

PHOGPT: GENERATIVE PRE-TRAINING FOR VIETNAMESE

Dat Quoc Nguyen, Linh The Nguyen, Chi Tran, Dung Ngoc Nguyen, Dinh Phung, Hung Bui
VinAI Research, Hanoi, Vietnam

{v.datnq9, v.linhnt140, v.chitb, v.dungnn28, v.dinhpq2, v.hungbh1}@vinai.io

ABSTRACT

We open-source a state-of-the-art 4B-parameter generative model series for Vietnamese, which includes 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. In addition, we also demonstrate its superior performance compared to previous open-source models. Our PhoGPT models are available at: <https://github.com/VinAIRResearch/PhoGPT>

1 INTRODUCTION

Undoubtedly, the success of large language models (LLMs), particularly decoder-only transformer-based generative models such as ChatGPT/GPT-4, LLaMA/LLaMA-2 (Touvron et al., 2023a;b), Mistral (Jiang et al., 2023) and Falcon (Almazrouei et al., 2023), stands out as one of the most significant achievements in recent AI research and development. However, that success has largely been limited to English.

For Vietnamese, we now release the 4B-parameter base pre-trained monolingual model, PhoGPT-4B, along with its chat variant, PhoGPT-4B-Chat, as open-source. We pre-train the base model PhoGPT-4B from scratch on a Vietnamese corpus of 102B tokens for two epochs, with an 8192 context length. Here, it employs a Vietnamese-specific byte-level BPE tokenizer with a vocabulary of 20480 tokens. We further fine-tune the base model on a dataset of 70K instructional prompts and their responses, along with an additional 290K conversations, resulting in the chat variant, PhoGPT-4B-Chat. We demonstrate its strong performance compared to previous closed-source and open-source 7B-parameter models.

Our goal is to provide comprehensive and powerful LLMs for Vietnamese, facilitating future research and applications in generative Vietnamese NLP. Our PhoGPT can be used with popular libraries such as “transformers” (Wolf et al., 2020), “vllm” (Kwon et al., 2023) and “llama.cpp”.¹

2 PHOGPT

2.1 PHOGPT-4B: MODEL ARCHITECTURE AND PRE-TRAINING

PhoGPT-4B is a Transformer decoder-based model (Brown et al., 2020; Vaswani et al., 2017), which incorporates (Triton) flash attention (Dao et al., 2022) and ALiBi (Press et al., 2022) for context length extrapolation. We train a Vietnamese-specific byte-level BPE tokenizer with a vocabulary of 20480 tokens using the “tokenizers” library.² In addition, we use a “max_seq_len” of 8192, “d_model” of 3072, “n_heads” of 24 and “n_layers” of 32, resulting in a model size of 3.7B param-

¹<https://github.com/ggerganov/llama.cpp>

²<https://github.com/huggingface/tokenizers>

eters (~4B). Utilizing the Mosaicml “llm-foundry” library (Team, 2023),³ we pre-train PhoGPT-4B from scratch on a 482GB deduplicated and cleaned pre-training corpus of Vietnamese texts (~102B tokens) for two epochs. Our pre-training Vietnamese corpus consists of:⁴

- 1GB of Wikipedia texts (version 20/05/2023);
- 1.5GB of medical-related texts crawled from a wide range of publicly available and medical domain-specific websites such as medical journals and universities;⁵
- 3GB of publicly available books spanning a range of genres;
- 12GB of legal data crawled from `thuvienphapluat.vn` and `lawnet.vn`;
- a 40GB variant of the “binhvq” news corpus (version 21/05/2021);⁶
- an 88GB variant of the Vietnamese OSCAR-2301 subset;⁷
- a 336GB variant of the Vietnamese mC4 subset.⁸

2.2 PHOGPT-4B-CHAT: SUPERVISED FINE-TUNING

We then fine-tune the base pre-trained PhoGPT-4B using a dataset consisting of 70K instructional prompts and their responses, along with an additional 290K conversations, constructed by concatenating the following sources:

- 500 instructional prompt and response pairs for poem writing, 500 for essay writing, 500 for spelling correction, 500 for single-document summarization and 1000 for context-based question answering;
- 67K instructional prompt and response pairs from the Vietnamese subset of Bactrian-X (Li et al., 2023);
- 20K Vietnamese-translated ChatAlpaca conversations;⁹
- 40K Vietnamese-translated ShareGPT conversations (without code and mathematics);¹⁰
- 230K Vietnamese-translated UltraChat conversations (Ding et al., 2023);¹¹

The resulting fine-tuned model is named PhoGPT-4B-Chat.

