

# Optimizing Question-Answering (QA) Performance in Large Language Models(LLMs): A Study in Parameter-Efficient Fine-Tuning (PEFT), Evaluating Retrieval-Augmented Generation (RAG) and Retrieval-Augmented Fine-Tuning (RAFT) Systems for Domain-Specific QAs

## Project Overview

This project evaluates Fine-Tuned, Retrieval-Augmented Generation (RAG) and Retrieval-Augmented Fine-Tuning (RAFT) systems for domain-specific QAs.

## Prerequisites

- Create a Hugging Face account
- Generate an access token with write permissions

## Setup Instructions

### 1. Creation of Q&A datasets

- i. This section includes this notebook:
  - i. QADataset.ipynb
- ii. Run all the cells in the “QADataset.ipynb” file in order to create the Q&A datasets, including HR general, and Advanced medical Questions and response. These files are created in two possible formats such as csv and json.

### 2. Dataset Creation

- i. This section includes this notebook:
  - i. NLPPProject\_Dataset\_MedicalQA.ipynb
- ii. Store the access token in the write mode in the secret keys in order to access it for pushing the dataset on the Hugging Face repository. We use this smaller dataset to fine tune our selected LLM.
- iii. Run the first cell to install all the required libraries in order to run this notebook.
- iv. Login to your Hugging Face by using the same access token stored in the colab notebook.
- v. Load the chosen dataset. You can change it to a preferred one.
- vi. Choose a LLM, the default model is Mistral-7b-v0.2. You can change it in the tokenizer section.
- vii. Run all the cells respectively.

- viii. You can change the token limits in “valid\_indices” in order to reduce the number of instructions and outputs.
- ix. In order to choose the top\_k rows, it's possible to change the k parameter to create a smaller or larger dataset.
- x. Push the created dataset to your Hugging Face repository by running the last cell.

### **3. Performing Domain-Specific Fine-Tuning**

- i. This section includes this notebook:
  - i. Mistral\_7B\_Instruct\_v0\_2\_finetuning\_medicaldb.ipynb
- ii. It follows the same procedure as before by running all the cells from top to down.
- iii. You can change the base\_model variable to use your preferred LLM.
- iv. It is possible to alter the configuration of training setups to fine tune the model in your preferred setups.
- v. Push the fine-tuned model to your Hugging Face by running the last cells.

### **4. Querying Baseline and fine-tuned models**

- i. This section includes these notebooks respectively. All these notebooks follow the same procedure. They are implemented by different configurations and system prompts to improve the generated responses.
  - i. NLPPProject\_Mistral\_7B\_Instrsuct\_v02\_Prompting\_QA(P1).ipynb
  - ii. NLPPProject\_Mistral\_7B\_Instrsuct\_v02\_Prompting\_QA(P2).ipynb
  - iii. NLPPProject\_Mistral\_7B\_Instrsuct\_v02\_Prompting\_QA(P3).ipynb
  - iv. NLPPProject\_Mistral\_7B\_Instrsuct\_v02\_Prompting\_QA(P4).ipynb
  - v. NLPPProject\_Mistral\_7B\_Instrsuct\_v02\_Prompting\_QA(P5).ipynb
- ii. In this section you need to upload the created dataset for HR general and advanced medical Q&A on the colab notebook before running the notebook cells.
- iii. After uploading the mentioned datasets, It follows the same procedure as before by running all the cells from top to down.
- iv. You can also change the load\_in parameter to load the model in 4 bit or 8 bit precision.
- v. It is possible to alter pipeline parameters to generate responses using different setups such as temperature and top\_p.
- vi. Save the generated csv and json files containing questions and generated answers by the model to evaluate them.

### **5. RAG System Implementation**

- i. This section includes these notebooks. The first notebook is implemented by using the base model and the second one uses the fine-tuned version. All these notebooks follow the same procedure, these two models implemented by the same setups and system prompts to generate responses.
  - i. NLPPProject\_Mistral2\_7B\_RAG.ipynb
  - ii. NLPPProject\_Mistral2\_7B\_RAG(FT).ipynb
- ii. Here, you need to upload the papers, from the papers directory in the zip file in this directory on the colab notebook “./content/papers” before running the notebook cells.
- iii. Then you need to upload the created dataset for HR general and advanced medical Q&A on the colab notebook as the same as before.
- iv. You can also change the load\_in parameter to load the model in 4 bit or 8 bit precision.
- v. It is possible to alter pipeline parameters to generate responses using different setups such as temperature and top\_p.
- vi. Save the generated csv and json files containing questions and generated answers by the model to evaluate them.

## 6. RAFT Integration and Evaluation

- i. This section includes this notebook:
  - i. Mistral\_7B\_Instruct\_v0\_2\_finetuning\_medicaldb.ipynb
- ii. In this section you need to upload the papers, from the papers directory in the zip file in a local directory on the colab notebook named “papers” before running the notebook cells.
- iii. Then, you need to upload the created dataset for HR general and advanced medical Q&A on the colab notebook as the same as before.
- iv. Save the generated csv files containing generated answers to evaluate them.

## Contributing and Contacts

Shakiba Farjood Fashalami

- shakiba.farjoodfashalami@studenti.unipd.it

Nima Daryabar

- nima.daryabar@studenti.unipd.it

Tobia Pavona

- tobia.pavona@studenti.unipd.it

Marcos Tidball

- marcos.tidball@studenti.unipd.it

Giacomo Ferrante

- giacomo.ferrante@studenti.unipd.it