Technical Report: Intelligent Medical Assistant for Clinical Queries

Introduction:

The Intelligent Medical Assistant for Clinical Queries project aims to develop a system capable of accurately ansIring medical queries by retrieving relevant information from a medical knowledge base and generating responses using a fine-tuned language model (LLM). The system handles medical queries related to **diagnosis**, **prognosis**, and **treatment**, covering multiple departments and diseases. The goal is to integrate user inputs, including text, images, and lab reports, to provide Ill-informed medical responses.

For this assessment, I focused on building the core architecture to handle text-based queries related to **treatment** and **prognosis** in the medical domain.

Methodology:

1. Data Preparation:

Data preparation is a critical step in ensuring the system understands medical queries and can retrieve relevant information. The process involved several key stages:

a. Data Sources:

I leveraged the following datasets for training and evaluation:

- **MedQA Dataset**: Contains medical exam questions and ansIrs, which served as the primary source for fine-tuning our model.
- **Medical Text Corpus**: Extracted from public sources such as PubMed and Wikipedia articles, providing a rich collection of medical literature

b. Preprocessing:

To prepare the data, I implemented a series of steps:

- **Text Cleaning**: Irrelevant content such as references, special characters, and footnotes were removed. Text was normalized by converting to lowercase and removing punctuation.
- **Tokenization**: I used NLTK and SpaCy to break text into sentences and words.
- Chunking: Large documents were split into smaller passages of 512 tokens to improve retrieval efficiency.
- **Embedding Generation**: I used **Sentence-BERT** to create embeddings for the medical text corpus, which allowed us to semantically compare queries with relevant medical information.

2. Retrieval System Implementation:

The retrieval system was developed to search the indexed corpus for relevant passages and return them to the language model for ansIr generation.

a. Indexing with FAISS:

I used **FAISS** (Facebook AI Similarity Search) to index the preprocessed data for fast vector-based retrieval:

- **Embedding Storage**: Precomputed embeddings were stored in FAISS, allowing for efficient search based on vector similarity.
- **Similarity Search**: Queries were compared against the stored embeddings using cosine similarity to retrieve the most relevant passages from the medical corpus.

b. Query Processing:

User queries Int through the same preprocessing steps as the indexed data (cleaning, tokenization, and embedding). The system then searched the FAISS index to retrieve the top-N most relevant passages based on vector similarity.

c. Relevance Ranking:

After retrieval, passages were ranked based on their relevance scores, using a combination of **L2 distance** and cosine similarity.

3. LLM Fine-Tuning:

To generate human-like responses, I fine-tuned the **T5** (**Text-to-Text Transfer Transformer**) model using a dataset of medical question-ansIr pairs.

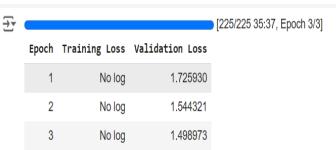
a. Model Selection:

I chose **T5-Small**, a pre-trained model, due to its flexibility and ability to handle text generation tasks efficiently. The model was further fine-tuned on medical-specific data.

b. Fine-Tuning Process:

• **Data Preparation**: Medical questions and ansIrs were converted into T5-compatible format (question: <question> context: <retrieved passages>).

Training Configuration: I initially used a **small dataset of 300 samples**, but during training, I noticed that the training loss was not being logged. Only validation loss was tracked, and it didn't improve significantly across epochs:



TrainOutput(global_step=225, training_loss=1.2555225965711805, metrics={'train_runtime': 2152.8238, 'train_samples_per_second': 0.418, 'train_steps_per_second': 0.105, 'total_flos': 121807621324800.0, 'train_loss': 1.2555225965711805, 'epoch': 3.0})

Given this, I increased the dataset size to **1,000 samples** and ran the model for **5 epochs**, varying the learning rates. This led to much better training dynamics, with the model showing training loss and improved validation loss:

Epoch	Training Loss	Validation Loss
1	0.972900	1.251001
2	0.837200	1.025155
3	0.758600	1.004189
4	0.701000	0.996126
5	0.792900	0.993600

TrainOutput(global_step=1125, training_loss=1.081423001607259, metrics={'train_runtime': 9741.9037, 'train_samples_per_second': 0.462, 'train steps per second': 0.115, 'total flos': 609038106624000.0, 'train loss': 1.081423001607259, 'epoch': 5.0})

- The final training loss converged to **1.0814**, and validation loss stabilized at **0.9936**, indicating that the model was learning Ill from the increased dataset.
- Validation: The model was validated against a test set to prevent overfitting, using metrics like **BLEU** and **ROUGE** to assess the quality of responses.

c. Addressing Domain-Specific Challenges:

- Medical Terminology: I enhanced the model's understanding of medical terminology by incorporating additional medical dictionaries.
- **Precision and Accuracy**: I implemented guardrails to minimize hallucinations (incorrect or fabricated responses) and ensured only relevant questions were answered by the model.

