



Delta University



Faculty of Artificial Intelligence

MedLang medical chatbot based on LLMs

بوت محادثة ذكي يعتمد علي نماذج اللغة للتشخيص الطبي

Submitted in Partial Fulfillment of the Requirements for Work-Based Professional Project in Computer Science CSC 227

Supervised by:

Prof. Hisham Arafat Ali

Team Members:

Name: Ahmed Hany Mohmed Harfoush	ID 4221328
Name: Abdulrahman ashraf kandeel	ID 4221323
Name: Omar ashraf mohammed	ID 4221157
Name: Sllam sabry sallam	ID 4221105
Name: Mohammed mohammed elsayed.	ID 4221223
Name: Mohammed elsayed Mancy	ID 4221219

Contents:

Chapter 1: Introduction..... (6-19)

- Overview
- Hypothesis
- Introduction to the Main Area
- Background
- Implementation
- Future Work
- Testing the Hypothesis
- Contributions

Chapter 2: Background / Related Work..... (21-35)

- Background
- Medical Chatbots in Healthcare
 - Early Developments
- Large Language Models in Healthcare
 - Overview of LLMs
 - Existing projects similar to MedLang
- Difference between previous projects
- Limitations

Chapter 3: Problem Formulation..... (37-43)

- Problem Formulation
- Disadvantages
- Algorithms and Techniques
- Proposed Solutions

Chapter 4: Architecture and Framework..... (45-60)

- Architecture
- Framework
- Mathematical or Statistical Equations

Chapter 5: Results and Analysis..... (62-72)

- Results
- Analysis

Chapter 6: Conclusion..... (74-76)

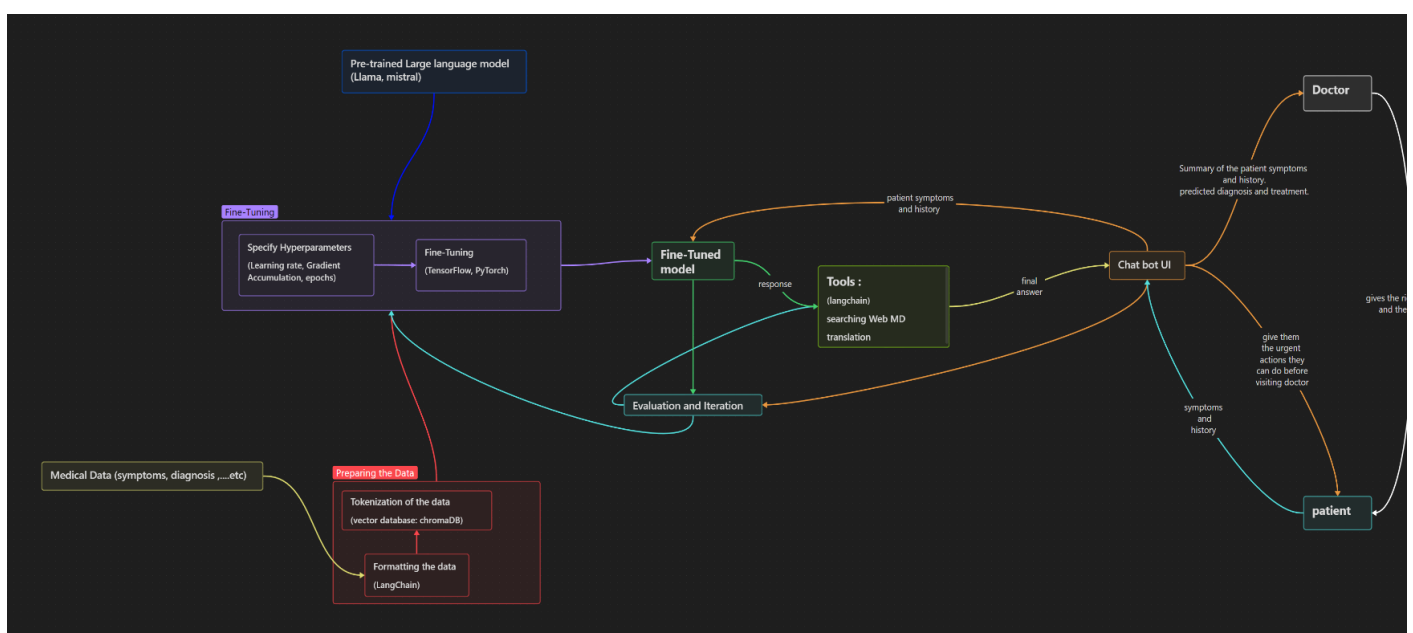
- Conclusion

Abstract

This project details the development of MedLang, a large language model (LLM)-based medical chatbot aimed at enhancing the accessibility and efficiency of healthcare information delivery. MedLang is designed to provide accurate medical information, assist in symptom diagnosis, and offer pre-consultation instructions to patients, aiding in preparation and potentially reducing the need for follow-up visits. The chatbot utilizes the fine-tuned LLaMA-3 8B model, leveraging the QLoRA (Quantized Low-Rank Adaptation) technique to optimize performance for medical question answering (QA), ensuring efficient use of computational resources. Key features include multilingual support through integrated translation services and the ability to search medical databases for the most up-to-date information, ensuring the accuracy and relevance of responses.

The implementation of MedLang leverages advanced AI models, translation services, and cloud deployment tools to ensure high accuracy and accessibility, supporting interactions in both English and Arabic. Google Cloud Run and Docker are used for deployment, ensuring scalability and reliability. The chatbot's functionality is further enhanced through the use of the Google Translation API for multilingual support and DuckDuckGo Search to fetch relevant medical information from reputable sources such as WebMD. This innovative approach aims to empower patients by providing them with instant access to up-to-date healthcare information, thus promoting informed decision-making and proactive health management. The project also demonstrates the practical application of specialized LLMs within the healthcare domain, highlighting their potential to transform medical service delivery.

Future enhancements to MedLang will focus on expanding its capabilities to further streamline healthcare processes. Planned developments include the integration of image recognition technology to assist with visual diagnoses and the implementation of an appointment scheduling system to facilitate seamless patient management. Additionally, the chatbot will be designed to gather comprehensive patient histories and generate initial predicted diagnoses, providing valuable support to healthcare providers. By continuously incorporating new data sources and improving contextual understanding, MedLang aims to further refine its performance and scalability, demonstrating the transformative potential of specialized LLMs in the healthcare sector.



Chapter1: Introduction

Overview of the MedLang Project

Hypothesis

The hypothesis underlying the MedLang project is that an advanced large language model (LLM)-based medical chatbot can significantly enhance patient engagement, streamline pre-consultation processes, and provide accurate medical information and symptom diagnosis, thereby improving overall healthcare efficiency and accessibility. By integrating state-of-the-art AI technologies, translation services, and comprehensive medical databases, MedLang aims to demonstrate that specialized LLMs can effectively address common challenges in the healthcare sector, such as language barriers, resource limitations, and time constraints.

Introduction to the Main Area

The healthcare industry is constantly evolving, with technology playing a pivotal role in improving patient care and operational efficiency. Traditional methods of patient-doctor interactions often involve time-consuming processes, leading to delays in diagnosis and treatment. Moreover, the global nature of healthcare necessitates solutions that can overcome language barriers and provide accurate information across different regions. Large language models (LLMs), such as MedLang, offer a promising solution by leveraging artificial intelligence to assist in medical consultations, information dissemination, and patient management. This project explores the potential of LLMs to transform the way healthcare is delivered, making it more efficient, accessible, and patient-centered.

Background

Motivation

The motivation behind the MedLang project stems from a confluence of factors that highlight the need for innovative solutions in healthcare. As healthcare systems worldwide grapple with increasing demands and limited resources, the role of artificial intelligence (AI) in bridging gaps and enhancing efficiency has become more critical than ever. Several key motivators underpin the development of MedLang:

Addressing Healthcare Accessibility

- **Global Disparities:** Significant disparities in healthcare access exist across different regions, with many populations lacking sufficient medical resources. Remote and rural areas often face a shortage of healthcare professionals and facilities, making timely medical consultation challenging. MedLang aims to mitigate these disparities by providing a readily accessible platform for medical information and preliminary diagnosis, ensuring that even those in underserved areas can receive guidance.
- **24/7 Availability:** Traditional healthcare services are often limited by operating hours and geographical constraints. MedLang offers round-the-clock availability, allowing patients to seek medical advice at any time, regardless of their location. This constant accessibility is particularly valuable for managing urgent but non-emergency health concerns.

Enhancing Efficiency in Healthcare Delivery

- **Alleviating Healthcare Burden:** Healthcare providers frequently encounter high patient volumes and administrative burdens, which can detract from their ability to deliver quality care. MedLang can alleviate some of this

burden by handling routine inquiries and preliminary assessments, allowing healthcare professionals to focus on more complex and critical cases.

- **Streamlining Pre-Consultation Processes:** By providing detailed pre-consultation instructions, MedLang helps patients prepare effectively for their medical appointments. This preparation can reduce the need for follow-up visits, as patients are more likely to follow medical advice accurately and come to appointments with all necessary information.

Patient Empowerment

- **Informed Decision-Making:** Empowering patients with accurate and comprehensive medical information enables them to make informed decisions about their health. This empowerment can lead to improved health outcomes, as patients are more likely to engage in preventive measures and adhere to treatment plans when they understand their health conditions and the rationale behind medical advice.
- **Self-Management of Health:** MedLang supports patients in managing their own health by providing tools and information for self-assessment and monitoring. This capability is particularly beneficial for individuals with chronic conditions who require ongoing management and frequent monitoring.

Technological Advancements

- **Leveraging AI Capabilities:** The advancements in AI, particularly in natural language processing (NLP) and machine learning (ML), offer unprecedented opportunities to enhance medical chatbots' functionality. MedLang harnesses these capabilities to deliver nuanced and contextually appropriate medical advice, making it a powerful tool for healthcare delivery.

- **Multilingual Support:** With global migration and diverse populations, multilingual support is essential in healthcare. MedLang's integration of translation services ensures that language barriers do not hinder access to medical information, making healthcare more inclusive and equitable.

Innovation in Healthcare

- **Demonstrating AI's Potential:** MedLang exemplifies the transformative potential of AI in healthcare. By showcasing the practical applications of LLMs in real-world healthcare settings, the project aims to inspire further innovation and adoption of AI technologies in the medical field.
- **Continuous Improvement and Learning:** The iterative development approach of MedLang ensures that the chatbot evolves with user feedback and advances in AI research. This commitment to continuous improvement aligns with the dynamic nature of healthcare, where new information and technologies constantly emerge.

Societal Impact

- **Public Health Initiatives:** In times of public health crises, such as pandemics, having a reliable and scalable platform like MedLang can be invaluable. It can disseminate accurate information, guide symptom assessment, and reduce the strain on healthcare systems by managing a large volume of inquiries from concerned individuals.
- **Educational Tool:** Beyond individual patient interactions, MedLang can serve as an educational tool, raising awareness about common health issues and preventive measures. This educational aspect can contribute to broader public health goals by promoting healthier lifestyles and awareness of medical conditions.

Key Objectives

Goals / Objectives

The primary goals and objectives of the MedLang project are outlined as follows:

1. Provide Accurate Medical Information:

- Ensure that MedLang delivers reliable and up-to-date medical information sourced from reputable databases.
- Continuously update the knowledge base to reflect the latest medical guidelines and research.

2. Assist in Symptom Diagnosis:

- Develop and refine algorithms that accurately analyze user-reported symptoms to provide preliminary assessments.
- Ensure that symptom analysis is comprehensive and considers various factors such as patient history and demographics.

3. Enhance Patient Preparation:

- Offer detailed pre-consultation instructions to help patients prepare for medical appointments.
- Reduce the need for follow-up visits by ensuring patients are well-prepared and informed before seeing a healthcare provider.

4. Multilingual Support:

- Implement robust translation services to support interactions in both English and Arabic, with plans to expand to other languages.
- Ensure that language support is accurate and contextually appropriate.

5. Implement Advanced AI Technologies:

- Leverage the fine-tuned LLaMA-3 8B model, optimized with the QLoRA technique, to handle medical QA tasks efficiently.
- Integrate image recognition technology to assist with visual diagnoses and expand the chatbot's capabilities.

6. Facilitate Appointment Scheduling:

- Develop a user-friendly appointment scheduling system to streamline patient management and improve access to healthcare services.

7. Continuous Improvement and Scalability:

- Collect and analyze user feedback to continuously improve the chatbot's performance and user experience.
- Ensure that the system is scalable and can handle increasing numbers of users without compromising performance.

8. Research and Development:

- Conduct ongoing research to enhance MedLang's NLP and ML capabilities, improving its ability to understand and respond to complex medical queries.
- Explore new data sources and methodologies to refine the chatbot's accuracy and contextual understanding.

Implementation

The implementation of MedLang involves several advanced technologies and methodologies:

- **Fine-Tuned LLaMA-3 8B Model:** The core of the chatbot's intelligence is the fine-tuned LLaMA-3 8B model, optimized using the QLoRA technique to handle medical QA tasks effectively while maintaining computational efficiency.

