

## PROGRESS PROJECT REPORT

### Members

Member 1: Tran Khanh Bang (V202100411)  
Member 2: Truong Gia Bao (V202100564)

### Date

August 21, 2024  
Advisor: Prof. Wray Buntine

## 1 Introduction

This report serves to document the ongoing progress and developments of the Chatbot project designed specifically for the Computer Science (CS) students at VinUniversity. The primary aim of this report is to provide a comprehensive update to stakeholders on the current status of the project, aligning with the final goals and demonstrating adherence to the outlined academic and practical objectives.

## 2 Challenges

Our project faces these major challenges:

- Leveraging OpenAI's API or undertaking fine-tuning on our infrastructure necessitates significant computational resources, which are currently beyond our project's budget. Additionally, the time required for fine-tuning these models is considerable, potentially impacting project timelines.
- Limited server resources, including insufficient disk space, hindered the installation of necessary packages and the deployment of the final model.
- While the CECS server provided GPU acceleration for model fine-tuning, limitations in deploying the large model to GitHub due to size constraints and LFS package issues posed additional obstacles (only the paid version can use more than 1GB)

```

(base) ubuntu@445ee971f5a2:~/21ba0$ cd 21ba0.tg@vinuni.edu.vn
(base) ubuntu@445ee971f5a2:~/21ba0.tg$ cd 21ba0.tg
(base) ubuntu@445ee971f5a2:~/21ba0.tg$ cd COMP3080-Course-Related-Project
(base) ubuntu@445ee971f5a2:~/21ba0.tg/COMP3080-Course-Related-Project$ git push
Username for 'https://github.com': bao-tg
Password for 'https://bao-tg@github.com':
Username for 'https://github.com': bao-tg
Password for 'https://bao-tg@github.com':
Username for 'https://bao-tg@github.com':
Password for 'https://bao-tg@github.com':
Uploading LFS objects: 100% (2/2), 1.5 GB | 8 B/s, done.
Enumerating objects: 32, done.
Counting objects: 100% (32/32), done.
Delta compression using up to 64 threads
Compressing objects: 100% (43/43), done.
Writing objects: 100% (48/48), 1.25 GiB | 8.88 MiB/s, done.
Total 48 (delta 17), reused 1 (delta 0)
remote: Resolving deltas: 100% (17/17), completed with 2 local objects.
remote: error: Trace: 9b29ae5921e566dd4d2d09effdd0719952b141f25a0897daf6dbdc49c792e85
remote: error: See https://gh.io/lfs for more information.
remote: error: File finetuning/results/checkpoint-3/model.safetensors is 4794.71 MB; this exceeds GitHub's file size limit
of 100.00 MB
remote: error: File finetuning/results/checkpoint-8/optimizer.pt is 949.52 MB; this exceeds GitHub's file size limit of
100.00 MB
remote: error: GH001: Large files detected. You may want to try Git Large File Storage - https://git-lfs.github.com.
To https://github.com/bao-tg/COMP3080-Course-Related-Project
! [remote rejected] master -> master (pre-receive hook declined)
error: failed to push some refs to 'https://github.com/bao-tg/COMP3080-Course-Related-Project'
(base) ubuntu@445ee971f5a2:~/21ba0.tg/COMP3080-Course-Related-Project$

```

Figure 1: LFS package issues of GitHub

- The inability to run the user interface locally on the CECS server using Streamlit impacted the development and testing process.

```

(base) ubuntu@445ee971f5a2:~/21ba0.tg/COMP3080-Course-Related-Project/finetuning$ streamlit run chatbot_finetune.py
Collecting usage statistics. To deactivate, set browser.gatherUsageStats to false.

You can now view your Streamlit app in your browser.

Local URL: http://localhost:8501
Network URL: http://172.17.0.12:8501
External URL: http://118.70.129.147:8501

```

Figure 2: Streamlit run locally

### 3 Current implementation

Our implementation strategy prioritizes cost-effectiveness and accessibility. To achieve this, we utilize a combination of open-source resources, the implementation in details can be found in the footnote<sup>1</sup>:

