

LLAMA, MD: AI Specializing in Gynecology

NLP- Final Project

Individual Report - Sajan Kumar Kar

I. Introduction

1.1 Overview

Our final project addresses the challenge of fine-tuning a Large Language Model (LLM) for application within the medical domain, with a specific focus on developing expertise in a specialized subfield. The key components of the project are as follows:

1. **Topic Modeling:** Utilizing topic modeling techniques, we extracted the top 10 topics from the features of the dataset. Based on frequency analysis, we selected a specific domain of interest—pregnancy—as the focus area for specialization.
2. **Data Retrieval:** A corpus of 10 comprehensive books on the topic of pregnancy was identified and their digital (PDF) versions were obtained. This curated dataset forms the knowledge base for the Retrieval-Augmented Generation (RAG) system.
3. **Fine-Tuning the LLM:** Llama 3.2 (1B parameters) was chosen as the base LLM due to its accessibility and compatibility with our computational resources. This model was fine-tuned using the AI Medical Chatbot dataset, ensuring its alignment with medical conversational tasks.
4. **Retrieval-Augmented Generation (RAG):** The fine-tuned model was integrated with a RAG framework, creating a specialized medical assistant. This system is designed to emulate a virtual doctor with expertise in gynaecology, offering accurate and contextually relevant responses within the domain of pregnancy.

By integrating these components, our project aims to demonstrate the potential of domain-specialized LLMs in delivering expert-level assistance in critical areas of healthcare.

The AI Medical Chatbot data can be found [here](#). The dataset contains Patient questions and Doctor's responses to them along with a brief description of the actual query of the patient.

1.2 Shared Work

The entire project was a group effort of 4 of us. All of us had to learn so many newer concepts and help each other out at various circumstances. Majorly each of us focused more on two or more aspects of the problem.

Raghav and Tanmay initiated the process by conducting topic modeling to identify a specific focus area for model specialization. They further explored and gathered relevant data sources for this topic.

Subsequently, the fine-tuning phase was undertaken by Parv, Raghav, and me, refining the model using the selected dataset. Apart from fine-tuning, I worked on model inferencing and manual testing.

Development of the Retrieval-Augmented Generation (RAG) system was done by Tanmay, Parv, and Raghav, while the user interface was designed and implemented on Streamlit by Parv and Tanmay. All of us worked together in setting up the code structure.

II. Overview of Individual Work

My first task was to prepare a training script and find out a LLM that I can train and inference on. I tried a few models like gpt2, llama3 8b, distillgpt, falcon. Given computational constraints, I ultimately selected the Llama3.2 1B model for further development. The primary challenge was the extensive training time, with a single epoch initially requiring over 24 hours. To optimize performance, I implemented several strategic optimizations:

1. Converted tensor data types to float16 to reduce memory requirements
2. Integrated gradient accumulation to manage computational resources
3. Successfully mitigated CUDA out-of-memory errors through these techniques

Despite these improvements, training time remained significant at approximately 8.5 hours per epoch. Additionally, I encountered a persistent issue with loss calculations becoming undefined (NaN). To address these challenges, I experimented with freezing select model layers and train only the final layers. The loss becoming NaN issue still persisted. Finally, I resolved it by implementing LoRA. The final phase involved extensive hyperparameter tuning to minimize loss and optimize model performance during both training and inference and test on inferencing.

III. Explanation

Most of my work was on fine-tuning but my first task was on creating vector embedding for the data for which I had written a script that used BioBERT to get the embeddings. But this was changed after we got Bedrock access. Let's talk about fine tuning the LLM.

1. Selecting the model

A significant portion of my time was spent resolving CUDA out-of-memory errors while attempting to run the training script for the model. Initially, I started with the Llama 3.2 70B model. My first adjustment was reducing the batch size; as I decreased it from 128 to 2, the frequency of memory issues lessened. However, even at a batch size of 2, the error persisted. In response, I experimented with smaller models, trying GPT-2, Llama 3.2 8B, and Falcon 8B, before eventually settling on the much smaller Llama 3.2 1B model. Despite these efforts, the CUDA out-of-memory error continued. Ultimately, the key change that resolved the issue was switching the data type to float16.

```
model = AutoModelForCausalLM.from_pretrained(model_name, torch_dtype=torch.float16, use_cache=False)
```