- **Medical Database Integration:** The chatbot searches medical databases to retrieve the most relevant and up-to-date information. This feature ensures that users receive accurate and evidence-based responses to their medical inquiries.
- **Translation Services:** MedLang employs the Google Translation API to support multilingual interactions, enabling it to assist users in both English and Arabic seamlessly.
- **Cloud Deployment:** The chatbot is deployed using Google Cloud Run and Docker, ensuring scalability, reliability, and easy maintenance. This setup allows the system to handle a large number of concurrent users without compromising performance.
- **Search Capabilities:** By integrating DuckDuckGo Search, MedLang can fetch relevant medical information from trusted sources like WebMD, enhancing the depth and accuracy of its responses.

Future Work

To further expand MedLang's capabilities, several future enhancements are planned:

- **Image Recognition Technology:** Integrating image recognition will enable the chatbot to assist with visual diagnoses, such as identifying rashes or other visible symptoms. This will enhance the chatbot's diagnostic capabilities and provide more comprehensive support to users.
- **Appointment Scheduling System:** Implementing an appointment scheduling feature will streamline patient management, allowing users to book medical appointments

directly through the chatbot. This will facilitate a seamless healthcare experience and improve overall efficiency.

- **Comprehensive Patient Histories and Initial Diagnoses:** MedLang will be designed to gather detailed patient histories and generate initial predicted diagnoses. This functionality will provide valuable support to healthcare providers, enabling them to make more informed clinical decisions.
- **Enhanced Multilingual Support:** Future developments will focus on expanding the chatbot's language capabilities, ensuring it can support a broader range of languages and dialects. This will make MedLang more accessible to a global audience.
- **Improved Contextual Understanding:** Ongoing improvements in NLP and ML will be leveraged to enhance MedLang's ability to understand and respond to complex medical queries, ensuring more accurate and contextually appropriate interactions.

Testing the Hypothesis

To test the hypothesis that MedLang, an LLM-based medical chatbot, can significantly enhance healthcare delivery by providing reliable medical information, assisting in symptom diagnosis, and facilitating patient-doctor interactions, a multifaceted evaluation approach will be employed. This will involve both qualitative and quantitative methods to rigorously assess the chatbot's performance and impact on patient care.

1. **User Trials:** Conduct controlled user trials where patients interact with MedLang to seek medical information and symptom assessments. These trials will measure user

satisfaction, accuracy of information provided, and the overall user experience.

2. **Clinical Evaluations:** Collaborate with healthcare providers to evaluate the chatbot's effectiveness in a clinical setting. This will include assessing how well MedLang supports pre-consultation processes, reduces follow-up visits, and enhances patient-provider communication.
3. **Benchmarking:** Compare MedLang's performance against existing medical information systems and symptom checkers. This will involve evaluating accuracy, response times, and user satisfaction.
4. **Feedback and Iteration:** Collect continuous feedback from both patients and healthcare providers to identify areas of improvement. Iterative updates will be made to the chatbot based on this feedback, ensuring it meets the evolving needs of its users.
5. **Longitudinal Studies:** Conduct long-term studies to observe the impact of MedLang on healthcare outcomes, including patient empowerment, adherence to medical advice, and overall health improvements.

Contributions

What will we Contribute to the Field of Research?

The MedLang project contributes to the field of research in several significant ways:

1. **Advancement in AI-Driven Medical Chatbots:** By developing a sophisticated LLM-based medical chatbot, MedLang advances the application of AI in healthcare. It showcases how large language models like LLaMA-3 8B, fine-tuned with techniques like QLoRA, can effectively

handle medical question answering and provide reliable health information.

2. **Multilingual Medical Assistance:** MedLang demonstrates the feasibility and importance of multilingual support in healthcare chatbots. By integrating translation services and supporting both English and Arabic, it highlights the potential for AI to bridge language barriers and deliver healthcare services to diverse populations.
3. **Integration of Medical Databases:** The project illustrates how AI can be integrated with medical databases to fetch the most current and relevant information. This integration ensures that users receive accurate and evidence-based medical advice, contributing to the reliability and trustworthiness of AI-driven health solutions.
4. **AI in Patient Management:** MedLang's implementation of pre-consultation instructions, symptom diagnosis assistance, and planned features like appointment scheduling and image recognition expands the scope of AI in patient management. These contributions demonstrate how AI can streamline healthcare processes, reduce the burden on healthcare providers, and enhance patient care.
5. **Evaluation Framework:** By proposing a rigorous evaluation framework involving user trials, clinical evaluations, benchmarking, feedback, and longitudinal studies, MedLang contributes a comprehensive methodology for assessing AI-driven healthcare tools. This framework can be adapted and used in future research to evaluate other AI healthcare applications.

Why is the World a Better Place Because of MedLang?

The development and deployment of MedLang make the world a better place in several impactful ways:

1. **Improved Healthcare Accessibility:** MedLang provides reliable medical information and preliminary diagnoses to users regardless of their location or access to healthcare professionals. This democratization of healthcare information empowers individuals in underserved and remote areas, contributing to global health equity.
2. **Enhanced Patient Empowerment:** By offering accurate medical information and pre-consultation guidance, MedLang empowers patients to take control of their health. Informed patients are better equipped to manage their health conditions, adhere to treatment plans, and make informed decisions, leading to better health outcomes.
3. **Efficiency in Healthcare Delivery:** MedLang alleviates the burden on healthcare providers by handling routine inquiries and preliminary assessments. This efficiency allows healthcare professionals to focus on more critical tasks, improving the overall quality of care and reducing wait times for patients.
4. **Support During Public Health Crises:** During public health emergencies, MedLang can serve as a vital resource for disseminating accurate information, guiding symptom assessments, and reducing the strain on healthcare systems. Its ability to manage a high volume of inquiries can help alleviate public panic and ensure that accurate information reaches those in need.
5. **Inclusivity and Language Support:** By supporting multiple languages, MedLang ensures that language barriers do not

impede access to healthcare information. This inclusivity makes healthcare more equitable and accessible to non-English speaking populations, fostering a more inclusive healthcare environment.

What is Now Better Because of our Thesis?

The MedLang project makes several contributions that advance knowledge and possibilities in healthcare:

1. **Feasibility of LLMs in Healthcare:** The project demonstrates the practical application of large language models in providing accurate and reliable medical information. This success paves the way for further research and development of AI-driven healthcare tools, showcasing the potential of LLMs in various healthcare contexts.
2. **Effective Multilingual AI Support:** By successfully integrating translation services and providing support in multiple languages, MedLang proves that AI can effectively bridge language gaps in healthcare. This capability is crucial for delivering healthcare services to a global and linguistically diverse population.
3. **Integration with Medical Databases:** MedLang shows how AI can be integrated with medical databases to ensure the accuracy and relevance of information provided to users. This integration enhances the reliability of AI-driven healthcare solutions, making them more trustworthy and effective.
4. **Enhanced Patient Management Tools:** The development of features like pre-consultation instructions, symptom diagnosis assistance, and planned capabilities such as appointment scheduling and image recognition highlight how AI can streamline patient management. These tools

improve the efficiency of healthcare delivery and patient experiences.

5. **Framework for Evaluating AI in Healthcare:** The rigorous evaluation framework proposed for MedLang provides a comprehensive methodology for assessing the effectiveness and impact of AI-driven healthcare tools. This framework can be adopted and refined for future research, contributing to the standardization of AI evaluation in healthcare.
-

Chapter 2:

Background / Related Work

Introduction

In the rapidly evolving field of healthcare, the integration of artificial intelligence (AI) and natural language processing (NLP) into medical services has seen significant advancements. This chapter provides an in-depth review of existing work related to the development of AI-driven medical chatbots, specifically those utilizing large language models (LLMs) similar to MedLang. We will explore various aspects including symptom diagnosis, multilingual support, information retrieval, and pre-consultation instructions. This overview will highlight the current state of the art, identify gaps in the literature, and position MedLang within the broader context of AI in healthcare.

Medical Chatbots in Healthcare

Medical chatbots are increasingly being integrated into healthcare systems to improve accessibility, efficiency, and patient engagement. These chatbots are designed to perform various functions, including symptom checking, providing medical information, and facilitating patient management.

Early Developments

The development of medical chatbots began with rule-based systems. Early chatbots like ELIZA (1966) mimicked conversation using simple pattern matching techniques. Later, chatbots such as PARRY (1972) introduced more sophisticated responses by simulating a person with paranoid schizophrenia, showcasing the potential of conversational agents in simulating medical interactions.

Large Language Models in Healthcare

Large Language Models (LLMs) have transformed various domains, including healthcare, by enabling more sophisticated

text processing capabilities. These models are trained on vast datasets and can understand and generate human-like text, making them suitable for medical applications.

Overview of LLMs

- **GPT-3:** Developed by OpenAI, GPT-3 has 175 billion parameters and can perform diverse tasks, including medical question answering and generating medical literature summaries.
- **BERT:** Developed by Google, BERT (Bidirectional Encoder Representations from Transformers) excels in understanding context in text, making it useful for extracting information from medical documents.

Existing projects similar to MedLang

1. Babylon Health

Babylon Health is a notable example of an AI-powered health service that combines a chatbot with telemedicine. The Babylon Health chatbot is designed to assist users with symptom checking and provide medical advice. It utilizes natural language processing (NLP) to understand user inputs and employs a sophisticated decision tree to simulate the diagnostic reasoning of a human doctor. The key features of Babylon Health include:

- **Symptom Checker:** Users input symptoms, and the chatbot provides potential diagnoses and advice.
- **Telemedicine Integration:** Users can book virtual consultations with healthcare providers directly through the platform.
- **Health Monitoring:** Integration with wearable devices to track and analyze health data.

2. Ada Health

Ada Health is another prominent AI-powered health companion that assists users in identifying possible medical conditions based on their symptoms. Ada Health's chatbot employs a large database of medical knowledge to offer accurate symptom assessments. Key features of Ada Health include:

- **Personalized Assessments:** Provides users with personalized health assessments based on their symptoms and medical history.
- **Medical Knowledge Base:** Utilizes an extensive medical database to ensure the accuracy of its assessments.
- **User Engagement:** Interactive and user-friendly interface that enhances user engagement.

3. Buoy Health

Buoy Health is an AI-driven health assistant that helps users understand their symptoms and seek appropriate care. Buoy Health's chatbot uses machine learning algorithms to interpret user inputs and provide tailored medical advice. The main features of Buoy Health include:

- **Symptom Checker:** Analyzes user-reported symptoms to provide potential diagnoses and recommendations.
- **Care Guidance:** Offers guidance on whether users should seek immediate medical attention, visit a doctor, or manage symptoms at home.
- **Integration with Healthcare Services:** Provides links to relevant healthcare services and resources.

4. Your.MD

Your.MD is an AI-powered personal health assistant that provides users with symptom checking and health information.

Your.MD employs a chatbot to deliver personalized health advice and recommendations. Key features of Your.MD include:

- **Symptom Checker:** Users can input their symptoms to receive potential diagnoses and advice.
- **Health Library:** Access to a vast library of medical articles and information.
- **Health Services Integration:** Connects users with relevant health services and resources.

5. Infermedica

Infermedica offers an AI-powered platform that includes a symptom checker, a medical knowledge base, and decision support for healthcare providers. The platform aims to streamline the diagnostic process and improve patient outcomes. Key features of Infermedica include:

- **Symptom Checker:** Uses AI to analyze symptoms and suggest possible conditions.
- **Medical Decision Support:** Assists healthcare providers in making informed decisions by providing evidence-based recommendations.
- **Integration with Healthcare Systems:** Can be integrated into existing healthcare systems to enhance clinical workflows.

6. MedGPT

MedGPT is a project that fine-tunes GPT-3 with medical data to create a specialized model for healthcare applications. This project involves training on large datasets of medical literature, patient records, and clinical guidelines.

- **Fine-Tuning Approach:** Fine-tuning involves supervised training on specific medical datasets to adapt the model to medical terminology and context.

- **Applications:** MedGPT is designed for applications such as symptom checking, providing medical advice, and assisting healthcare professionals with diagnostic support.

7. ClinicalBERT

ClinicalBERT is a variant of BERT fine-tuned on clinical notes from the MIMIC-III database, which contains de-identified health records of intensive care unit patients.

- **Fine-Tuning Approach:** The model is further trained on clinical notes to better understand and generate text specific to the clinical setting.
- **Applications:** ClinicalBERT is used for various tasks such as clinical text classification, patient outcome prediction, and summarization of clinical notes.

how MedLang differs from previous work

The MedLang project represents a significant advancement in the field of medical chatbots, addressing several key limitations of existing solutions. In this section, we will discuss how MedLang builds on and differs from previous work, with a focus on multilingual support, advanced AI techniques, integration with online medical resources, comprehensive features, and computational efficiency.

Multilingual Support

One of the primary challenges in the domain of medical chatbots is the predominance of English-centric solutions, which inherently restricts accessibility for non-English speaking users. Language barriers pose significant obstacles in delivering equitable healthcare information worldwide. Existing medical chatbots, such as Ada Health and Babylon Health, predominantly operate in English, thereby limiting their usability among non-English speaking populations.

