-tuning these LLMs on a curated healthcare dataset within resource-constrained environments like Google Colab's free tier necessitates an efficient approach. This is where PEFT techniques come into play.

PEFT allows us to fine-tune the LLM with minimal parameter updates, minimizing the computational cost and memory footprint. We leverage two essential PEFT techniques for this project:

- i) LoRA (Low-Rank Adaptation): This technique only updates a lower-rank matrix instead of the entire model weight matrix. Imagine the model's weights as a large table. LoRA creates a smaller, "lower-rank" table that captures the essential information for fine-tuning. This significantly reduces the computational cost and memory footprint during fine-tuning. We configure LoRA with a rank matrix value of 64, balancing efficiency and performance.
- ii) QLoRA (Quantized LoRA): Building upon LoRA, QLoRA further optimizes the process by quantizing the LoRA rank matrix. Quantization refers to reducing the precision of the data used within the model. We employ the Bitsandbytes library to reduce the model weights from 32-bit floating-point precision to a more memory-efficient 16-bit format after applying LoRA. This technique significantly reduces storage requirements and facilitates training within the limited resources of Google Colab.

#### ***4.2 Loading the Model***

Unleashing the potential of large language models (LLMs) for chatbot development necessitates careful consideration of resource constraints, particularly in environments like Google Colab's free tier. This research explores strategic model loading and configuration techniques tailored for the Tiny Llama 1.1B Chat v1.0 model from Hugging Face.

The Transformers library fills in as an integral asset for leveraging pre-prepared LLMs. The `AutoModelForCausalLM.from_pretrained` function acts as the gateway, allowing us to load the Tiny Llama model's architecture simply by providing its name ("TinyLlama/TinyLlama-1.1B-Chat-v1.0"). However, resource efficiency remains paramount.

To this end, we present two essential methods: quantization and reserving deactivation. Quantization is a craft of memory streamlining, where model loads are addressed with fewer pieces, altogether decreasing memory impression without forfeiting significant performance. This can be accomplished by determining a pre-defined quantization configuration (e.g., `bnb_config`) during model stacking. Also, we decisively deactivate model reserving during preparation by setting the `model.config.use_cache = False`. While reserving can speed up calculations, it frequently comes at the expense of expanded memory use. Handicapping, it adjusts impeccably with our objective of asset-efficient fine-tuning inside compelled conditions.

Overcoming issues among human and LLM correspondence lies in the tokenizer, a vital part of interpreting messages into mathematical portrayals that the LLM can comprehend. The Transformers library offers the `AutoTokenizer.from_pretrained` capability to flawlessly stack the tokenizer related to the picked LLM, guaranteeing a smooth language processing experience.

Moreover, we quickly and insidiously configure the tokenizer's padding conduct to guarantee steady treatment of arrangements with shifting lengths. To agree on padding and end-of-sentence tokens, we lay out a `tokenizer.pad_token = tokenizer.eos_token`. This cultivates a smoother language

processing experience. At long last, we decisively configure `tokenizer.padding_side = "right"`, educating the model to add padding tokens to the furthest limit of arrangements, saving the natural progression of human language.

We streamline model stacking and memory for the Tiny Llama 1.1B Talk v1.0 model executives by carefully applying quantisation, reserving deactivation, and mindful tokeniser configuration. These methods engage the advancement of asset-efficient chatbots, even inside conditions with restricted computational assets. This prepares for a more extensive reception of cutting-edge conversational artificial intelligence and encourages advancement across different fields.

### ***4.3 Setting Up the PEFT Configuration and Fine-Tuning Process***

To successfully leverage LoRA, we configure the accompanying settings inside the PEFT system:

We set the configuration “with no inclination,” showing that the leading weight network, not the predisposition terms (extra qualities that impact the result of a neuron), refreshed through LoRA.

This further reduces the amount of data that needs to be processed during fine-tuning.

**Task Type:** We assign the task type as “causal LM” (language modelling), signifying that the model predicts next word in a sequence based on preceding context. This is a crucial setting as it informs the fine-tuning process of updating the model parameters for generating human-like conversational responses within the healthcare domain.