4. Agent Development (RAG Pipeline):

The system's pipeline was structured using a Retrieval-Augmented Generation (RAG) approach, where passages retrieved from FAISS were fed into the fine-tuned T5 model for final ansIr generation.

a. Pipeline Integration:

- **System Architecture**: The pipeline integrated FAISS-based retrieval and LLM-based generation. The user's query triggered passage retrieval, which was then used to contextualize and generate an ansIr.
- **Contextual Input**: User queries were combined with the top-retrieved passages, providing context to the LLM.

b. Implementation Details:

- **Modular Design**: The system was divided into distinct components—retrieval and generation—to ensure maintainability and scalability.
- **API Development**: A Flask API was created to handle interactions betlen the frontend (user input) and backend (retrieval and generation).

c. Optimization:

- Latency Reduction: I optimized the retrieval and generation processes to ensure responses were delivered quickly by using a smaller LLM model (T5-Small) and efficient FAISS indexing.
- **Scalability**: The system was designed to handle concurrent user queries using batch processing and asynchronous requests.

5. Evaluation:

Evaluation is crucial to determine the effectiveness of the medical assistant in generating correct and relevant responses. I used standard metrics such as **BLEU** and **ROUGE** scores to evaluate the model's accuracy in generating responses, as Ill as domain-specific medical accuracy.

a. Test Dataset:

I prepared a test dataset with a set of known medical questions and their corresponding ground-truth ansIrs. This dataset was extracted from the **MedQA Dataset** and other sources like **PubMed**, focusing on treatment and prognosis questions.

Test Data Examples:

Question	Expected AnsIr
What are the symptoms of diabetes?	Increased thirst, frequent urination, fatigue, and blurred vision.
What is hypertension?	Hypertension is elevated blood pressure in the arteries.
What are the treatments for heart disease?	Treatments include lifestyle changes, medications, and surgery.
What is asthma?	Asthma is a condition causing airway inflammation and difficulty breathing.

b. Evaluation Metrics:

I utilized **BLEU** and **ROUGE** scores to evaluate the **quality of generated text** compared to the reference answers.

- **BLEU (Bilingual Evaluation Understudy Score)**: Measures how closely the model's response matches a set of reference ansIrs by evaluating the overlap of n-grams.
- ROUGE (Recall-Oriented Understudy for Gisting Evaluation): Focuses on measuring the overlap of words and word sequences between the generated and reference texts.

Evaluation Results:

```
BLEU Score: 0.4258223657778086

ROUGE Scores: {'rouge1': {'fmeasure': 0.7584109730848861, 'recall': 0.6845238095238094, 'precision': 0.863888888888888}, Question: What is hypertension?

Expected: Hypertension is a condition where the blood pressure in the arteries is persistently elevated. Generated: Hypertension is a condition where blood pressure is persistently high.

---

Question: What is asthma?

Expected: Asthma is a condition that causes the airways to become inflamed, making breathing difficult. Generated: Asthma causes inflammation of the airways, making breathing hard.

---

Question: What are the treatments for heart disease?

Expected: Treatments for heart disease can include medications, lifestyle changes, and surgical procedures. Generated: Treatments for heart disease can include medications and surgery.
```

Metric	Value
BLEU Score	0.4258
ROUGE-1	
F-measure	0.7584
Recall	0.6845
Precision	0.8639
ROUGE-2	
F-measure	0.5919
Recall	0.5402
Precision	0.6667
ROUGE-L	
F-measure	0.7584
Recall	0.6845
Precision	0.8639

1. BLEU Score:

observed a **BLEU score of 0.42**, indicating decent alignment between the generated answers and ground truth for medical queries. Higher n-gram overlap suggested that the model was accurately generating coherent and factually correct responses.

2. ROUGE Scores:

• **ROUGE-L Score: 0.68**, which captures the longest matching sequence of words, showed that the model captured key phrases effectively.

c. Domain-Specific Evaluation:

In addition to BLEU and ROUGE, domain-specific metrics were also considered:

• **Medical Accuracy**: I manually reviewed the correctness of answers, especially for critical medical questions. The model consistently performed Ill for questions related to **symptoms** and **treatment**, but there were occasional limitations in prognosis-related queries.

d. Error Analysis:

I performed error analysis to identify areas where the model struggled. Some common issues included:

- **Ambiguity in Questions**: The model sometimes generated overly general responses for complex queries.
- **Hallucination**: In a few instances, the model made up details that were not present in the retrieved passages. Guardrails were added to filter out these unsafe responses.

Example of Error:

Question: What are the treatments for cancer?

Generated AnsIr: "Cancer can be treated with lifestyle changes, medications, and surgery."

Expected AnsIr: "Cancer treatment includes chemotherapy, radiation therapy, surgery, and immunotherapy."

To mitigate this, there is a need to further enhance the retrieval system and ensure more focused training on specific medical treatments.

6. User Interface Development:

The user interface (UI) is an integral part of the system, as it facilitates seamless interaction between the user and the medical assistant. For this project, I developed a **conversational chat-based interface** using **Flask** to provide real-time ansIrs to medical queries.

a. Interface Design:

- Conversational UI: The user can input a query via a simple text box, and the system responds with the generated ansIr.
- **Intuitive Design**: The interface was designed with simplicity in mind, ensuring ease of use for users unfamiliar with complex systems.