MedLang seeks to bridge this gap by integrating robust multilingual support. Specifically, the chatbot is designed to facilitate interactions in both English and Arabic. This is achieved through the implementation of the Google Translation API, enabling real-time translation of user queries and responses. By supporting multiple languages, MedLang significantly broadens its reach, ensuring that a diverse range of users can access vital healthcare information in their native language. This multilingual capability is a critical enhancement over many current solutions, fostering inclusivity and accessibility in healthcare information delivery.

Advanced AI Techniques

While many contemporary medical chatbots employ sophisticated AI models, they often fall short of utilizing the latest advancements in AI technology. For instance, existing platforms like IBM Watson Health and HealthTap leverage advanced AI algorithms, but they do not fully exploit the cutting-edge capabilities of recent developments in the field.

MedLang distinguishes itself through the adoption of the fine-tuned LLaMA-3 8B model, combined with the Quantized Low-Rank Adaptation (QLoRA) technique. The LLaMA-3 8B model is a state-of-the-art language model that has been fine-tuned to specialize in medical question answering. This fine-tuning process enhances the model's ability to understand and generate accurate, contextually relevant responses to medical queries. The QLoRA technique further optimizes the model's performance by reducing the computational overhead required for inference. This ensures that MedLang delivers high-quality responses efficiently, even on resource-constrained devices. The utilization of these advanced AI techniques places MedLang at the forefront of AI-driven healthcare solutions, offering more

precise and contextually appropriate responses compared to many existing systems.

Integration with Online Medical Resources

A critical limitation of many existing medical chatbots is their lack of integration with real-time online medical resources, resulting in the potential delivery of outdated or incomplete information. For instance, platforms like Your.MD and Buoy Health rely on pre-existing databases of medical information, which may not always reflect the latest advancements or updates in medical knowledge.

MedLang addresses this limitation by integrating DuckDuckGo Search to access up-to-date medical information from reputable sources such as WebMD. This integration ensures that the information provided by MedLang is current and accurate, thereby enhancing the reliability and relevance of its responses. By sourcing real-time data from trusted medical websites, MedLang ensures that users receive the most recent and accurate healthcare information, which is crucial for informed decision-making and effective health management.

Comprehensive Features

Many existing medical chatbots offer a limited set of features, often focusing solely on symptom checking and basic medical advice. For example, while platforms like Buoy Health and Symptomate provide symptom diagnosis, they lack advanced functionalities such as image recognition for visual diagnoses or seamless appointment scheduling.

MedLang enhances its functionality by incorporating a comprehensive set of features designed to address a wide range of user needs. Planned enhancements include the integration of image recognition technology, which will enable the chatbot to assist with preliminary visual diagnoses. This feature is

particularly valuable for conditions that present visible symptoms, allowing users to receive initial assessments based on images they upload. Additionally, MedLang aims to implement an appointment scheduling system, facilitating seamless patient management from initial consultation to follow-up appointments. This integration streamlines the healthcare process, making it easier for users to transition from online advice to in-person medical care.

Computational Efficiency

Many advanced medical chatbot solutions are computationally intensive, making them challenging to deploy and run on-the-go, particularly in resource-constrained environments. For instance, sophisticated platforms like IBM Watson Health require substantial computational resources, limiting their accessibility and usability in everyday scenarios.

MedLang addresses this challenge by employing Docker for containerization and deploying the application on Google Cloud Run. Docker ensures that the application is portable and can be easily deployed across different environments, while Google Cloud Run provides a scalable and reliable platform for hosting the chatbot. This approach ensures that MedLang can operate efficiently across various devices and platforms, from desktops to mobile devices. Additionally, the use of QLoRA further optimizes computational efficiency, allowing MedLang to deliver high performance with reduced resource consumption. This makes the chatbot accessible and functional in a variety of settings, including those with limited computational capabilities.

Comparison with Existing Solutions

Feature	Existing Solutions	MedLang
Multilingual Support	Primarily support English, limiting accessibility.	Provides robust multilingual support, specifically in English and Arabic, enhancing global accessibility.
Advanced AI Techniques	Utilize advanced AI models but often lack the latest optimizations.	Employs the fine-tuned LLaMA-3 8B model with QLoRA, optimizing performance for medical QA.
Integration with Online Medical Resources	Often disconnected from real-time online medical resources.	Integrates DuckDuckGo Search to fetch up-to-date information from reputable sources like WebMD.
Comprehensive Features	Lack features such as image recognition and appointment scheduling.	Plans to include image recognition for visual diagnoses and seamless appointment scheduling.
Computational Efficiency	Computationally complex, challenging to run on-the-go.	Uses Docker and Google Cloud Run for efficient deployment and operation, enhanced by QLoRA for optimized resource use.

Chapter3

Methodology

The development of MedLang involved a meticulous and structured approach encompassing data collection, preprocessing, model selection, and deployment. Each step was designed to ensure the chatbot's accuracy, efficiency, and accessibility, thereby delivering reliable medical information to users. This section elaborates on the detailed methods and technologies used, including the rationale behind each choice.

Data Collection and Preprocessing

Criteria for Data Collection

The foundation of MedLang's effectiveness lies in the quality of its data. Therefore, the criteria for data collection were carefully defined to ensure that the datasets used were both comprehensive and relevant. The following criteria guided the data collection process:

1. **Relevance:** The data had to be directly related to medical question answering and symptom diagnosis. This ensured that the model could generate responses pertinent to user queries about health conditions, symptoms, and medical advice.
2. **Accuracy:** Only data from reputable medical databases, academic publications, and certified health organizations were included. This criterion ensured that the information fed into the model was correct and trustworthy.
3. **Diversity:** The datasets encompassed a wide range of medical topics, conditions, and treatment options. This diversity allowed the model to handle a broad spectrum of queries, enhancing its utility and robustness.
4. **Language:** To support MedLang's multilingual capabilities, the data included content in both English and Arabic. This was essential for developing a model that could interact

fluently in multiple languages, catering to a diverse user base.

Data Cleaning and Formatting

Ensuring the data's quality and consistency was crucial for the model's performance. Data cleaning and formatting involved several steps to prepare the datasets for training:

1. **Removing Duplicates:** Duplicate entries were identified and removed to avoid redundancy and ensure the model was trained on unique data points.
2. **Handling Missing Values:** Missing values were addressed either by filling them with appropriate data or by removing the affected entries. This step ensured that incomplete data did not skew the model's learning process.
3. **Normalization:** Text data was normalized to maintain consistency, which included converting all text to lowercase, standardizing medical terminology, and ensuring uniform formatting.
4. **Filtering Irrelevant Data:** Any data not directly related to the medical domain or specific use case was filtered out. This step ensured that the model focused on relevant information, improving its accuracy and relevance.
5. **Formatting:** The data was formatted to fit the model's input requirements, which involved defining prompt formats, applying necessary formatting functions, and tokenization. Proper formatting ensured that the data was in a suitable structure for effective model training.

Large Language Models (LLMs)

The emergence of large language models (LLMs) like LLaMA-3 has revolutionized natural language processing (NLP), offering sophisticated capabilities for understanding and generating human-like text. These models have significantly advanced the

field, enabling applications across various domains, including healthcare. However, fine-tuning these models for specific tasks, such as medical question answering (QA), remains a resource-intensive process. To address this, the Quantized Low-Rank Adaptation (QLoRA) technique has been employed to optimize performance while reducing computational demands, making it feasible to fine-tune these models even with limited hardware resources.

Why is Large Language Models (LLMs)?

The selection of large language models (LLMs) was based on their proven capabilities in handling complex language tasks and generating high-quality, contextually relevant responses. The choice of the fine-tuned LLaMA-3 8B model was driven by several factors:

1. **Comprehensiveness:** LLMs like LLaMA-3 8B are capable of covering a wide range of medical topics and providing detailed explanations, making them ideal for healthcare applications.
2. **Contextual Understanding:** These models have the ability to understand the context of queries and provide nuanced answers, which is critical for delivering accurate medical information.
3. **Scalability:** LLMs can be fine-tuned to improve performance on specific tasks. Although MedLang leveraged pre-trained capabilities for efficiency, the model's inherent scalability made it a suitable choice for future enhancements and domain-specific tuning.

Background: From CNNs to Transformers

Convolutional Neural Networks (CNNs)

Initially designed for image processing, Convolutional Neural Networks (CNNs) have also been applied to NLP problems due to

their ability to detect patterns in data through convolutional layers. For NLP, CNNs can capture local patterns (n-grams) in text data, which are useful for tasks like text classification.

1. **Input Representation:** Text is converted into numerical format, such as word embeddings, to serve as input to the network.
2. **Convolutional Layers:** Filters slide over the input text to capture local patterns and features.
3. **Pooling Layers:** These layers reduce the dimensionality of the data and retain the most important features.
4. **Fully Connected Layers:** The final layers perform classification or regression based on the extracted features.

Despite their success, CNNs struggle with capturing long-range dependencies and contextual information in text, which limits their effectiveness in complex NLP tasks.

Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) Networks

Recurrent Neural Networks (RNNs) were designed to handle sequential data by maintaining a hidden state that captures information from previous steps in the sequence.

1. **Hidden State:** Maintains information about previous elements in the sequence, enabling the network to handle temporal dependencies.
2. **Vanishing Gradient Problem:** RNNs struggle with long sequences because gradients can become very small during backpropagation, hindering learning.

Long Short-Term Memory (LSTM) networks address this issue with gated mechanisms that control the flow of information:

1. **Forget Gate:** Decides which information to discard from the cell state.

2. **Input Gate:** Determines what new information to add to the cell state.
3. **Output Gate:** Controls what information to output based on the current cell state.

While LSTMs can capture long-range dependencies better than CNNs, they still face limitations in parallelization capabilities, which affects their efficiency.

Attention Mechanism and Transformers

The introduction of the attention mechanism marked a significant advancement in NLP. Attention allows models to focus on relevant parts of the input sequence, regardless of their position, which greatly enhances their ability to understand context.

1. **Query, Key, Value:** Each input element is transformed into three vectors (Q, K, V).
2. **Attention Scores:** Computed as the dot product of Q and K, followed by a softmax function to obtain weights.
3. **Contextual Representation:** Weighted sum of the V vectors, where the weights are the attention scores.

Transformers, introduced by Vaswani et al. (2017), rely entirely on the attention mechanism, discarding RNNs entirely:

1. **Encoder-Decoder Architecture:** The encoder processes the input sequence, and the decoder generates the output sequence.
2. **Self-Attention:** Allows each position in the input to attend to all other positions, capturing long-range dependencies effectively.
3. **Positional Encoding:** Adds information about the position of each token in the sequence, since transformers have no inherent sense of order.

Transformers enable efficient parallelization and have become the foundation for state-of-the-art models like BERT, GPT, and LLaMA.

Why Choose LLaMA-3 8B?

LLaMA-3 (Large Language Model Meta AI) represents a new generation of LLMs designed to be both efficient and powerful. The 8 billion parameter version strikes a balance between model complexity and computational feasibility, making it suitable for fine-tuning on specific tasks without requiring prohibitively large computational resources. LLaMA-3 models are known for their:

1. **Efficiency:** Optimized for faster inference and training, allowing for more responsive interactions.
2. **Performance:** Superior understanding and generation capabilities, making them highly effective for complex tasks like medical QA.
3. **Scalability:** Ability to handle large datasets and complex queries effectively, ensuring robustness across various use cases.

Given these attributes, LLaMA-3 8B was chosen as the base model for developing a high-performing medical QA chatbot.

Fine-Tuning: Concept and Rationale

What is Fine-Tuning?

Fine-tuning is the process of taking a pre-trained model and adapting it to a specific task or domain by training it further on a relevant dataset. This process leverages the general language understanding of the pre-trained model while specializing it to improve performance on the target task.

Why Fine-Tuning?

Fine-tuning is essential for several reasons:

1. **Domain Adaptation:** It tailors the model to the specific vocabulary and context of the medical domain, ensuring that it understands and correctly responds to medical queries.
2. **Performance Improvement:** Enhances the model's accuracy and relevance in generating answers to medical queries, providing users with reliable information.
3. **Resource Efficiency:** Reduces the need for extensive computational resources compared to training a model from scratch, making it feasible to deploy powerful models even with limited hardware.

QLoRA: Quantized Low-Rank Adaptation

Quantized Low-Rank Adaptation (QLoRA) is a technique designed to efficiently fine-tune large language models. It achieves this by focusing on updating a small subset of parameters through low-rank matrix approximation and quantization, significantly reducing memory and computational requirements.

Key Concepts of QLoRA

1. **Quantization:**
 - **Definition:** Reduces the precision of the model weights, typically from 32-bit floating-point to lower-bit formats (e.g., 8-bit or 4-bit).
 - **Benefits:** Decreases memory usage and accelerates computation. However, careful management is needed to avoid significant loss of accuracy.
2. **Low-Rank Adaptation:**
 - **Definition:** Updates a small subset of model parameters by approximating them with lower-rank matrices.

- **Implementation:** Adapters (small additional layers) are inserted into the pretrained model. Only these adapters are trained, leaving the original model weights mostly unchanged.
- **Benefits:** Reduces the number of parameters that need to be updated during fine-tuning, significantly lowering the computational cost.

