

**VIETNAM NATIONAL UNIVERSITY, HANOI
UNIVERSITY OF ENGINEERING AND TECHNOLOGY**



NGO DUC HUY

TITLE

SUMMER INTERNSHIP REPORT

Major: Computer Science

Supervisor: Dr. Ta Viet Cuong

HANOI - 2024

**VIETNAM NATIONAL UNIVERSITY, HANOI
UNIVERSITY OF ENGINEERING AND TECHNOLOGY**

INTERNSHIP REPORT

TITLE

Student: Ngo Duc Huy
Student ID: 21020046
Class: QH-2021-I/CQ-I-CS2
Supervisor: Dr. Ta Viet Cuong

HANOI - 2024

Acknowledgements

I would like to express my sincere gratitude to my supervisor, Dr. Ta Viet Cuong, for his guidance, support, and encouragement throughout the internship. His valuable advices and feedbacks have helped me a lot in completing the research and experiments.

I would also like to the Faculty of Information Technology, University of Engineering and Technology, VNU, and HMI lab for providing me with the opportunity to participate in the summer internship program.

Without the support, I would not have been able to complete the internship and this report.

Abstract

Abstract: This is the report on the research and experiments on Large Language Models (LLMs) during the summer internship at the HMI Lab, Faculty of Information Technology, University of Engineering and Technology, VNU, under the guidance of Dr. Ta Viet Cuong.

Keywords: *Keyword*

Table of Contents

Acknowledgements	iii
Abstract	iv
Table of Contents	v
List of Figures	vi
1 Introduction	1
1.1 Large Language Models and Tiny Language Models	1
1.2 Research Objectives	2
2 Related Works	3
2.1 PhoGPT	3
2.2 Tiny Language Models	3
3 Research methodology	5
3.1 Model architecture modification	5
3.2 Metrics and performance	6
4 Experiments and Results	7
Conclusion	9
References	10

List of Figures

Chapter 1

Introduction

1.1 Large Language Models and Tiny Language Models

Large Language Models (LLMs) are machine learning models that are capable of processing and understanding natural language through learning from a large amount of text data. They can predict, generate text, answer questions, translate, and perform many other complex language tasks. Thanks to the computational power and large-scale data, these models have made significant breakthroughs in Natural Language Processing (NLP).

LLMs can come in various sizes, from a few million to hundreds of billions of parameters. The larger the model, the better the performance, but also the higher the computational cost. For example, GPT-3, one of the largest LLMs, has 175 billion parameters, which requires a large amount of memory and computational resources to train and deploy.

Therefore, one of the challenges when using LLMs is the high computational cost and resource requirements. To address this issue, Tiny Language Models are introduced, which are smaller versions of LLMs that could be acquired by creating from scratch or through processes such as knowledge distillation, model finetuning, and other training techniques. These models still retain the core language capabilities like other larger model but with lower resource requirements, faster query processing speed, at a small exchanged cost of performance.

The research and use of tiny language models is a new and relatively important research direction in the field of NLP, helping to optimize the performance and cost

of NLP applications in practice, or in training step, it may reduce the amount of data or energy required, which is beneficial for the environment. This could be the way for small organizations or individuals lacking resources to enter the field of language models research and development.

1.2 Research Objectives

This report presents the research and experiments on a tiny language model based on PhoGPT, a state of the art language models for Vietnamese. The process including explore the process of modifying the model, evaluating the performance of the model in specific evaluations. In overall, the report aims to answer the following questions:

- Can we modify a LLM like PhoGPT to create a smaller version with fewer parameters and lower resource requirements?
- What is the tradeoff between the performance and the resource requirements of the SLM compared to the original LLM?

Chapter 2

Related Works

2.1 PhoGPT

PhoGPT [1] is a open-source state-of-the-art Transformer decoder-based model series for Vietnamese, including the base pre-trained monolingual model PhoGPT-4B and its chat variant, PhoGPT-4B-Chat. The base model, PhoGPT-4B, with exactly 3.7B parameters, is pre-trained from scratch on a Vietnamese corpus of 102B tokens, with an 8192 context length, employing a vocabulary of 20480 token types. The chat variant, PhoGPT-4B-Chat, is the modeling output obtained by fine-tuning PhoGPT-4B on a dataset of 70K instructional prompts and their responses, along with an additional 290K conversations.

The PhoGPT series has achieved state-of-the-art results on various Vietnamese NLP tasks.

<TABLE here>

2.2 Tiny Language Models

The work of Yehui Tang et al. [2] introduces a series of experiments on Tiny Language Models (TLMs) based on Pangu, a large-scale language model on Chinese and English. The authors propose a method to compress Pangu into Pangu- π -1B-pro and Pangu- π -1.5B-pro with 1B and 1.5B parameters, respectively. The experiments show that the TLMs can achieve competitive performance with a much smaller number of parameters compared to the original model.

The authors propose many methods for creating the models, from the model architecture, parameter initialization, to the optimized training strategies. Firstly, for the model architecture, there are two main components: the tokenizer and the model itself:

- As for the tokenizer, the authors after conducting experiments with different vocabulary sizes, they came to the conclusion that over 50% vocabularies may be redundant as they cater to less than 3% of the corpus. Therefore, they propose to choose the tokenizer that cover 90% of the corpus vocabulary, which is the reduced one from Pangu- π model, resulting in a vocabulary size of 48K instead of 100K tokens.
- And for the model architecture tweaking, the authors explored the impact of depth, width and the expanding rate of Feed-Forward Networks (FFN) on the performance of a 1B-size language model. They found that the depth of the model has a significant impact on the performance and inference speed, while the width and the expanding rate of FFN have a minor impact.