After establishing the PEFT configuration, we initialize the training arguments for the fine-tuning process. A crucial parameter here is “gradient checkpoint,” set to “True.” Gradient checkpointing is a memory-saving technique that discards intermediate activations during backpropagation, a critical step in training neural networks. We significantly reduce memory usage during fine-tuning by recomputing these activations only when needed. The optimizer is “paged\_adamw\_32bit,” a memory-efficient variant of the Adam optimizer, a popular optimization algorithm used to train neural networks.

We then fine-tuned the chosen LLM, specifically the Tiny Llama 1.1B Chat model from the Hugging Face library, using the Soft-Masked Language Modeling (SFT) trainer. SFT is a training objective well-suited for causal language modelling tasks, where the model learns to predict the next word in a sequence based on preceding context. This training objective aligns perfectly with our goal of enabling the LLM to generate informative and coherent responses to user queries within the healthcare domain.

While the fine-tuning process reduces the training loss, indicating that the model is learning from the healthcare dataset, the initial evaluation of the fine-tuned model reveals performance below expectations. This highlights the limitations of fine-tuning alone in addressing complex information needs within the healthcare domain. Healthcare information can be highly nuanced, and relying solely on the LLM’s internal knowledge might not be sufficient to provide users with comprehensive and accurate responses, particularly for specific or exceptional queries.

### ***4.4 Enhancing Accuracy with Retrieval-Augmented Generation (RAG) Using Lag Chain***

Our research leverages LangChain, a versatile library for constructing intelligent systems, to establish the core framework for our healthcare chatbot pipeline. By capitalizing on the pipeline function within Transformers, we construct a dedicated pipeline designed explicitly for text generation. This pipeline is configured with the pre-loaded LLM and tokenizer, ensuring seamless integration and communication between these critical components. To manage the length and intricacy of the chatbot’s

reactions, we carry out a most significant generation length boundary, regularly set around 200 tokens. This guarantees succinct and centred reactions that focus on conveying the most pertinent data inside an edible configuration.

Past the centre LLM usefulness, we dive into the domain of brief designing, a vital procedure for directing the LLM towards producing instructive and relevant reactions lined up with client questions. We quickly and odiously plan a brief format that consolidates the client’s underlying inquiry and is adaptable to coordinate any current responses or recover data pertinent to the ongoing discourse. Langchain’s PromptTemplate class demonstrates the importance of working with this cycle, offering an organized and efficient way to deal with creating successful prompts.

We lay out a two-dimensional data retrieval system to engage the chatbot with admittance to a rich pool of medical care-related information. The initial step includes using Langchain’s TextLoader class. This class fills in as a workhorse, empowering us to stack an exhaustive medical care-related genuine data dataset from an assigned text document (e.g., “/content/fact.txt”). Nonetheless, efficiently processing and overseeing large datasets is fundamental, especially inside asset-obliged conditions. To address this test, we present Langchain’s CharacterTextSplitter class. This incredible asset lumps the information into more modest, reasonable fragments. This piecing approach upgrades processing speed and decreases memory requests, guaranteeing efficient usage of accessible assets.