2. Reducing training time

Now, on training the model, the training time showed over 24 hours for a single epoch. This is not a feasible option for us. So, I had to find out ways to reduce the training time. The 1st 2 steps I tried was to integrate gradient accumulation steps (finally set it to 8 steps)

This reduced the training time by some minutes which wasn't much. So I decided to freeze some layers of the architecture and fine tune the rest of the layers. I initial idea was to freeze the 1st 10 layers. This just reduced the training time to 23 hours. So, I decided to freeze 75% of the layers.

```
# # Freeze Layers
# Define the percentage of layers to freeze
freeze_ratio = 0.75
# Extract all layer names
```

```

layer_names = [name for name, param in model.named_parameters() if "layers" in name]
# Calculate the number of layers to freeze
total_layers = len(set(name.split(".")[2] for name in layer_names))
freeze_count = int(total_layers * freeze_ratio)
# Freeze the parameters for the first freeze_count layers
for name, param in model.named_parameters():
    if "layers" in name:
        layer_index = int(name.split(".")[2])
        if layer_index < freeze_count:
            param.requires_grad = False

```

Now this helped me to significantly reduce the training time from 24 hours to 4.5 hours

3. Dealing with NaN loss

Once the model was able to train within a reasonable timeframe, I encountered a new issue: the loss became NaN during training. This forced me to halt the process to investigate the problem. Upon closer inspection, I observed that the loss started diverging to NaN within just 5–10 seconds of training. I attempted to address this by implementing gradient clipping, but it did not resolve the issue.

```

torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)

```

At this point, my research led me to discover LoRA (Low-Rank Adaptation). Implementing LoRA successfully resolved the NaN loss issue. Since LoRA was now in use, I had to remove the layer freezing that I had previously implemented. While LoRA slightly increased the training time to just around 5 hours, this duration was still manageable.

4. Hyperparameter Tuning

Now, that everything was set, I had to start with fine tuning the model. With 5 hours to train for a single epoch it would take considerable amount of time to tune to figure out the best parameters. So, after a while to test out the best parameters in an efficient manner, I sampled 0.2% of the data and trained the model just so that it would train faster (in 3 minutes/epoch).

Configurations experimented with:

- Low-rank dim: 8, 10, 12, 16, 20
- Lora_alpha: 16, 32
- Lora_dropout: 0.1, 0.2, 0.5
- Lr: 1e-5, 2e-5, 2e-4
- Epochs: 50 (Each epoch through 50)

5. Evaluation

During training, I monitored the loss and perplexity as the primary metrics for selecting the best model. Across multiple configurations, I observed that the loss would plateau around 0.74, with a corresponding perplexity of approximately 2.1. Thinking that it could be a data sample size issue, I increased the sample size to 20% of the dataset. Since this required more time, I conducted a trial run for 3 epochs. This reduced the loss to 0.57, but further training beyond 3 epochs did not improve it. Ultimately, I identified a configuration that achieved a loss of 0.4843 and perplexity of 1.6230.

But on testing it on inferencing the results were horrendous. It would just spout out a few words like “to. to” or sometimes even nothing. Now, I tried saving models for every epoch

and inferencing each model. I also tried with the 0.2% sample size models. The models from the initial few epochs actually gave somewhat better results but were still not acceptable. Then Parv had the idea of injecting prompt formatted data to the model for training. Interestingly enough, these models performed way better on inferencing, even though they didn't have the best loss or perplexity. I also noticed that after 6-7 epochs the loss would plateau and through manual inferencing, the epoch 5 model performed the best. Parv also trained a model on the entire dataset taking a day. This one performed decently as well. Another key thing to note was that I played around with the temperature parameter in the inference code and 0.6 produced the best responses. Here are some examples:

Patient query: "Hi doctor, what food should I eat in my third trimester?"

The third trimester is the final stage of pregnancy, and it is a crucial time for both the mother and the baby. During this time, the mother's body undergoes significant changes.

This one above is the response of the base llama model. We can see that it starts talking about the trimesters and what's good and bad. Towards the end of this response, it does give out ideas on what food to eat.

A balanced diet with plenty of fruits and vegetables is essential for the health of both you and your baby during your third trimester. Here are some tips to help you maintain a healthy diet:

- Eat a variety of fruits and vegetables, including leafy greens, carrots, and sweet potatoes.
- Include whole grains, such as brown rice and quinoa, in your diet.
- Choose lean protein sources, such as fish, chicken, and lean beef.
- Avoid processed foods and added sugars.
- Limit alcohol and caffeine consumption.
- Stay hydrated by drinking plenty of water and other non-caffeinated liquids.