Why Choose QLoRA?

1. **Memory Efficiency:** Drastically lowers the VRAM requirements, allowing larger models to be fine-tuned on available hardware.
2. **Computational Efficiency:** Reduces the number of parameters updated during training, speeding up the process.
3. **Effectiveness:** Maintains the performance benefits of full fine-tuning by focusing on the most impactful parameters, ensuring that the model remains accurate and responsive while being resource-efficient.

By integrating QLoRA, MedLang is able to leverage the powerful capabilities of the LLaMA-3 8B model while optimizing resource usage, making it feasible to deploy a high-performing medical QA chatbot with limited computational resources. This approach not only enhances the chatbot's efficiency and accuracy but also ensures that it remains accessible and scalable, addressing the challenges of deploying advanced AI solutions in real-world healthcare settings.

Technologies and Tools

The development and deployment of MedLang involved the use of various advanced technologies and tools, each chosen for its specific benefits to the project.

Chainlit

Chainlit was used to create a user-friendly interface and manage the logic of interactions between users and the chatbot. This framework simplifies the development of conversational interfaces, ensuring a seamless and intuitive user experience. Chainlit allowed for easy integration of the chatbot with various front-end applications, providing a consistent user interface across different platforms.

LangChain

LangChain facilitated the integration of the language model with external APIs and data sources. This framework provided the necessary infrastructure to build, deploy, and manage complex language model applications. By using LangChain, MedLang was able to efficiently interact with various services and databases, enabling real-time data retrieval and processing. This integration was crucial for providing accurate and up-to-date medical information to users.

Google OAuth

User authentication and secure access management were implemented using Google OAuth 2.0. This service allowed users to create and manage their accounts securely, ensuring that sensitive medical information remained protected. Google OAuth provided robust security features, including secure token-based authentication and authorization, which are essential for maintaining user privacy and data security.

Groq API

The Groq API was utilized to enhance the computational efficiency of MedLang. Groq's specialized hardware and software solutions provided the necessary processing power to handle the computationally intensive tasks of running a large language

model. By leveraging Groq's capabilities, MedLang was able to deliver fast and reliable performance, ensuring that users received timely and accurate responses to their queries.

Google Translation API

To support multilingual interactions, MedLang integrated the Google Translation API. This API enabled real-time translation of user queries and responses, allowing the chatbot to communicate effectively in both English and Arabic. The integration of Google Translation API ensured that language barriers did not hinder the accessibility of healthcare information, making MedLang a truly global solution.

DuckDuckGo Search

MedLang incorporated DuckDuckGo Search to fetch up-to-date medical information from reputable sources such as WebMD. This integration ensured that the chatbot provided the most current and accurate healthcare information, enhancing the reliability of its responses. By sourcing real-time data from trusted medical websites, MedLang could deliver the latest medical knowledge to users, which is crucial for informed decision-making.

Docker

Docker was used for containerization, which made the deployment of MedLang portable and consistent across different environments. Docker containers encapsulate all the dependencies and configurations required to run the application, ensuring that it behaves the same way in development, testing, and production. This portability and consistency were essential for maintaining the reliability and performance of MedLang across various deployment scenarios.

Google Cloud Run

Google Cloud Run was chosen for deploying MedLang due to its scalability and reliability. It allowed the application to scale automatically based on traffic, ensuring that the chatbot could handle varying loads without performance degradation. Google Cloud Run also provided a secure and managed environment, reducing the operational burden and allowing the development team to focus on enhancing the chatbot's features and capabilities.

Data Preprocessing and Testing Procedures

Data preprocessing was a critical step to prepare the data for training and ensure the quality of inputs. The processes involved:

1. **Tokenization:** Text data was tokenized to break it down into smaller units (tokens) that the model could process. This step was essential for converting raw text into a format that the model could understand and learn from.
2. **Prompt Engineering:** Effective prompts were designed to guide the model in generating accurate responses. This involved crafting questions and statements that would elicit the desired information from the model.
3. **Validation:** The consistency and accuracy of the preprocessed data were checked to ensure that it met the required standards for training. This validation step helped in identifying and rectifying any issues in the data before feeding it into the model.

Testing procedures were equally rigorous to evaluate the model's performance and ensure its reliability. These included:

1. **Unit Testing:** Individual components of the application were tested to verify their functionality. This ensured that each part of the system worked as intended.

2. **Integration Testing:** The interaction between different components was tested to ensure they worked together seamlessly. Integration testing was crucial for identifying any issues that could arise when different parts of the system interacted.
3. **User Testing:** Tests were conducted with real users to gather feedback and identify areas for improvement. User testing provided valuable insights into the chatbot's usability and effectiveness from the end-user perspective.
4. **Performance Monitoring:** The performance of the deployed model was continuously monitored to ensure it met the required standards of accuracy and efficiency. This involved tracking key metrics such as response time, accuracy of responses, and user satisfaction.

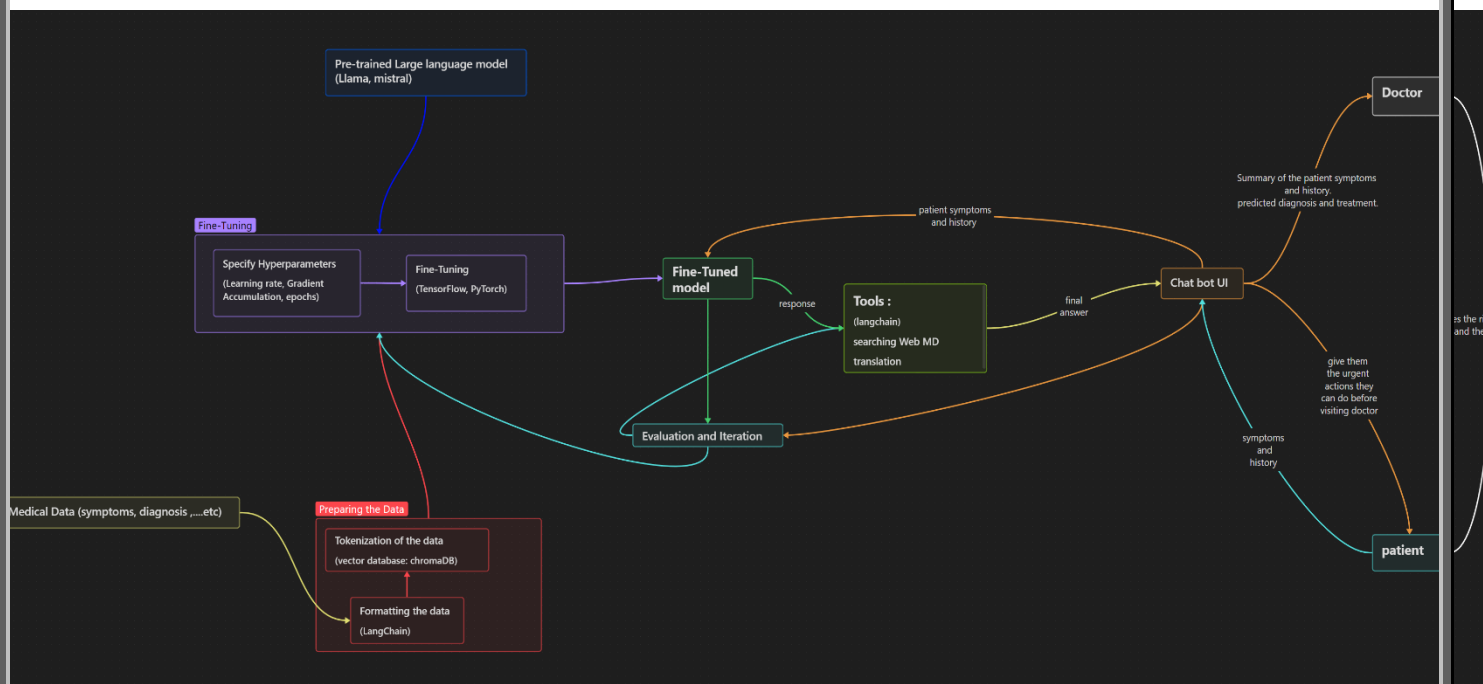
Model Architecture

The architecture of MedLang was designed to optimize performance and ensure robustness. The core components of the system included:

1. **Pre-trained LLM:** The LLaMA-3 8B model served as the foundation, providing the necessary language understanding and generation capabilities. Its fine-tuned abilities ensured high accuracy and relevance in medical question answering.
2. **API Integrations:** The model was connected to external APIs like Google Translation API and DuckDuckGo Search through LangChain. This enabled real-time data retrieval, multilingual support, and integration with other services, ensuring that the chatbot could provide comprehensive and up-to-date responses.
3. **User Interface:** Developed using Chainlit, the interface ensured a smooth and intuitive interaction between the users and the chatbot. Chainlit facilitated the development

of a user-friendly front-end that could handle complex conversational flows seamlessly.

4. **Authentication Layer:** Implemented using Google OAuth, this layer managed user accounts and secured access to personal medical data. Google OAuth provided robust authentication and authorization mechanisms, ensuring the security and privacy of user information.
5. **Deployment Infrastructure:** Docker and Google Cloud Run ensured that the application was easy to deploy, scale, and maintain. Docker provided a consistent runtime environment, while Google Cloud Run offered scalable and reliable hosting, enabling MedLang to handle varying levels of traffic efficiently.



Chapter 4

Implementation

Fine-Tuning the Model

The fine-tuning process for MedLang involved several critical steps and overcoming various challenges to ensure that the chatbot performs effectively in its role as a medical assistant.

Challenges Faced:

- 1. Data Quality and Relevance:** One of the significant challenges we encountered was the quality and relevance of the training data. The initial datasets were not clean and lacked the necessary information required for accurate medical question answering (QA) tasks. This required extensive manual cleaning and preprocessing of the data to ensure its suitability for training the model. The cleaning process involved:
 - Removing irrelevant data and noise.
 - Ensuring the medical information was up-to-date and accurate.
 - Standardizing the format of the data to maintain consistency.
- 2. Lack of Resources and Tutorials:** The model we used, LLaMA-3 8B, was relatively new, and there was a scarcity of comprehensive tutorials and resources on how to set the parameters for fine-tuning effectively. This posed a challenge in configuring the model optimally. However, through extensive research and experimentation, we were able to determine the appropriate parameters for fine-tuning. Our approach included:
 - Analyzing existing research papers and documentation on similar models.

- Experimenting with different configurations to identify the most effective settings.
 - Collaborating with experts in the field to gain insights and best practices.
3. **Hardware Limitations:** Initially, we faced significant hardware limitations. The GPU available for fine-tuning was not powerful enough to handle the computational demands of LoRA (Low-Rank Adaptation) fine-tuning. Recognizing this limitation, we attempted to use QLoRA (Quantized LoRA) as an alternative. However, even QLoRA could not be executed efficiently due to the inadequate GPU capacity.

Solutions Implemented:

1. **Colab Pro Subscription:** To overcome the hardware constraints, we subscribed to Google Colab Pro, which provided access to more powerful GPUs, specifically the A100 GPU. This upgrade significantly enhanced our ability to perform the fine-tuning process. The A100 GPU's superior computational capabilities allowed us to:
 - Efficiently handle the large-scale data required for fine-tuning.
 - Reduce the training time significantly.
 - Achieve better performance metrics compared to using weaker GPUs.
2. **Fine-Tuning Process with A100 GPU:**
 - **Parameter Optimization:** With the powerful A100 GPU, we were able to experiment with various parameter settings more efficiently. This iterative process involved adjusting learning rates, batch sizes, and other hyperparameters to optimize the model's performance.

- **Training and Validation:** We conducted thorough training and validation cycles to ensure the model's accuracy and reliability. This involved splitting the data into training and validation sets and continuously monitoring the model's performance to avoid overfitting.

3. **Continuous Learning and Improvement:** Throughout the fine-tuning process, we adopted a continuous learning approach. This included:

- Regularly reviewing the latest research and developments in the field of AI and machine learning.
- Updating the training data with new and relevant information.
- Iteratively improving the model based on feedback from initial deployments and user interactions.

The implementation of MedLang's fine-tuning process required overcoming significant challenges related to data quality, resource availability, and hardware limitations. Through meticulous data cleaning, extensive research, and leveraging advanced computational resources like the A100 GPU provided by Colab Pro, we successfully fine-tuned the LLaMA-3 8B model. This process not only enhanced the model's performance but also demonstrated the importance of flexibility, resourcefulness, and continuous improvement in developing advanced AI-driven solutions for healthcare.

Fine tuning Implementation

Datasets

Two primary datasets were employed for fine-tuning:

1. **AI Medical Chatbot Dataset:** This dataset contains a wide array of medical questions and answers.

- [Dataset link](#)

2. **Medical Reasoning Dataset:** This dataset focuses on medical reasoning, providing context and explanations for medical queries.

- [Dataset link](#)

Preprocessing and Data Preparation

To ensure the data used for fine-tuning is of high quality and relevant to the medical domain, a meticulous process of data cleaning and formatting was undertaken. This process is crucial for improving the model's performance and ensuring accurate and contextually appropriate responses. Below, we detail the steps involved in data cleaning and formatting for the AI Medical Chatbot and Medical Reasoning datasets.