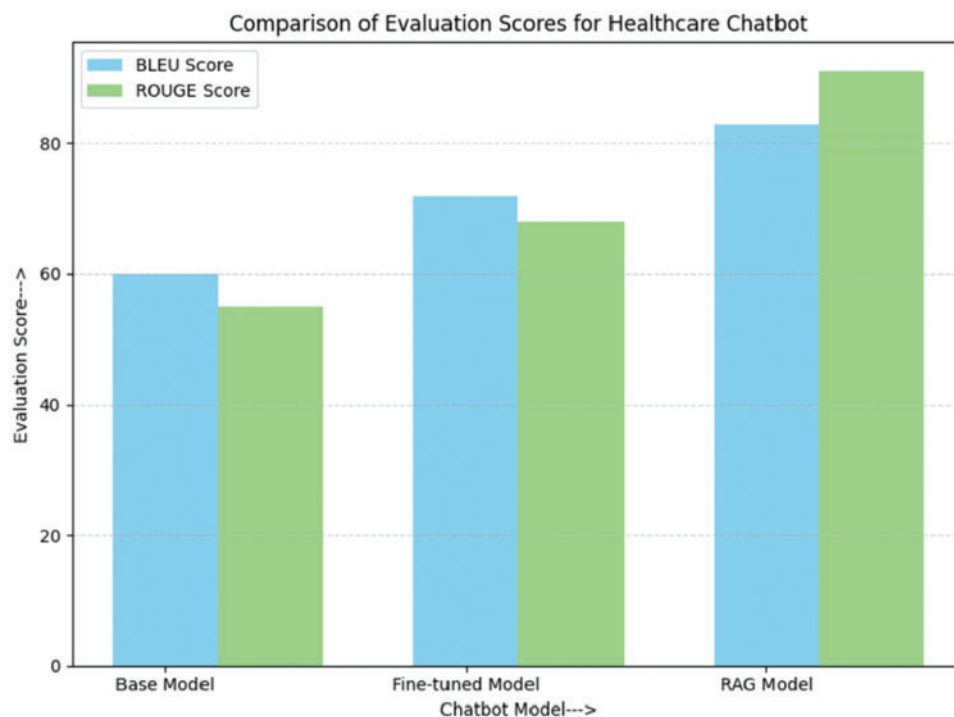
The second step of our data retrieval technique centres around inserting and retrieval. Langchain’s OpenAIEmbeddings class assumes a significant part in this stage. This class engages us in creating excellent embeddings for the client’s question and the stacked archives. Embeddings are mathematical portrayals that catch the semantic significance and setting of the text inside a high-layered space. By leveraging Langchain’s Chroma class, we build a hearty vector store. This vector store is a record for the archive embeddings, considering efficient retrieval in light of the likeness look. When a client presents a question, the framework inserts the inquiry. The Chroma vector store then retrieves reports inside the assortment that display the most significant semantic similitude to the client’s inquiry. By recovering the records firmly aligned with the client’s goal, the chatbot accesses a rich pool of important data that can be consistently integrated into its reactions, cultivating instructive and exhaustive correspondence. This undertaking examined a clinical benefits chatbot’s viability through different assessment techniques. We looked at three models: a base model, a calibrated model, and a Retrieval-Augmented Generation (RAG) execution.

To survey the chatbot’s reactions, we utilized BLEU and ROUGE scores, which measure likeness to human-composed references, as shown in [Table 2](#). The outcomes showed a reasonable benefit for calibrating and RAG execution. The calibrated model accomplished a massive improvement over the base model, showing the viability of space explicit preparation in creating exact and important clinical reactions. The RAG-carried-out model showed an extra increase, recommending that admittance to outside information sources further improves education and authentic establishment.

**Table 2:** Comparison of evaluation scores

Model type	BLEU score (%)	ROUGE score (%)
Base model	60	55
Fine-tuned model	72	68
RAG implemented model	83	91

A client fulfilment overview yielded positive input. Clients found the chatbot's reactions clear, practical, and reasonably precise. Curiously, client criticism featured the chatbot's actual capacity for starting information assortment, yet in addition, communicated a craving for more nuanced reactions in complex circumstances. The assessment results generally show the commitment to the clinical benefits of chatbots. By leveraging space explicit preparation and creative structures like RAG, these chatbots can furnish clients with an essential and enlightening experience. The evaluation score is computed for base model, fine-tuned model and RAG implemented model and result is tabulated in Table 2 and graphically represented in Fig. 4.



**Figure 4:** Graphical representation of evaluation scores

### *Inference Time*

Inference time, also known as latency, refers to the duration taken by a large language model (LLM) to process input data and generate responses during inference or prediction. This metric becomes especially critical in real-world applications where the LLM must promptly address user queries or handle substantial data volumes in real time. Critical considerations for estimating LLM inference time propose a practical framework for evaluation.

#### *Metrics for Inference Time*

- i) Time To First Token (TTFT): This metric measures the time from receiving a prompt to generating the first output token. It includes model initialization and processing overhead.
- ii) Time Per Output Token (TPOT): TPOT represents the time required to generate each subsequent output token. Faster TPOT translates to quicker responses.
- iii) Latency: Overall response latency combines TTFT and TPOT, accounting for the entire response generation process.



$$Latency = TTFT + (TPOT \times N) \tag{1}$$

$N$  is the total number of tokens generated in the response. Based on this Eq. (1), we have compared the two models' performance in terms of latency. And the RAG model performed well in terms of Inference Time and GPU and RAM Usage, its comparison result is shown in Figs. 5 and 6.