- **Fine-Tuning Dataset:** We leverage the "ArtifactAI/arxiv-beir-cs-ml-generated-queries" dataset available on Hugging Face (huggingface.co). This pre-existing dataset, specifically focused on computer science and machine learning, provides a strong foundation for fine-tuning the chosen language model.

```
[{"eval_loss": 5.296281814575195, "eval_runtime": 1.2021, "eval_samples_per_second": 166.375, "eval_steps_per_second": 4.159, "epoch": 0.94}, {"train_runtime": 22.0772, "train_samples_per_second": 36.236, "train_steps_per_second": 0.362, "train_loss": 8.944538116455978, "epoch": 0.94}
done
(base) ubuntu@f5ee971f5e2:~/21bao.tg/COMP3080-Course-Related-Project/finetuning$
```

Figure 3: Finetuning 1000 data, 8 batch, 1 epoch(toy example)

- **Language Model and Tokenizer:** We employ the widely used GPT-2 model ("model = GPT2LMHeadModel.from\_pretrained('gpt2'), tokenizer = GPT2Tokenizer.from\_pretrained('gpt2')") from Hugging Face. This combination offers a good balance between performance and resource efficiency, making it suitable for deployment with Streamlit.
- **High-Performance Fine-Tuning:** To expedite the fine-tuning process, we leverage the CECS server's robust infrastructure, including 6 NVIDIA A5000 GPUs, 64 CPUs, and 230GB RAM. This computational power enables efficient training of the language model on our dataset, accelerating development timelines.
- **User Interface Development:** The user interface of the chatbot is developed using **Streamlit**, we also found the way to deploy the model as an application in the Streamlit repository. This choice was made to ensure that the interface is not only intuitive and user-friendly but also capable of supporting the dynamic requirements of interactive chatbot functionalities.

### 4 Teamwork

- **Truong Gia Bao (Project Manager - Software Engineer):**
  - Led the project planning and coordination efforts.
  - Oversaw the overall project timeline and ensured timely delivery.
  - Developed and maintained core system functionalities.
- **Tran Khanh Bang (Software Engineer):**
  - Contributed to project development by implementing features and fixing bugs.
  - Collaborated with the team to ensure successful app delivery.
  - Demonstrated technical expertise in connecting the UI components.

### 5 Goals in the final phase

To further optimize the chatbot's capabilities and user experience, we plan to implement several enhancements:

- **Dataset Expansion:** We will significantly expand the training dataset by incorporating a wider range of CS-related materials, including textbooks, research papers, and code repositories. This will enrich the chatbot's knowledge base and improve its ability to provide comprehensive and informative responses.

<sup>1</sup>See [here](#).

- **Model Upgrading:** The current model size may limit the chatbot’s capacity for complex reasoning and context understanding. We aim to explore larger language models to enhance these capabilities, potentially leading to more sophisticated and nuanced interactions.
- **Resolution of the LFS Package Issues on GitHub:** In the final phase, the objective is to address and resolve the size constraint limitations of the LFS package.

## 6 Conclusion

This report outlines the development of a CS student chatbot, addressing challenges related to resource constraints and model performance. By leveraging open-source tools, high-performance computing, and a user-centric design, we have created a foundation for a valuable educational resource. Future enhancements will focus on expanding the dataset, upgrading the language model, and implementing user data management to further improve the chatbot’s capabilities.

## 7 Appendices

1. Hugging Face. (n.d.). ArtifactAI/arxiv-beir-cs-ml-generated-queries [Dataset]. Retrieved from <https://huggingface.co/datasets/ArtifactAI/arxiv-beir-cs-ml-generated-queries/tree/main>
2. Hugging Face. (n.d.). transformers: Pre-trained Deep Learning Models for NLP [Software]. Retrieved from <https://huggingface.co/transformers>
3. Streamlit. (n.d.). *Streamlit Official Website*. Retrieved from <https://streamlit.io/>