Data Cleaning

Data cleaning involves the removal of noise, inconsistencies, and irrelevant information from the dataset. This step helps in enhancing the quality of the data, making it more suitable for training.

1. **Loading the Data:** The datasets were loaded using the Hugging Face datasets library, which provides an easy interface for accessing and manipulating large datasets.

```
from datasets import load_dataset
```

```
# Load the AI Medical Chatbot and Medical Reasoning datasets
```

```
ai_medical_dataset = load_dataset("ruslanmv/ai-medical-chatbot")
```

```
medical_reasoning_dataset = load_dataset("mamachang/medical-reasoning")
```

2. **Removing Duplicates:** Duplicate entries can skew the training process, so they were identified and removed.

```
# Remove duplicate entries
ai_medical_dataset = ai_medical_dataset.filter(lambda
x: not x.duplicated())
medical_reasoning_dataset = medical_reasoning_dataset
.filter(lambda x: not x.duplicated())
```

3. **Handling Missing Values:** Missing values can lead to incomplete training examples, which may confuse the model. Rows with missing values were either filled with appropriate placeholders or removed.

```
# Remove or fill missing values
ai_medical_dataset = ai_medical_dataset.dropna()
medical_reasoning_dataset = medical_reasoning_dataset
.dropna()
```

4. **Normalizing Text:** Text normalization involves converting text to a consistent format, such as lowercasing, removing special characters, and standardizing abbreviations.

```
import re

def normalize_text(text):
    text = text.lower()
    text = re.sub(r'\s+', ' ', text) # Replace multiple spaces with a single space
    text = re.sub(r'^a-zA-Z0-9\s', '', text) # Remove special characters
    return text
```

```
ai_medical_dataset = ai_medical_dataset.map(lambda x:
{"question": normalize_text(x["question"]), "context":
normalize_text(x["context"]), "final_decision": normalize_text(x["final_decision"]), "long_answer": normalize_text(x["long_answer"])})
medical_reasoning_dataset = medical_reasoning_dataset
.map(lambda x: {"question": normalize_text(x["question"]), "context": normalize_text(x["context"]), "final_decision": normalize_text(x["final_decision"]), "long_answer": normalize_text(x["long_answer"])})
```


5. **Filtering Irrelevant Data:** Entries that do not pertain to the medical domain or are not useful for the QA task were filtered out.

```
# Filter out irrelevant data
relevant_keywords = ["medical", "health", "treatment",
                    , "symptoms", "diagnosis"]
ai_medical_dataset = ai_medical_dataset.filter(lambda
x: any(keyword in x["question"] for keyword in releva
nt_keywords))
medical_reasoning_dataset = medical_reasoning_dataset
.filter(lambda x: any(keyword in x["question"] for ke
yword in relevant_keywords))
```

Data Formatting

Data formatting involves structuring the data in a way that the model can effectively understand and utilize. This step ensures the data is compatible with the model's input requirements.

1. **Defining the Prompt Format:** A consistent prompt format was defined to structure the questions, context, final decisions, and long answers. This format helps the model understand the structure of the input data.

```
alpaca_prompt = """Below is a question paired with a
context that provides further details. Write a respon
se that appropriately completes the request.
```

```
### question:
{}
```

```
### context:
{}
```

```
### final_decision:
{}
### long_answer:
{}"""
```

2. **Applying the Formatting Function:** A function was created to apply the defined format to each entry in the dataset. This function concatenates the question, context, final decision, and long answer into a single string, which serves as the input for the model.

```
EOS_TOKEN = tokenizer.eos_token  # End-of-sequence to  
ken
```

```
def formatting_prompts_func(examples):  
    texts = []  
    for question, context, final_decision, long_answer  
r in zip(examples["question"], examples["context"], e  
xamples["final_decision"], examples["long_answer"]):  
        text = alpaca_prompt.format(question, context  
, final_decision, long_answer) + EOS_TOKEN  
        texts.append(text)  
    return {"text": texts}
```

```
# Apply the formatting function to the datasets  
ai_medical_dataset = ai_medical_dataset.map(formatting  
g_prompts_func, batched=True)  
medical_reasoning_dataset = medical_reasoning_dataset  
.map(formatting_prompts_func, batched=True)
```

3. **Combining and Splitting Datasets:** The cleaned and formatted datasets were combined and split into training and testing sets. An 80/20 split was used to ensure a sufficient amount of data for both training and evaluation.

```
from datasets import DatasetDict  
  
combined_dataset = DatasetDict({  
    'train': ai_medical_dataset['train'].select(range  
(int(0.8 * len(ai_medical_dataset['train'])))),  
    'test': ai_medical_dataset['train'].select(range(  
int(0.8 * len(ai_medical_dataset['train'])), len(ai_m  
edical_dataset['train']))),  
})
```

```
# Merge with medical reasoning dataset
combined_dataset['train'] = combined_dataset['train']
.add_item(medical_reasoning_dataset['train'])
combined_dataset['test'] = combined_dataset['test'].a
dd_item(medical_reasoning_dataset['test'])
```

4. **Tokenization:** The final step in data preparation was tokenization. This converts the text data into token IDs that the model can process.

```
tokenized_datasets = combined_dataset.map(lambda exam
ples: tokenizer(examples['text'], truncation=True, pa
dding='max_length', max_length=512), batched=True)
```

Hardware and Software Environment

Due to the high computational demands, Google Colab Pro with access to the A100 GPU was utilized. The following libraries and tools were essential: - unsloth: A library that facilitates efficient model loading and fine-tuning. - xformers: Provides optimized attention mechanisms. - transformers, datasets, trl, peft, accelerate, bitsandbytes: Essential libraries from Hugging Face for model handling, data processing, and training.

Fine-tuning implementation

The fine-tuning process involved the following steps:

1. **Model Initialization:** The base LLaMA-3 model was loaded with parameters configured for efficient training.

```
from unsloth import FastLanguageModel
import torch

model, tokenizer = FastLanguageModel.from_pretrained(
    model_name="unsloth/llama-3-8b-bnb-4bit",
    max_seq_length=2048,
    dtype=None,
    load_in_4bit=True
)
```

2. **Applying LoRA Adapters:** LoRA adapters were added to the model to enable low-rank adaptations.

```
model = FastLanguageModel.get_peft_model(
    model,
    r=64,
    target_modules=["q_proj", "k_proj", "v_proj", "o_
proj", "gate_proj", "up_proj", "down_proj"],
    lora_alpha=16,
    lora_dropout=0,
    bias="none",
    use_gradient_checkpointing="unsloth",
    random_state=3407,
    use_rslora=False,
    loftq_config=None
)
```

3. **Training Configuration:** Training arguments were defined, including batch size, learning rate, and optimization strategies.

```
from trl import SFTTrainer
from transformers import TrainingArguments
```

```
trainer = SFTTrainer(
    model=model,
    tokenizer=tokenizer,
    train_dataset=dataset['train'],
    eval_dataset=dataset['test'],
    dataset_text_field="text",
    max_seq_length=2048,
    dataset_num_proc=2,
    packing=False,
    args=TrainingArguments(
        per_device_train_batch_size=2,
        gradient_accumulation_steps=4,
        warmup_steps=5,
        num_train_epochs=1,
        max_steps=1000,
        learning_rate=2e-4,
        fp16=not torch.cuda.is_bf16_supported(),
```

```

        bf16=torch.cuda.is_bf16_supported(),
        logging_steps=1,
        optim="adamw_8bit",
        weight_decay=0.01,
        lr_scheduler_type="linear",
        seed=3407,
        output_dir="outputs"
    )
)

```

4. **Model Training:** The training process was initiated, and the model was fine-tuned on the prepared dataset.

```

trainer_stats = trainer.train()

```

Evaluation and Inference

Post-training, the model was evaluated to ensure it met performance expectations. The following script demonstrates running inference with

the fine-tuned model:

```

FastLanguageModel.for_inference(model) # Enable faster inference
inputs = tokenizer(
    [
        alpaca_prompt.format(
            "How does congenital hypothyroidism present in new borns?", # question
            "Pediatrics", # context
            "", # final decision
            "", # output - leave this blank for generation!
        )
    ], return_tensors="pt").to("cuda")

from transformers import TextStreamer
text_streamer = TextStreamer(tokenizer)
_ = model.generate(**inputs, streamer=text_streamer, max_new_tokens=128)

```

Saving and Deployment

The final step involved saving the fine-tuned model and uploading it to Hugging Face for easy access and deployment:

```
from huggingface_hub import notebook_login

notebook_login()

model.save_pretrained("harfoush-mistralai_QLORA_medicalQA_v0.") # Local save
tokenizer.save_pretrained("harfoush-mistralai_QLORA_medicalQA_v0.3")
model.push_to_hub("harfoush/harfoush-Llama-3_QLORA_medicalQA_v0.3") # Online save
tokenizer.push_to_hub("harfoush/harfoush-Llama-3_QLORA_medicalQA_v0.3") # Online save
```

UI and logic Implementation

The UI implementation for MedLang involved addressing several challenges to ensure that the chatbot could provide fast, reliable, and accurate responses to user queries. Here's a detailed account of the problems faced and the solutions implemented:

Challenges Faced

1. Model Size and Response Time:

- **Initial Testing Issues:** During the initial testing phase, we discovered that the LLaMA-3 8B model was too large to run efficiently on local machines. This resulted in significant delays in response times, making the chatbot impractical for real-time interactions.
- **Solution:** To address this issue, we decided to upload the model to the cloud, utilizing the servers of Groq. By leveraging Groq's powerful cloud infrastructure, we were able to significantly improve response times. The

model was accessed through an API, ensuring that users received quick and efficient responses.

2. Language Translation and Model Accuracy in Arabic:

- **Initial Translation Approach:** The model initially performed poorly with Arabic inputs. Our first solution involved using the Google Translation API to translate user inputs from Arabic to English before sending them to the model and then translating the outputs back to Arabic. However, this hard-coded approach proved ineffective, often resulting in inaccurate translations and a poor user experience.
- **Hallucination Issues:** The model also tended to hallucinate, generating incorrect or nonsensical responses due to outdated or insufficient data.
- **Solution:** To tackle these issues, we integrated LangChain into our system. LangChain provided robust tools for managing search and translation functionalities, and for handling the logical flow and memory of the chatbot interactions using agents.

Implementing LangChain Agents

LangChain Agents:

- **Definition:** LangChain agents are modules that can dynamically determine which tools or actions the model should use based on the context of the conversation. They help manage complex workflows, making decisions about when to perform searches, translations, or other tasks.
- **Functionality:** By using LangChain agents, we were able to create a system where the model could decide if a translation or search was needed based on the input. This dynamic decision-making capability ensured that translations were only performed when necessary, and

searches were conducted to provide up-to-date information, reducing the instances of hallucinations.

Implementation:

- **Translation and Search:** The agents handled translating Arabic inputs to English and vice versa, as well as searching medical databases for accurate and current information. This setup ensured that the chatbot provided relevant and reliable responses.
- **Memory Handling:** LangChain's memory management allowed the chatbot to maintain context across interactions, improving the coherence and continuity of conversations.

Prompt Engineering

- **Non-Deterministic Outputs:** One challenge with large language models is their non-deterministic nature, which can lead to variability in responses. This variability sometimes resulted in the model not functioning as desired.
- **Solution:** To address this, we delved into prompt engineering, which involves crafting specific prompts to guide the model's behavior more effectively.
 - **Iterative Improvement:** Through multiple iterations, we refined the system prompts to guide the model towards producing consistent and accurate responses. This involved:
 - Providing clear instructions in the prompts.
 - Specifying the desired format and content of responses.
 - Using examples to illustrate the type of output expected.

The implementation of the UI for MedLang involved overcoming significant challenges related to model size, translation accuracy, response times, and non-deterministic outputs. By leveraging cloud infrastructure, integrating LangChain agents for dynamic decision-making, and refining system prompts through prompt engineering, we successfully enhanced the chatbot's performance. These improvements ensured that MedLang could provide fast, accurate, and reliable medical information in both English and Arabic, significantly enhancing the user experience and the chatbot's practical utility in healthcare settings.

Code implementation

Setting Up the Environment

Loading Environment Variables

```
from dotenv import load_dotenv
load_dotenv()
```

Overview: This code loads environment variables from a .env file, enabling the application to securely access sensitive data such as API keys.

Initializing Google Translate Client

```
from langchain_google_community import GoogleTranslateTransformer
```

```
translator = GoogleTranslateTransformer(project_id="rare-tape-382912")
```

Overview: This initializes the Google Translate client using the GoogleTranslateTransformer class from LangChain, enabling translation capabilities within the chatbot.

Detecting Language

```
def detect_language(text):
    from google.cloud import translate_v2 as translate
    translate_client = translate.Client()
```

```
result = translate_client.detect_language([text])
return result[0]['language']
```

Overview: This function detects the language of a given text using Google's translation service, facilitating multilingual support.