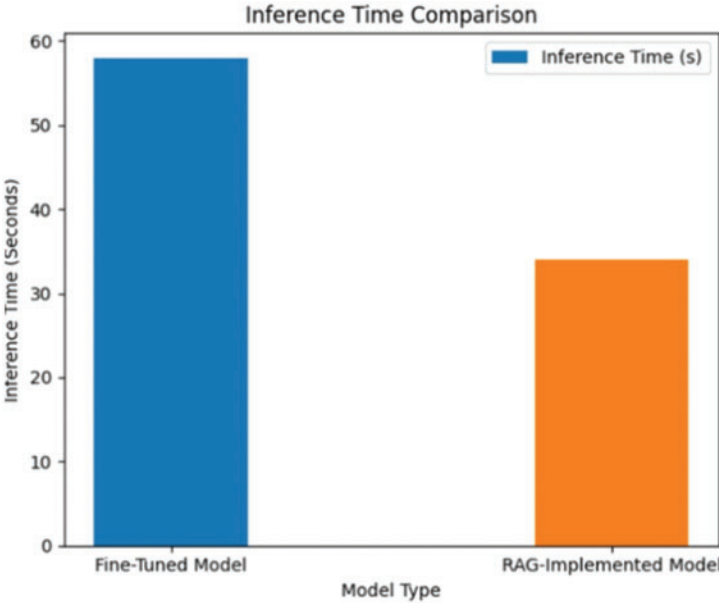


Figure 5: Inference time comparison

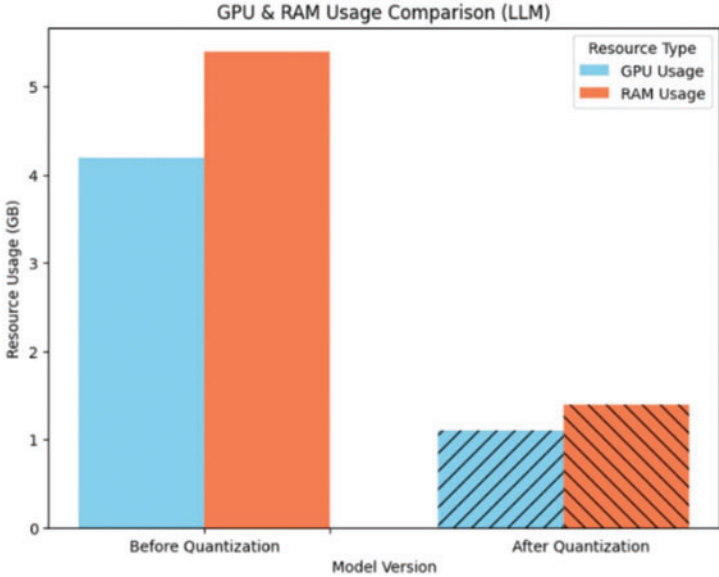


Figure 6: GPU and RAM usage comparison

## 5 Conclusion and Future Work

This examination has investigated the improvement of an asset-efficient medical services chatbot leveraging large language models (LLMs), boundary-efficient fine-tuning (PEFT), and efficient information processing strategies. The LangChain library was a powerful establishment for developing the chatbot pipeline, while key brief designing strategies directed the LLM towards producing educational and important reactions to client inquiries. A two-dimensional data retrieval methodology was used to engage the chatbot with admittance to a rich vault of medical services information, including record stacking and piecing close by installing and retrieval utilizing vector stores. The assessment interaction utilized a multifaceted methodology incorporating automatic measurements like BLEU and ROUGE scores, emphasizing accurate precision, client fulfilment studies, and evaluations by clinical experts. The outcomes were encouraging, exhibiting the chatbot's capacity as a significant instrument for patient instruction and introductory data arrangement inside asset-obliged conditions.

This examination offers critical commitments to the field of medical care chatbots. It can leverage PEFT and quantization strategies to develop an LLM-controlled chatbot that works efficiently even in conditions with restricted computational assets. Furthermore, the essential data retrieval approach guarantees admittance to significant medical services information, cultivating valuable and extensive reactions. Lastly, the assessment philosophy provides a vital structure for surveying the viability of medical care chatbots, stressing both client experience and genuine exactness of the data conveyed.