Secondly, when working with parameter initialization methods like random initialization and parameter inheritance, the authors found that the latter method can help the model perform better by preserving the knowledge of the pre-trained model and reducing the training time.

By using both important inter-layer selection and intra-layer parameter selection strategies, the authors can choose which layers and neurons to inherit to the model, keeping the most important elements that contribute to the larger model’s performance while vastly reducing the number of parameters.

Finally, the authors propose a series of optimized training strategies after experimenting with batch size and learning rate, in combination with multiple-round training. This results in less hardware consumption while rectifying the problem of data forgetting due to the limitation in capacity of the model.

Chapter 3

Research methodology

In this chapter, the methods used to conduct the research and answer to the 2 research questions mentioned previously will be presented. Firstly, the method for modifying the model architecture will be discussed. And then the performance metrics and the results of the modified model will be presented.

3.1 Model architecture modification

The model architecture modification is done by changing the model's architecture to a smaller one. One of the most common ways to reduce the size of the model is to reduce the number of layers in the model. The layers in the concept of a generative Transformer LLM is referred to as the number of Transformer blocks placed consecutively in the model. The Transformer block is the basic building block of the Transformer model, which consists of a multi-head self-attention mechanism and a feed-forward neural network. The number of Transformer blocks in a model is the most important factor that determines the model's size. Therefore, reducing the number of Transformer blocks in the model will reduce the model's size significantly.

Since the blocks are placed consecutively in the model, i.e., the output of the previous block is the input of the next block, the process happens layer by layer automatically. Therefore layer reduction can be done by simply removing the layers from the model.

3.2 Metrics and performance

The performance of the modified model is evaluated using the perplexity metric. Perplexity is a common metric used to evaluate the performance of a language model. It is calculated as the exponential of the cross-entropy loss of the model on the test dataset. The lower the perplexity, the better the model's performance.

The perplexity is calculated as follows:

$$\text{Perplexity} = \exp \left(\frac{1}{N} \sum_{i=1}^N -\log P(x_i) \right)$$

where N is the number of tokens in the test, and $P(x_i)$ is the probability of the i -th token in the test dataset.

The performance of the model is evaluated by measuring the perplexity of the model on the test dataset consisting of input prompts sequences. The process is done by feeding the input sequences to the model and calculating the perplexity of the model on the output sequences. Then the perplexity is averaged over the whole test dataset to get the final perplexity score of the model.

Chapter 4

Experiments and Results

The experiments are conducted first by using the original PhoGPT-4B-Chat model as the baseline model. With 24 Transformer blocks or layers, each with 24 attention heads, the model results in roughly 3.7B parameters.

The model is then modified in multiple ways to reduce the model size by removing different numbers and positions of layers from the model. We had experimented with removing 1, 2, 4, 8 layers consecutively, each time removing the layers beginning from the first layer to the last layer possible of the model. The modified model are then evaluated using the perplexity metric on the test prompt dataset. The average perplexity is calculated as metrics to evaluate the performance of the model.

The experiments are conducted on a single NVIDIA RTX 3090Ti GPU with 24GB of VRAM, given the following results:

<TABLE here>

As can be seen from the table, the baseline model PhoGPT-4B-Chat has a perplexity of ... on the test dataset. When removing 1 layer from the model, the perplexity increases to ... as average, to ... when removing the critical layers The results pattern continues as more layers are removed from the model. The perplexity increases as the number of layers removed increases, the first and last layers have much more significant impact on the model's performance compared to the middle layers.

The runtime or inference speed of the model is also measured during the experiments. The runtime of the model is measured by calculating the average time taken to generate the output sequences for the test dataset. As the model size decreases, the run-

time of the model also decreases. The runtime of the model is inversely proportional to the model size, as the trade-off between model size or performance and inference speed is a common problem in the battlefield of LLMs.

Conclusion

In this internship report, we have presented the work on the modification of the PhoGPT-4B-Chat model to reduce the model size while not significantly affecting the model's performance. The experiments are conducted by removing different numbers and positions of layers from the model and evaluating the model's performance using the perplexity metric on the test prompt dataset. The results show that the trade-off between model size and performance, as well as the inference speed. The results also show one important feature of the model: the first and last layers have much more significant impact on the model's performance compared to the middle layers.

References

- [1] D. Q. Nguyen, L. T. Nguyen, C. Tran, D. N. Nguyen, D. Phung, and H. Bui, “Phogpt: Generative pre-training for vietnamese,” 2024. [Online]. Available: <https://arxiv.org/abs/2311.02945>
- [2] Y. Tang, F. Liu, Y. Ni, Y. Tian, Z. Bai, Y.-Q. Hu, S. Liu, S. Jui, K. Han, and Y. Wang, “Rethinking optimization and architecture for tiny language models,” 2024. [Online]. Available: <https://arxiv.org/abs/2402.02791>



Supervisor's comment:

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

Mark: _____

In words: _____

Hanoi, ____/____/2024

Supervisor

(Signature and full name)