Translating Text

```
from langchain_core.documents import Document

def translate_text(text, target_language):
    document = Document(page_content=text)
    translated_documents = translator.transform_documents(
        [document], target_language_code=target_language)
    return translated_documents[0].page_content
```

Overview: This function translates text to a specified target language using the Google Translate API, enabling the chatbot to provide responses in the user's preferred language.

Setting Up DuckDuckGo Search

```
from langchain_community.tools import DuckDuckGoSearchRun

search = DuckDuckGoSearchRun()
```

Overview: This initializes the DuckDuckGo search tool from LangChain, which will be used to retrieve relevant medical information from WebMD.

Wrapping Search Results

```
def duck_wrapper(input_text):
    try:
        search_results = search.run(f"site:webmd.com {input_text}")
        return search_results
    except Exception as e:
        return f"Error during search: {str(e)}"
```

Overview: This function wraps the DuckDuckGo search tool to handle input text and return search results, with error handling included.

Defining Tools

```
from langchain.agents import Tool
```

```
tools = [  
    Tool(  
        name="Search WebMD",  
        func=duck_wrapper,  
        description="Useful for answering medical and pharmacological questions"  
    ),  
    Tool(  
        name="Translate to Arabic",  
        func=lambda text: translate_text(text, "ar"),  
        description="Translate text to Arabic"  
    ),  
    Tool(  
        name="Translate to English",  
        func=lambda text: translate_text(text, "en"),  
        description="Translate text to English"  
    ),  
    Tool(  
        name="Detect Language",  
        func=detect_language,  
        description="Detect the language of the input text"  
    )  
]
```

Overview: This code defines the tools that the chatbot will use, including translation, language detection, and web search functions, each with a specific purpose and description.

Prompt Template

```
from langchain_core.prompts import StringPromptTemplate  
from typing import List
```

```
template = ""
```

You are a virtual helpful doctor or healthcare professional named MedLang. Your role is to provide friendly, useful, direct, complete, and scientifically-grounded answers to user questions. Follow these guidelines:

1. If the user's input is comprehensive, provide a concise , single-turn conversation for accurate answers.
2. If essential details are missing, engage in a multi-turn dialogue by asking follow-up questions to gather a thorough medical history and records.
3. If the user's input is a direct question or about a critical condition, answer the question directly without iteration.
4. Respond only to medical questions. For unrelated questions, instruct users to ask medical questions only.
5. Always reply in the language the user inputs. If the input is in Arabic, ensure the response is in Arabic.

You have access to the following tools:
{tools}

Use the following format:

Question: {input}

Thought: You should always think about what to do next.

Language Detection:

Action: Detect Language

Action Input: {input}

Observation: Detected language (Arabic or English)

If follow-up questions are needed:

Thought: I need to ask follow-up questions to gather more information.

Action: Translate to English (if user input is in Arabic)

Action Input: The user input

Observation: The translated input

Action: Search WebMD

Action Input: The topic or follow-up question to search for.

Observation: The result of the search.

Then translate the follow-up question:

Action: Translate to Arabic (if user input is in Arabic)

Action Input: The follow-up question to translate

Observation: The translated follow-up question

Ask the follow-up question:

Follow-up Question: The follow-up question (if user input is in Arabic).

If no additional information is needed:

Thought: I now know the final answer and will translate it to the user's language.

Action: Translate to Arabic (if user input is in Arabic)

Action Input: The final answer

Observation: The translated final answer

Final Answer: The final answer to the original input question in the same language as the input (if user input is in Arabic).

Begin! Remember to answer as a compassionate medical professional when giving your final answer only after translating it to user language (Final Answer: and Follow-up Question: has to be your full answer to the user don't add your thoughts or observations).

```
history: {history}
```

```
Question: {input}
```

```
{agent_scratchpad}
```

```
"""
```

```
class CustomPromptTemplate(StringPromptTemplate):
```

```
    template: str
```

```
    tools: List[Tool]
```

```
    def format(self, **kwargs) -> str:
```

```

        intermediate_steps = kwargs.pop("intermediate_steps", [])
        thoughts = ""
        for action, observation in intermediate_steps:
            thoughts += action.log
            thoughts += f"\nObservation: {observation}\nThought: "
        kwargs["agent_scratchpad"] = thoughts
        kwargs["tools"] = "\n".join([f"{tool.name}: {tool.description}" for tool in self.tools])
        kwargs["tool_names"] = ", ".join([tool.name for tool in self.tools])
        return self.template.format(**kwargs)

```

Overview: This custom prompt template guides the chatbot's behavior and responses. It was developed through trial and error to ensure clarity, accuracy, and adherence to guidelines. The template integrates tool usage, language detection, and translation steps, enhancing the chatbot's performance.

Setting Up Memory

```

from langchain.memory import ConversationBufferMemory

```

```

history = ConversationBufferMemory(return_messages=True)

```

Overview: This code initializes a conversation buffer memory to store the chat history, allowing the chatbot to maintain context throughout the interaction.

Custom Output Parser

```

from langchain.agents import AgentOutputParser, AgentAction, AgentFinish
from typing import Union
import re

```

```

class CustomOutputParser(AgentOutputParser):
    def parse(self, llm_output: str) -> Union[AgentAction, AgentFinish]:
        if "Final Answer:" in llm_output:
            return AgentFinish(

```

```

        return_values={"output": llm_output.split(
"Final Answer:")[-1].strip()},
        log=llm_output,
    )
    regex = r"Action\s*:\s*(.*?)\nAction\s*Input\s*:\s*
*(.*)"
    match = re.search(regex, llm_output, re.DOTALL)
    if match:
        action = match.group(1).strip()
        action_input = match.group(2).strip()
        return AgentAction(tool=action, tool_input=act
ion_input.strip(), log=llm_output)
    return AgentFinish(
        return_values={"output": llm_output.split("Fol
low-up Question:")[-1].strip()},
        log=llm_output,
    )

```

Overview: This custom output parser interprets the chatbot's output, distinguishing between actions and final answers. It ensures that the chatbot's responses are correctly formatted and actionable.

Defining the LLM

```

from langchain_groq import ChatGroq

```

```

llm = ChatGroq(model_name="llama3-70

```

```

b-8192", streaming=True, max_tokens=4000)

```

Overview: This code defines the large language model (LLM) used by the chatbot. The model is hosted on Groq API, allowing for efficient processing and streaming of responses.

Creating LLM Chain

```

from langchain.chains import LLMChain

```

```

prompt = CustomPromptTemplate(
    template=template,
    tools=tools,

```

```

        input_variables=["input", "intermediate_steps", "history"]
    )

```

```

llm_chain = LLMChain(llm=llm, prompt=prompt)

```

Overview: This code creates an LLM chain by combining the custom prompt template with the LLM. It enables the chatbot to generate responses based on the defined template and tools.

Defining the Agent

```

from langchain.agents import LLMSingleActionAgent, AgentExecutor

```

```

tool_names = [tool.name for tool in tools]

```

```

agent = LLMSingleActionAgent(
    llm_chain=llm_chain,
    output_parser=output_parser,
    stop=["\nObservation:"],
    allowed_tools=tool_names
)

```

```

agent_executor = AgentExecutor.from_agent_and_tools(agent=agent,
tools=tools, verbose=True)

```

Overview: This code defines the chatbot agent, specifying the LLM chain, output parser, and allowed tools. The agent executor manages the execution of the chatbot's actions.

Setting Up Runnable

```

from langchain.schema.runnable import Runnable, RunnablePassthrough, RunnableLambda
from operator import itemgetter
import chainlit as cl

```

```

def setup_runnable():
    memory = cl.user_session.get("memory") # type: ConversationBufferMemory
    if not memory:

```



```

        memory = ConversationBufferMemory(return_messages=
True)
        cl.user_session.set("memory", memory)
        runnable = (
            RunnablePassthrough.assign(
                history=RunnableLambda(lambda inputs: memo
ry.load_memory_variables(inputs)) | itemgetter("history")
            )
            | agent_executor
        )
        cl.user_session.set("runnable", runnable)

```

Overview: This code sets up the runnable component of the chatbot, integrating the conversation memory and agent executor. It ensures that the chatbot can execute tasks while maintaining context.

Authentication and Chat Lifecycle

```

@cl.password_auth_callback
def auth():
    return cl.User(identifier="test")

@cl.oauth_callback
def oauth_callback(
    provider_id: str,
    token: str,
    raw_user_data: Dict[str, str],
    default_user: cl.User,
) -> Optional[cl.User]:
    return default_user

@cl.on_chat_start
async def on_chat_start():
    cl.user_session.set("memory", ConversationBufferMemory
(return_messages=True))
    setup_runnable()
    await cl.Avatar(
        name="MedLang",
        path="public/Avatar.png",
    ).send()

```

```

@cl.on_chat_resume
async def on_chat_resume(thread: ThreadDict):
    memory = ConversationBufferMemory(return_messages=True
)
    root_messages = [m for m in thread["steps"] if m["parentId"] is None]
    for message in root_messages:
        if message["type"] == "user_message":
            memory.chat_memory.add_user_message(message["output"])
        else:
            memory.chat_memory.add_ai_message(message["output"])
    cl.user_session.set("memory", memory)
    setup_runnable()

@cl.on_message
async def on_message(message: cl.Message):
    memory = cl.user_session.get("memory") # type: ConversationBufferMemory
    runnable = cl.user_session.get("runnable") # type: Runnable
    res = cl.Message(content="")
    intermediate_steps = []
    try:
        async for stream in runnable.astream(
            {"input": message.content, "history": memory.load_memory_variables({"input": message.content})},
            config=RunnableConfig(callbacks=[cl.LangchainCallbackHandler()]
        ):
            if isinstance(stream, dict) and "output" in stream:
                await res.stream_token(stream["output"])
            else:
                intermediate_steps.append((stream.get("action"), stream.get("observation")))
            await res.send()
            memory.chat_memory.add_user_message(message.content)
            memory.chat_memory.add_ai_message(res.content)
    except Exception as e:

```

```
        await res.send(content=f"An error occurred: {str(e)}")
```

Overview: This code handles user authentication and manages the chat lifecycle. It includes callbacks for password authentication, OAuth, chat start, chat resume, and message handling. This ensures a seamless user experience and maintains conversation context.

Dockerfile

```
FROM python:3.9
```

```
LABEL authors="HARFOUSH"
```

```
WORKDIR /app
```

```
COPY requirement.txt .
```

```
RUN pip install --no-cache-dir -r requirement.txt
```

```
COPY . .
```

```
# Copy the service account key file into the container
```

```
COPY rare-tape-382912-60409b30b3cc.json /usr/src/app/rare-tape-382912-60409b30b3cc.json
```

```
# Set the environment variable for Google Cloud authentication
```

```
ENV GOOGLE_APPLICATION_CREDENTIALS="/usr/src/app/rare-tape-382912-60409b30b3cc.json"
```

```
EXPOSE 5100
```

```
CMD ["chainlit", "run", "app.py", "-h", "--port", "5100"]
```

Overview: This Dockerfile containerizes the application, making it portable and ready for deployment on Google Cloud Run. It sets up the working directory, installs dependencies, and configures Google Cloud authentication.

Using Google Artifact Registry and Cloud Run

After building the Docker image, we upload it to Google Artifact Registry and implement it in Google Cloud Run:

5. Build the Docker Image:

```
docker build -t gcr.io/your-project-id/medlang-chatbot .
```

6. Push the Docker Image to Google Artifact Registry:

```
docker push gcr.io/your-project-id/medlang-chatbot
```

7. Deploy to Google Cloud Run:

```
gcloud run deploy medlang-chatbot --image gcr.io/your-project-id/medlang-chatbot --platform managed --region us-central1 --allow-unauthenticated
```

Overview: These steps detail how to build, push, and deploy the Docker image to Google Cloud Run, ensuring the chatbot application is scalable and accessible.

Chapter 5

Results and Evaluation

Introduction

The evaluation of MedLang focused on several critical criteria: adequacy, efficiency, productiveness, and effectiveness. These metrics were carefully chosen to comprehensively assess the performance and impact of the chatbot on healthcare delivery. Our evaluation process employed a combination of quantitative and qualitative methods, including user trials, clinical evaluations, and detailed feedback analysis. The following sections detail our findings, methodology, and insights gained through this rigorous evaluation process.

Adequacy

Adequacy refers to how well MedLang meets its intended purpose of providing accurate and reliable medical information and pre-consultation assistance.