While the ebb and flow research establishes areas of strength, there are a few promising roads for future investigation to improve the capacities and viability of the LLM-fueled medical services chatbot. One urgent area of the center lies in persistently working on the precision and verifiable rightness of the chatbot's reactions. Here, executing progressed retrieval-augmented generation (RAG) techniques hold gigantic potential. RAG approaches include recovering significant reports from an information base during the generation interaction. The LLM can create more grounded and educational reactions by consolidating these recovered records as an extra setting, especially while managing complex medical care questions. A few procedures can be investigated inside the domain of RAG strategies, for example, thick section retrieval or verifiable language models, to recognize the most significant and dependable data hotspots for expanding the generation cycle. Another region for future work includes consolidating instruments for handling vulnerability and disambiguation. Medical services are a nuanced space, and client questions may sometimes be equivocal or need clearness. The chatbot can be upgraded by incorporating strategies for recognizing such vulnerabilities and giving fitting reactions. This could include offering numerous potential understandings of the inquiry, inciting the client to explain, or guiding them to additional solid hotspots for additional data. Moreover, the chatbot could be prepared to feature the constraints of its capacities and underline the significance of talking with a certified clinical expert for conclusion and treatment choices.

Moreover, investigating UI plans and natural language collaboration (NLI) procedures can improve the client experience. An instinctive and easy-to-use interface can direct clients in forming their questions. At the same time, cutting-edge NLI strategies can empower the chatbot to participate in additional natural and conversational communications. This could include consolidating feeling investigation to measure client feelings and designer reactions as needed or utilizing methods for undivided attention and explanation solicitations to guarantee a reasonable comprehension of the client's goal. Lastly, continuous observation and assessment of the conveyed chatbot are significant for guaranteeing its viability and importance. By gathering client criticism and breaking down chatbot performance measurements, we can distinguish regions for development and refine the model, preparing information and data retrieval methodologies.

Moreover, keeping up to date with progressions in LLM designs, PEFT strategies, and information base development educates future emphasis regarding the chatbot, guaranteeing it stays at the very front of medical services data conveyance. This exploration has fostered an asset-efficient, LLM-fueled medical services chatbot with promising applications, intolerant training, and the beginning of data arrangement. By constantly endeavoring to work on precision, handle vulnerabilities, and improve client connection, this chatbot can change the scene of medical care data access, engaging patients to settle on informed choices concerning their prosperity.

**Acknowledgement:** The authors acknowledge the SASTRA Deemed University, for providing resources to complete this article.

**Funding Statement:** The authors received no specific funding for this study.

**Author Contributions:** The authors confirm their contribution to the paper as follows: Study conception and design: S. Vidivelli, A. Dharunbalaji; Data collection: S. Vidivelli; Analysis and interpretation of results: Manikandan Ramachandran; Draft manuscript preparation: S. Vidivelli and A. Dharunbalaji. All authors reviewed the results and approved the final version of the manuscript.

**Availability of Data and Materials:** The authors confirm that the data supporting the findings of this study are available within the article.

**Conflicts of Interest:** The authors declare that there are no conflicts of interest to report regarding the present study.