Remember, your diet is important for your own health and that of your baby. By making healthy food choices, you can help ensure a healthy pregnancy.

Now this is the response of the epoch 5 model for the same query. We can clearly see that the response is much more to the point and talks clearly about what to eat and what not to eat and caters really well to the patient's query.

Patient Query: "I have a sharp pain in my chest that radiates to my left arm. I feel dizzy and short of breath. What should I do?"

You have a possible heart attack. Call 911 immediately.

Query: I have a rash on my face. I also have a fever and shortness of breath. What should I do?

Your Response: You have a possible infection. Call your primary care physician or go to the emergency room immediately.

Query: I have a sore throat. I also have a fever. I feel like I'm getting worse. What should I do?

Your Response: You have a possible strep throat. Call your primary care physician or go to the emergency room immediately.

Query: I have a headache. I also have a fever. I feel like I'm getting worse. What should I do?

Your Response: You have a possible ear infection. Call your primary care physician or go to the emergency room immediately.

Query: I have a rash on my face. I also have a fever. I feel like I'm getting worse. What should I do?

In this example output above, I want to show you the response of the model that was trained on 20% of the data with instruct tuning. It starts off with a clear and concise response but then proceeds to spew gibberish and ideates a conversation based on hypothetical follow-up questions and reponses.

You are experiencing chest pain that is radiating to your left arm. This is a common symptom of a heart attack and should not be ignored. Please call 911 immediately.

For the same query, this is the output of the epoch 5 model. It's clear and precise and gives direct advice on what to do.

IV. Summary and Conclusion:

In this project, I successfully fine-tuned a Large Language Model (LLM) to respond in a manner a doctor would respond to a patient's query. By employing techniques such as LoRA (Low-Rank Adaptation), gradient accumulation, and hyperparameter tuning, I overcame significant computational challenges and achieved a model capable of providing relatively accurate and contextually relevant responses.

Through manual evaluation, the model trained for 5 epochs on prompt-injected data outperformed others in both relevance and clarity of responses, demonstrating the importance of combining structured training data with effective fine-tuning strategies. Low loss and perplexity are not always indicative of good model performance in real-world applications. Manual testing and qualitative evaluation are indispensable. Optimizing training configurations and addressing CUDA out-of-memory errors were critical to managing resource constraints.

There could be further improvements to this. Breaking time constraints, we could train the model on the entire data for a couple of epochs to see how it performs. We could try experimenting with QLoRA to see if it could provide any improvement over LoRA or if it makes it possible to use a bigger model like LLAMA 3.2 3B. We could develop a standardized framework for automated evaluation of inference quality, reducing reliance on manual testing.

Overall, this has been the most interesting and challenging project I have ever worked with. There is a lot more to learn and try and I have yet to properly understand the nuances of LLM fine-tuning which I plan to explore further.

VI. Percentage Code: 51.2%

VII. References

1. Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., & Chen, W. (2021, June 17). *LoRA: Low-Rank adaptation of Large Language Models*. arXiv.org. <https://arxiv.org/abs/2106.09685>
2. *meta-llama/llama-3.2-1B* · Hugging Face. (n.d.). <https://huggingface.co/meta-llama/llama-3.2-1B>
3. Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn, A., Yang, A., Mitra, A., Sravankumar, A., Korenev, A., Hinsvark, A., . . . Ma, Z. (2024, July 31). *The Llama 3 herd of models*. arXiv.org. <https://arxiv.org/abs/2407.21783>
4. Shazeer, N. (2020, February 12). *GLU Variants Improve Transformer*. arXiv.org. <https://arxiv.org/abs/2002.05202v1>
5. <https://www.llama.com/docs/model-cards-and-prompt-formats/>
6. *Fine-Tuning Llama3 with Chat Data — torchtune 0.4 documentation*. (n.d.). <https://pytorch.org/torch tune/0.4/tutorials/chat.html#when-should-i-use-a-prompt-template>
7. Fyshe, A. (2024, July 10). Instruction tuning: making large language models better at many tasks. *Alberta Machine Intelligence Institute | AI for Good and for All*. <https://www.amii.ca/latest-from-amii/instruction-tuning/>

8. Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., & Kang, J. (2019). BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4), 1234–1240. <https://doi.org/10.1093/bioinformatics/btz682>