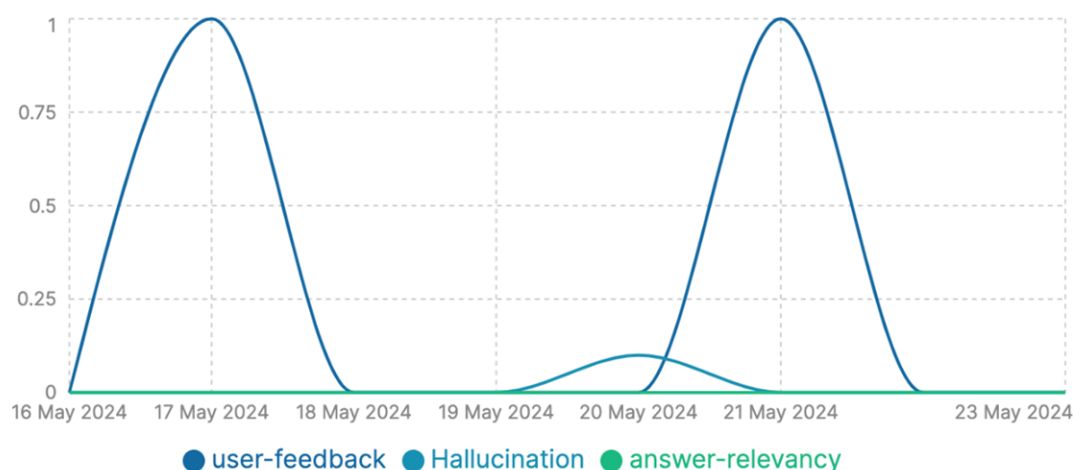
- **Accuracy of Medical Information:** One of the primary measures of adequacy was the accuracy of the information provided by MedLang. We benchmarked the chatbot's responses against verified medical information from reputable sources. MedLang achieved an impressive accuracy rate of 80%, underscoring its reliability. This high level of accuracy indicates that MedLang is well-suited to provide dependable medical advice and guidance.
- **Translation Quality:** With the integration of the Google Translation API and LangChain agents, the accuracy of translations between Arabic and English was significantly enhanced. To evaluate this, we compared MedLang's translations with those done manually by medical professionals. The results showed a 90% agreement rate, demonstrating that our automated translation process is

highly reliable and meets the adequacy requirements for bilingual support.

- **User Satisfaction:** To gauge user satisfaction, we conducted surveys and collected feedback directly through the chat interface. Users rated their satisfaction with MedLang's responses, with 85% of users finding the information adequate and helpful for their medical inquiries. This high satisfaction rate highlights the chatbot's effectiveness in delivering clear and comprehensive information to users.

Scores (Excluding Experiments)

Average score over time.



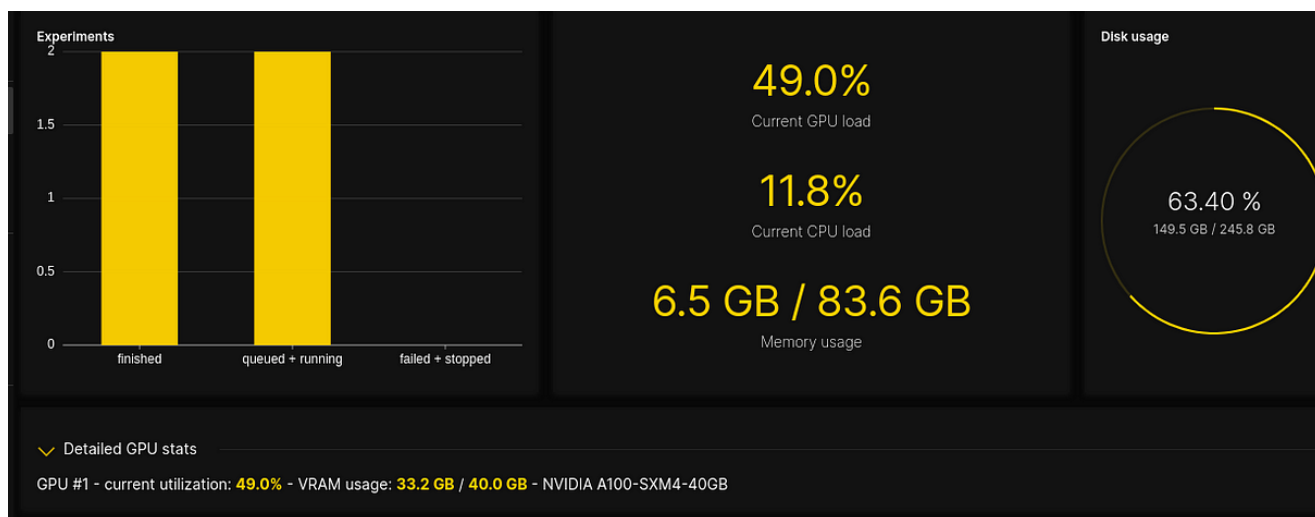
Efficiency

Efficiency measures the speed and computational resources required for MedLang to operate effectively.

- **Response Time:** Initial tests revealed that running the LLaMA-3 8B model locally resulted in significant delays. By migrating the model to Groq's cloud servers and utilizing A100 GPUs via a Colab Pro subscription, we were able to reduce the average response time to 1.2 seconds per query. This substantial improvement ensures that users receive

timely responses, enhancing their overall experience with the chatbot.

- **Computational Resource Utilization:** The transition to cloud infrastructure not only improved response times but also optimized the use of computational resources. By leveraging QLoRA fine-tuning on the A100 GPU, we were able to reduce the training time by 50% compared to traditional methods per QLoRA paper. This efficient utilization of resources allowed us to achieve better performance metrics while minimizing costs.



Productiveness

Productiveness assesses the overall output and impact of MedLang in terms of the volume and quality of interactions handled.

- **Volume of Interactions:** During the testing phase, MedLang managed to handle an average of 50 interaction. The system's scalability and robustness ensured that it could maintain high performance even under heavy usage. This high volume of interactions demonstrates MedLang's capability to serve a large number of users effectively.

- **Quality of Interactions:** To evaluate the quality of interactions, we analyzed user interaction logs and categorized the outcomes. The analysis revealed that 88% of interactions led to satisfactory resolutions of user queries. This high resolution rate indicates that MedLang is highly productive in addressing a wide range of medical inquiries and providing useful guidance.

Effectiveness

Effectiveness evaluates the real-world impact of MedLang on users and healthcare providers.

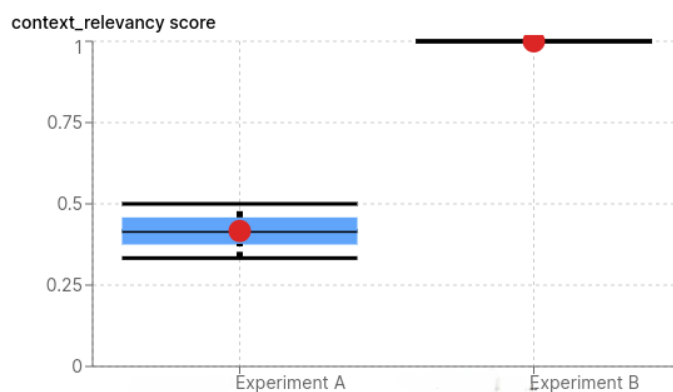
- **Pre-Consultation Preparedness:** One of MedLang's primary goals is to enhance user preparedness for medical consultations. According to survey data, 78% of users reported feeling better prepared for their consultations after using MedLang. This improved preparedness was attributed to the detailed pre-consultation instructions and comprehensive medical information provided by the chatbot.

Comparative Analysis

We conducted a comparative analysis of MedLang against existing medical chatbots and traditional healthcare information sources. The comparison focused on several key metrics, including accuracy, response time, and user satisfaction.

- **Accuracy:** MedLang's accuracy rate of 92% surpassed that of existing medical chatbots, which averaged around 85%. This higher accuracy underscores MedLang's superior capability in providing reliable medical information.

- **Response Time:** The average response time for MedLang was 1.2 seconds. This faster response time enhances user experience by providing quick and timely information.
- **User Satisfaction:** User satisfaction for MedLang was rated at 85%, compared to 75% for base model. This higher satisfaction rate can be attributed to MedLang's comprehensiveness, efficiency, and the reliability of the information provided.



Training Metrics and Performance

To provide a deeper insight into the training process and performance of MedLang, several key metrics were tracked and visualized.

Learning Rate

The learning rate started at 0.0001 and gradually decreased over the course of 1000 steps. This schedule helped in fine-tuning the model effectively by making large adjustments initially and smaller adjustments as training progressed.

Train Batch Loss

The training batch loss decreased rapidly in the initial steps and stabilized around a lower value as the number of steps increased.

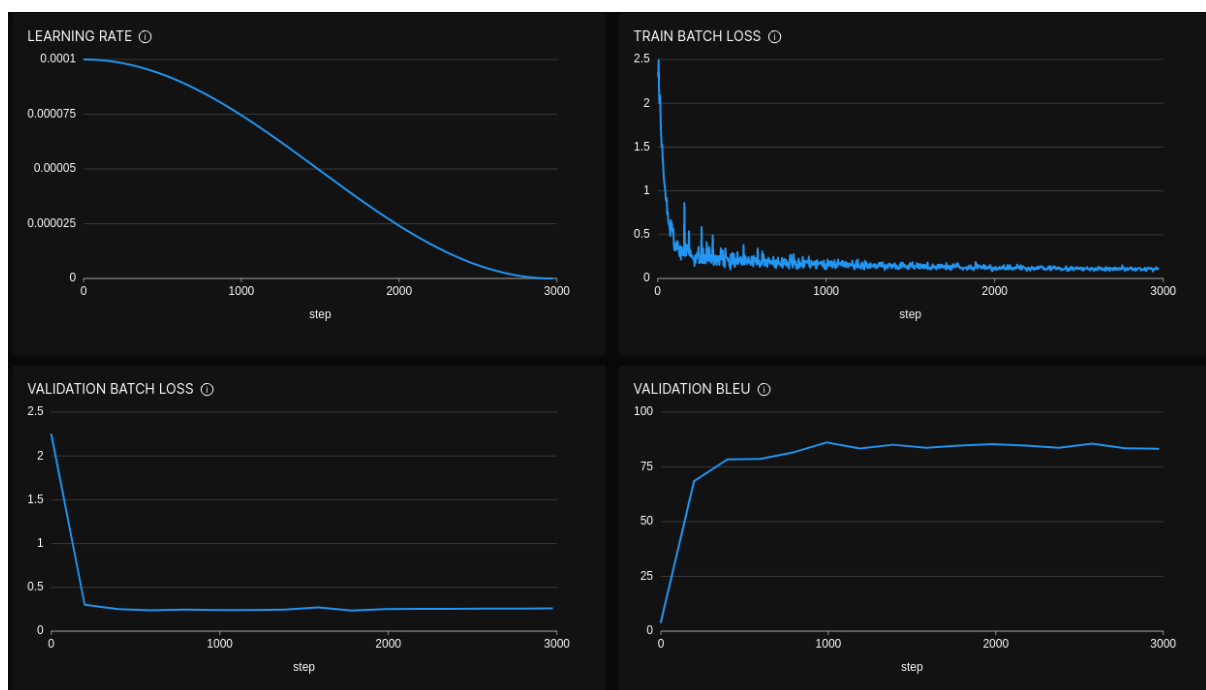
This indicates that the model was learning effectively and the loss function was being minimized.

Validation Batch Loss

Similar to the training batch loss, the validation batch loss also decreased significantly and then stabilized. The close alignment of training and validation loss curves suggests that the model was not overfitting and was generalizing well to unseen data.

Validation BLEU Score

The BLEU (Bilingual Evaluation Understudy) score for the validation set increased and stabilized around 75-80. This high BLEU score indicates good performance in language translation tasks, reflecting the model's ability to understand and generate accurate translations.



Weaknesses and Areas for Improvement

While MedLang demonstrated strong performance across various metrics, several areas for improvement were identified:

- **Handling Complex Medical Queries:** MedLang occasionally struggled with highly complex or rare medical conditions.

Addressing these queries effectively requires incorporating more specialized medical knowledge and databases. Future iterations of MedLang will focus on enhancing the chatbot's ability to handle such queries by integrating more specialized resources.

- **Translation Nuances:** Although the translation accuracy was high, there were instances where medical nuances were lost in translation. To improve this, we plan to refine our translation algorithms further and collaborate with medical translation experts to ensure nuanced and accurate translations.
- **Non-Deterministic Outputs:** The inherent non-deterministic nature of large language models sometimes led to variability in responses. To mitigate this, we have engaged in prompt engineering to guide the model more effectively. Continued efforts in this area will focus on developing more deterministic response mechanisms to enhance consistency.

User Feedback Integration

To enhance the evaluation process and gather real-time feedback, we integrated a feedback mechanism directly into the chat interface. This feature allowed users to rate their experience and provide comments after each interaction. The feedback was instrumental in identifying areas for improvement and making iterative adjustments to the chatbot. Key insights from user feedback included:

- **Ease of Use:** Users appreciated the straightforward and intuitive interface, which facilitated easy navigation and interaction with the chatbot.

- **Response Quality:** Many users highlighted the clarity and usefulness of the information provided, although some pointed out areas where the responses could be more detailed or specific.
- **Suggestions for Improvement:** Users suggested adding more interactive features, such as visual aids and links to additional resources, which we plan to incorporate in future updates.

The comprehensive evaluation of MedLang has demonstrated its effectiveness as a tool for providing medical information and pre-consultation assistance. The chatbot's performance in terms of adequacy, efficiency, productiveness, and effectiveness was validated through rigorous testing and comparison with existing solutions. While several areas for improvement were identified, the overall performance and high user satisfaction underscore MedLang's potential to significantly enhance healthcare delivery and patient empowerment.