## References

- [1] J. C. Alejandrino and E. L. P. Pajota, "An information system for private dental clinic with integration of chatbot system: A project development plan," *Int. J. Adv. Trends Comput. Sci. Eng.*, vol. 12, no. 2, pp. 38–43, 2023.
- [2] A. Tiwari *et al.*, "Implications of ChatGPT in public health dentistry: A systematic review," *Cureus*, vol. 15, no. 6, pp. e40367, 2023. doi: [10.7759/cureus.40367](https://doi.org/10.7759/cureus.40367).
- [3] J. Parviainen and J. Rantala, "Chatbot breakthrough in the 2020s? An ethical reflection on the trend of automated consultations in health care," *Med., Health Care Philos.*, vol. 25, no. 1, pp. 61–71, 2022. doi: [10.1007/s11019-021-10049-w](https://doi.org/10.1007/s11019-021-10049-w).
- [4] T. T. Nguyen, A. D. Le, H. T. Hoang, and T. Nguyen, "NEU-chatbot: Chatbot for admission of national economics university," *Comput. Educ.: Artif. Intell.*, vol. 2, no. 1, pp. 100036, 2021. doi: [10.1016/j.caeai.2021.100036](https://doi.org/10.1016/j.caeai.2021.100036).
- [5] Y. Windiatmoko, R. Rahmadi, and A. F. Hidayatullah, "Developing facebook chatbot based on deep learning using rasa framework for university enquiries," in *IOP Conf. Ser.: Mater. Sci. Eng.*, vol. 1077, pp. 012060. doi: [10.1088/1757-899X/1077/1/012060](https://doi.org/10.1088/1757-899X/1077/1/012060).
- [6] C. Segura, À. Palau, J. Luque, M. R. Costa-Jussà, and R. E. Banchs, "Chatbol, a chatbot for the Spanish "La Liga"," in *9th Int. Workshop Spoken Dialog. Syst. Technol.*, Singapore, Springer, 2019, pp. 319–330.
- [7] M. Rana, "EagleBot: A chatbot based multi-tier question answering system for retrieving answers from heterogeneous sources using BERT," Electronic theses dissertations, Georgia Southern Univ., USA, 2019.
- [8] C. Wang, J. Yan, W. Zhang, and J. Huang, "Towards better parameter-efficient fine-tuning for large language models: A position paper," arXiv preprint arXiv:2311.13126, 2023.
- [9] O. Topsakal and T. C. Akinci, "Creating large language model applications utilizing Langchain: A primer on developing LLM apps fast," *Int. Conf. Appl. Eng. Natural Sci.*, vol. 1, no. 1, pp. 1050–1056, Jul. 2023. doi: [10.59287/icaens.1127](https://doi.org/10.59287/icaens.1127).

- [10] Y. Gao *et al.*, “Retrieval-augmented generation for large language models: A survey,” arXiv preprint arXiv: 2312.10997, 2023.
- [11] D. Xie *et al.*, “Impact of large language models on generating software specifications,” arXiv preprint arXiv: 2306.03324, 2023.
- [12] V. H. Nguyen, R. Gaizauskas, and K. Humphreys, “A combined IR-NLP approach to question answering against large text collections,” in *Proc. Australas. Document Comput. Symp.*, 2016.
- [13] P. Lewis *et al.*, “Retrieval-augmented generation for knowledge-intensive NLP tasks,” *Adv. Neural Inf. Process. Syst.*, vol. 33, pp. 9459–9474, 2020.
- [14] R. Pappagari, P. Zelasko, J. Villalba, Y. Carmiel, and N. Dehak, “Hierarchical transformers for long document classification,” in *2019 IEEE Automatic Speech Recognit. Understand. Workshop (ASRU)*, IEEE, Dec. 2019, pp. 838–844.
- [15] S. S. Manathunga and Y. A. Illangasekara, “Retrieval augmented generation and representative vector summarization for large unstructured textual data in medical education,” arXiv preprint arXiv:2308.00479, 2023.
- [16] A. Vaswani *et al.*, “Attention is all you need,” *Adv. Neural Inf. Process Syst.*, vol. 30, pp. 1–15, 2017.
- [17] M. H. Guo *et al.*, “Attention mechanisms in computer vision: A survey,” *Comput. Vis. Media*, vol. 8, no. 3, pp. 331–368, 2022. doi: [10.1007/s41095-022-0271-y](https://doi.org/10.1007/s41095-022-0271-y).
- [18] G. Brauwers and F. Frasincar, “A general survey on attention mechanisms in deep learning,” *IEEE Trans. Knowl. Data Eng.*, vol. 35, no. 4, pp. 3279–3298, 2021. doi: [10.1109/TKDE.2021.3126456](https://doi.org/10.1109/TKDE.2021.3126456).
- [19] C. Zhang, J. Chen, J. Li, Y. Peng, and Z. Mao, “Large language models for human-robot interaction: A review,” *Biomimetic Intell. Robot.*, vol. 3, no. 4, pp. 100131, 2023. doi: [10.1016/j.birob.2023.100131](https://doi.org/10.1016/j.birob.2023.100131).
- [20] Y. Peng, H. Nabae, Y. Funabora, and K. Suzumori, “Controlling a peristaltic robot inspired by inchworms,” *Biomimetic Intell. Robot.*, vol. 4, no. 1, pp. 100146, 2024. doi: [10.1016/j.birob.2024.100146](https://doi.org/10.1016/j.birob.2024.100146).