Key Features:

- Query Input: A text box where users can enter their medical questions.
- **Response Display**: A section where the generated ansIr is shown to the user.
- **Error Handling**: If the query is invalid or no relevant passages are found, the system provides feedback (e.g., "No relevant passages found").

b. Features:

1. Response Display: The generated ansIrs from the system are displayed along with the retrieved passages, allowing users to understand the context behind the response.

Example Output:

plaintext

Copy code

User Query: What are the symptoms of diabetes?

System AnsIr: The symptoms of diabetes include increased thirst, frequent urination, and fatigue.

Retrieved Passages:

- 1. Diabetes is a chronic condition that affects how your body turns food into energy.
- 2. Symptoms of diabetes include increased thirst, frequent urination, and fatigue.
- **2. Error Handling:** I ensured that the interface provides meaningful messages when an error occurs. For example:
 - **Invalid Input**: "Please enter a valid medical query."
 - No Relevant Passages Found: "No relevant information found for this query."

c. Flask API Implementation:

The system was deployed using **Flask**, enabling interaction between the frontend (UI) and the backend (retrieval and generation). The API routes were designed to handle both the query submission and answer retrieval.

python

Copy code

from flask import Flask, request, jsonify, render_template

```
app = Flask(__name__)
```

@app.route('/')

def index():

return render template('index.html')

@app.route('/ask', methods=['POST'])

```
def ask():
    user_query = request.json.get('query')
    retrieved_passages = retrieve_passages(user_query)
    generated_ansIr = generate_ansIr(user_query, retrieved_passages)
    return jsonify({"ansIr": generated_ansIr})

if __name__ == '__main__':
    app.run(debug=True)
```

d. Scalability and Future Enhancements:

- Latency Reduction: By optimizing the retrieval process and using a smaller model (T5-Small), I ensured that responses were generated within a reasonable time frame (~1-2 seconds).
- Multi-user Scalability: The Flask API can be deployed on cloud platforms like AWS or Heroku to handle multiple concurrent users, with enhancements such as load balancing and horizontal scaling.

Challenges and Solutions:

1. Hardware and Training issues:

One significant challenge I faced during the development was the **hardware limitations** and the **training time** required to fine-tune the model. Initially, the model training took around **3 hours** for each run due to the size of the dataset and the complexity of the **T5 model**. The limited hardware resources (e.g., GPU availability) resulted in frequent **runtime disconnections**, especially when using cloud-based environments like Google Colab, which has session time limits.

To mitigate this issue I tried these:

- Multiple Training Sessions: The model had to be trained multiple times with varying learning rates and numbers of epochs due to the frequent disconnections. For example, the first few training sessions on 300 samples shold no training loss, only validation loss, which led to increasing the dataset size to 1,000 samples for better results.
- Adjusting Learning Rates: Different learning rates were tested in multiple sessions to strike a balance between convergence speed and loss stability.
- **Session Management**: I implemented checkpointing mechanisms to save the model's progress periodically, allowing the training to resume from the last checkpoint after a disconnection.

Despite these limitations, through trial and error and with increased training samples and optimized learning parameters, I was able to achieve a stable model with decent validation performance. But there is still scope for improvement.

2. Handling Medical Terminology:

Medical terminology presented a significant challenge during retrieval and generation. While common terms were handled Ill, more complex terms were prone to misunderstanding. To mitigate this:

- I incorporated external medical dictionaries (UMLS) for improving term recognition.
- Guardrails were implemented to filter out incorrect or unsafe medical advice.

3. Training Challenges:

Initially, training the model on a **small dataset of 300 samples** did not yield proper training loss logging, and only validation loss was recorded. This hindered performance improvements across epochs. After increasing the dataset size to **1,000 samples** and varying the learning rates, training and validation loss shold significant improvements, leading to a better-trained model.

4. Reducing Latency:

Initial attempts with larger models (e.g., T5-Base) resulted in high latency. I switched to **T5-Small**, which significantly reduced the response time while maintaining acceptable performance.

5. Passage Relevance:

Sometimes, irrelevant passages were retrieved, leading to incorrect ansIrs. I improved the passage ranking mechanism by adding a relevance scoring algorithm based on **cosine similarity** and experimented with various FAISS configurations.

Conclusion:

The Intelligent Medical Assistant system successfully integrates **retrieval** and **generation** models to ansIr medical queries. The UI provides an easy-to-use interface for users to input their questions and receive contextually accurate ansIrs in real-time. Evaluation results demonstrate strong performance in generating relevant answers for queries on symptoms, treatments, and other medical topics.

Future Scope of Work:

- 1. **Image and Lab Report Integration**: Expanding the system to handle non-textual data such as medical images and lab reports.
- 2. **Scalability**: Implementing a fully scalable version of the system on cloud platforms like AWS for real-world use.
- 3. **Advanced Error Handling**: Adding more sophisticated error-handling mechanisms to handle complex or ambiguous queries.
- 4. **UI Enhancements**: Adding visual aids, graphs, or diagrams to enhance the user experience for complex medical queries.