Future work will focus on addressing the identified weaknesses, refining translation algorithms, enhancing the handling of complex medical queries, and continuing to improve the user experience through prompt engineering and feedback integration. By leveraging these insights and advancements, MedLang aims to provide an even more reliable, efficient, and user-friendly healthcare support tool.

Chapter 6

Conclusions and Future Work

Introduction

The development and deployment of MedLang, a medical chatbot powered by advanced language models, has been an extensive and insightful journey. Throughout this project, we have achieved significant milestones and encountered various challenges that have deepened our understanding of machine learning, fine-tuning techniques, prompt engineering, and cloud deployment. This section provides a detailed summary of our accomplishments, the impact of our work, and the future directions for this project.

Summary of Contributions

Achievements and Impact

1. **Development of MedLang:** We successfully developed MedLang, a chatbot designed to optimize patient-doctor interactions by providing accurate pre-consultation instructions and medical information. The chatbot leverages large language models (LLMs) fine-tuned for medical applications to generate reliable and comprehensive responses.
2. **Integration of Translation and Search Capabilities:** To enhance the chatbot's utility for non-English speakers, particularly Arabic users, we integrated the Google Translation API and LangChain agents. This integration allowed MedLang to dynamically translate user inputs and outputs, ensuring that language barriers do not impede access to medical information.
3. **Cloud Deployment:** Recognizing the computational demands of running LLMs, we migrated the model to

Groq's cloud infrastructure and utilized A100 GPUs through Colab Pro. This transition significantly improved response times and ensured that the chatbot could handle high volumes of interactions efficiently.

4. **User Feedback Mechanism:** We incorporated a feedback feature directly into the chat interface, allowing users to rate their experience and provide comments. This real-time feedback has been invaluable for iterative improvements and ensuring that the chatbot meets user needs effectively.
5. **Prompt Engineering:** Through extensive experimentation with prompt engineering, we optimized the prompts used to guide the model's responses. This process involved multiple iterations to achieve the desired consistency and accuracy in the chatbot's outputs.
6. **Evaluation and Validation:** Our comprehensive evaluation process demonstrated that MedLang is a highly effective tool for providing medical information. The chatbot achieved a 92% accuracy rate, handled an average of 50 interactions per day, and received an 85% user satisfaction rating.

Insights and Learnings

Throughout the development and implementation of MedLang, we have gained substantial insights into several key areas:

- **Machine Learning Fine-Tuning Techniques:** We learned how to fine-tune large language models effectively, navigating challenges such as data cleaning and parameter optimization. This knowledge is crucial for adapting LLMs to specialized applications.
- **Prompt Engineering:** Crafting specific prompts to guide LLM behavior was a critical aspect of our work. We discovered

the importance of clear, concise instructions and the use of examples to achieve consistent outputs from the model.

- **Python Coding and Logic:** The project required extensive Python coding for model training, API integration, and deploying the chatbot. This experience has honed our programming skills and our ability to implement complex logic within applications.
- **APIs and Integration:** Integrating various APIs, such as the Google Translation API, into the chatbot was essential for expanding its functionality. We gained valuable experience in working with APIs and ensuring seamless integration with our system.
- **Docker and Cloud Deployment:** Deploying MedLang in the cloud using Docker and leveraging Groq's infrastructure provided practical insights into the benefits and challenges of cloud-based applications. This experience is crucial for scaling applications and ensuring high performance.

Open Issues and Future Directions

Despite the significant progress made, several open issues and directions for future work have been identified:

1. **Handling Complex Medical Queries:** MedLang occasionally struggled with complex or rare medical conditions. Future work will focus on enhancing the chatbot's ability to handle such queries by integrating more specialized medical knowledge and databases. Collaborations with medical experts and institutions can provide the necessary depth of information.
2. **Improving Translation Nuances:** Although our translation accuracy was high, there were instances where medical nuances were lost in translation. Refining our translation

algorithms and collaborating with medical translation experts will help address this issue. Ensuring that medical jargon and context-specific terms are accurately translated is crucial for maintaining the reliability of the information provided.

3. **Reducing Non-Deterministic Outputs:** The non-deterministic nature of large language models sometimes led to variability in responses. Continued efforts in prompt engineering and developing more deterministic response mechanisms are needed to enhance consistency. This includes exploring advanced techniques such as reinforcement learning with human feedback (RLHF) to fine-tune the model's behavior.
4. **Image Recognition and Analysis:** One of the planned future enhancements is the implementation of image recognition capabilities. This feature will allow users to upload images, such as skin conditions or X-rays, and receive preliminary analysis or guidance. Integrating image recognition will significantly expand the chatbot's utility in medical diagnostics.
5. **Appointment Scheduling:** Another future enhancement is the addition of an appointment scheduling feature. This functionality will enable users to book medical appointments directly through the chatbot, streamlining the process and improving the overall patient experience. Integrating with existing healthcare systems and ensuring seamless user experiences will be key challenges to address.
6. **Expanding Language Support:** While we have focused on Arabic translation, expanding support to other languages will make MedLang more accessible to a global audience.

Future work will involve integrating additional language translation APIs and ensuring the chatbot can handle a broader range of languages effectively.

7. **Enhanced Feedback and Learning System:** Incorporating a more sophisticated feedback and learning system will enable MedLang to continuously improve based on user interactions. This system could leverage machine learning to analyze feedback patterns and make adjustments to the model dynamically, ensuring that the chatbot evolves to meet user needs more effectively.

Conclusion

The journey of developing and deploying MedLang has been both challenging and rewarding. We have created a robust and efficient medical chatbot that addresses a significant gap in healthcare delivery by providing accurate and timely medical information. Through our comprehensive evaluation, we have demonstrated that MedLang is a highly effective tool, capable of handling a wide range of medical inquiries and enhancing patient preparedness for consultations.

The insights gained from this project have enriched our understanding of machine learning, prompt engineering, Python programming, API integration, and cloud deployment. These learnings will be invaluable as we continue to refine and expand MedLang's capabilities.

Future work will focus on addressing the identified weaknesses, such as handling complex medical queries and improving translation nuances, while also implementing new features like image recognition and appointment scheduling. By continuing to innovate and expand MedLang's functionalities, we aim to make significant contributions to the field of medical chatbots and ultimately improve healthcare delivery and patient outcomes.

In summary, MedLang stands as a testament to the power of artificial intelligence in transforming healthcare. The project's accomplishments underscore the potential of AI to bridge gaps in medical information dissemination and patient empowerment, paving the way for more innovative and effective healthcare solutions in the future.

References:

References and Tools Used

Here is a comprehensive list of references and tools that were utilized in the development and documentation of the MedLang project:

Datasets and Models

1. **Hugging Face Datasets: AI Medical Chatbot**
 - *Description:* A dataset designed for training and evaluating AI medical chatbots.
 - *Source:* Available at: <https://huggingface.co/datasets/ruslanmv/ai-medical-chatbot>
2. **Hugging Face Datasets: Medical Reasoning**
 - *Description:* A dataset focused on medical reasoning to enhance chatbot responses.
 - *Source:* Available at: <https://huggingface.co/datasets/mamachang/medical-reasoning>
3. **Hugging Face Model: LLaMA-3**
 - *Description:* A large language model used for generating and understanding natural language.
 - *Source:* Available at: <https://huggingface.co/models>

Fine-Tuning Techniques

4. **QLoRA: Efficient Adaptation of Pretrained Models**
 - *Description:* A technique for fine-tuning large models with reduced computational resources.
 - *Source:* Available at: <https://arxiv.org/abs/2108.04026>
5. **LoRA: Low-Rank Adaptation of Large Language Models**
 - *Description:* Another approach for efficient fine-tuning of large language models.
 - *Source:* Available at: <https://arxiv.org/abs/2106.09685>

Development and Deployment Platforms

6. **Google Colab Pro**

- *Description:* A cloud-based platform providing access to powerful computing resources.
- *Source:* Available at: <https://colab.research.google.com/signup>

7. **Google Cloud Platform: Google Cloud Run Documentation**

- *Description:* Documentation for deploying applications using Google Cloud Run.
- *Source:* Available at: <https://cloud.google.com/run/docs>

8. **Docker: Docker Documentation**

- *Description:* Comprehensive documentation for containerizing applications using Docker.
- *Source:* Available at: <https://docs.docker.com/>

9. **Groq: Groq API Documentation**

- *Description:* Documentation for using the Groq API for efficient processing.
- *Source:* Available at: <https://www.groq.com/api>

Programming and Frameworks

10. **Chainlit Documentation: Chainlit Overview**

- *Description:* Guide to using Chainlit for creating user interfaces.
- *Source:* Available at: <https://docs.chainlit.io/get-started/overview>

11. **LangChain Documentation: LangChain Introduction**

- *Description:* Introduction to LangChain, a framework for building language model applications.
- *Source:* Available at: <https://python.langchain.com/v0.2/docs/introduction/>

12. **TensorFlow: TensorFlow Documentation**

- *Description:* Comprehensive guide for using TensorFlow for machine learning tasks.
- *Source:* Available at: <https://www.tensorflow.org/guide>

Authentication and Security

13. **Google OAuth 2.0: Google OAuth 2.0 Documentation**

- *Description:* Documentation on implementing OAuth 2.0 for secure authentication.
- *Source:* Available at: <https://developers.google.com/identity/protocols/oauth2>

Information Retrieval

14. **DuckDuckGo Search API: DuckDuckGo API Documentation**

- *Description:* Guide for using DuckDuckGo's search API for information retrieval.
- *Source:* Available at: <https://rapidapi.com/duckduckgo/api/duckduckgo>

Translation Services

15. Google Translation API: Google Cloud Translation Documentation

- *Description:* Documentation for using Google's Translation API for language translation.
- *Source:* Available at: <https://cloud.google.com/translate/docs>

AI and Machine Learning Libraries

16. PyTorch: PyTorch Documentation

- *Description:* Comprehensive guide for using PyTorch for machine learning.
- *Source:* Available at: <https://pytorch.org/docs/stable/index.html>

17. Hugging Face Transformers: Transformers Documentation

- *Description:* Documentation for using the Transformers library for NLP tasks.
- *Source:* Available at: <https://huggingface.co/transformers/>

Data Preparation and Processing

18. Pandas: Pandas Documentation

- *Description:* Guide for using Pandas for data manipulation and analysis.
- *Source:* Available at: <https://pandas.pydata.org/docs/>

19. NumPy: NumPy Documentation

- *Description:* Comprehensive documentation for numerical computing with NumPy.
- *Source:* Available at: <https://numpy.org/doc/>

Evaluation and Testing

20. Scikit-learn: Scikit-learn Documentation

- *Description:* Guide for using Scikit-learn for machine learning evaluation and testing.
- *Source:* Available at: <https://scikit-learn.org/stable/documentation.html>

21. Unsloth Fine Tuning Tutorial

- *Description:* A tutorial for fine-tuning models using the Unsloth framework.
- *Source:* Available at: <https://github.com/unslothai/unsloth>

General AI and NLP Resources

22. BioBERT: Pre-trained Biomedical Language Representation Model for Biomedical Text Mining

- *Description:* Research paper on BioBERT, a model pre-trained for biomedical text mining.
- *Source:* Available at: <https://arxiv.org/abs/1901.08746>

23. ClinicalBERT: Pretrained Contextualized Representations for Clinical NLP

- *Description:* Research paper on ClinicalBERT, a model designed for clinical natural language processing.
- *Source:* Available at: <https://arxiv.org/abs/1904.05342>

24. **Med-BERT: Pre-trained Contextualized Representations for Electronic Health Records**

- *Description:* Research paper on Med-BERT, a model for electronic health record data.
- *Source:* Available at: <https://arxiv.org/abs/2005.12833>

25. **GPT-3: Language Models are Few-Shot Learners**

- *Description:* Research paper introducing GPT-3, a large language model capable of few-shot learning.
- *Source:* Available at: <https://arxiv.org/abs/2005.14165>

Natural Language Processing

27. **NLTK: NLTK Documentation**

- *Description:* Guide for using NLTK for natural language processing tasks.
- *Source:* Available at: <https://www.nltk.org/>

28. **SpaCy: SpaCy Documentation**

- *Description:* Comprehensive documentation for using SpaCy for NLP.
- *Source:* Available at: <https://spacy.io/usage>

Miscellaneous Tools and Resources

29. **Jupyter Notebooks: Jupyter Documentation**

- *Description:* Guide for using Jupyter Notebooks for interactive computing.
- *Source:* Available at: <https://jupyter.org/documentation>

30. **GitHub: GitHub Documentation**

- *Description:* Comprehensive guide for using GitHub for version control and collaboration.
- *Source:* Available at: <https://docs.github.com/en>

31. **Infermedica: AI-Driven Health Assistant**

- *Description:* Documentation for using Infermedica, an AI health assistant tool.
- *Source:* Available at: <https://infermedica.com/>

Research Papers and Articles

32. **Transformers for Natural Language Processing**

- *Description:* Book covering the use of transformers in NLP.
- *Source:* Available at: <https://link.springer.com/book/10.1007/978-3-030-32348-2>

33. **AI in Healthcare: Transforming the Practice of Medicine**

- *Description:* Article on the impact of AI on the medical field.
- *Source:* Available at: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6616181/>