3 EVALUATION

We compare PhoGPT-4B-Chat with the closed-source models GPT-4-0125-preview, GPT-3.5-turbo and Gemini Pro 1.0, as well as other open-source models, including:

- Vistral-7B-Chat is the modeling output obtained by continually pre-training Mistral-7B (Jiang et al., 2023) on a diverse corpus of Vietnamese texts and then performing supervised fine-tuning using a diverse instructional and conversational dataset.¹²
- SeaLLM-7B-v2 is the modeling output obtained by continually pre-training Mistral-7B on a multilingual corpus from Southeast Asian (SEA) languages, including Vietnamese, and then performing supervised fine-tuning using instructional question and answer pairs.¹³

³<https://github.com/mosaicml/llm-foundry>

⁴Last crawling/cutoff date: 31/05/2023.

⁵The content extracted from these sources contains no private data about the patients.

⁶<https://github.com/binhvq/news-corpus>

⁷<https://huggingface.co/datasets/oscar-corpus/OSCAR-2301>

⁸https://huggingface.co/datasets/allenai/c4/tree/mC4_3.1.0

⁹<https://github.com/cascip/ChatAlpaca>

¹⁰https://huggingface.co/datasets/anon8231489123/ShareGPT_Vicuna_unfiltered

¹¹https://huggingface.co/datasets/HuggingFaceH4/ultrachat_200k (including both train & test sets)

¹²<https://huggingface.co/Viet-Mistral/Vistral-7B-Chat>

¹³<https://huggingface.co/SeaLLMs/SeaLLM-7B-v2>

Model	All truthful questions	Vietnam-specific
PhoGPT-4B-Chat	41.7 (83 / 199)	43.5 (64 / 147)
GPT-4-0125-preview	44.7 (89 / 199)	39.5 (58 / 147)
GPT-3.5-turbo	29.1 (58 / 199)	22.4 (33 / 147)
Gemini Pro 1.0	39.7 (79 / 199)	34.7 (51 / 147)
Vistral-7B-Chat	41.2 (82 / 199)	42.9 (63 / 147)
Sailor-7B-Chat	28.6 (57 / 199)	27.9 (41 / 147)
Sailor-4B-Chat	15.6 (31 / 199)	14.3 (21 / 147)
SeaLLM-7B-v2	20.6 (41 / 199)	13.6 (20 / 147)
VBD-Llama2-7B-50B-Chat	15.6 (31 / 199)	10.9 (16 / 147)
Vinallama-7B-Chat	11.1 (22 / 199)	8.2 (12 / 147)
Gemma-7B-it	8.0 (16 / 199)	6.1 (9 / 147)

Table 1: Obtained results.

- Sailor-7B-Chat¹⁴ and Sailor-4B-Chat¹⁵ are the modeling outputs obtained by continually pre-training Qwen1.5-7B and Qwen1.5-4B (Bai et al., 2023) on 400B tokens from SEA languages, including Vietnamese, and then performing supervised fine-tuning using instructional question and answer pairs.
- VBD-LLaMA2-7B-50b-Chat is the modeling output obtained by continually pre-training LLaMA-2-7B (Touvron et al., 2023b) on a data combination of 40B Vietnamese tokens and 16B English tokens and then performing supervised fine-tuning using 2M instructional and conversational samples.¹⁶
- VinaLLaMA-7B-Chat is the modeling output obtained by continually pre-training LLaMA-2-7B on a data combination of 100B English tokens, 230B Vietnamese tokens (from books and news), and 500B automatically generated Vietnamese tokens, and then performing supervised fine-tuning using 1M instructional and conversational samples.¹⁷
- Gemma-7B-it is the 7B instruct version of the Gemma model (Gemma-Team et al., 2024).¹⁸

Our empirical study employs the Vietnamese truthful question-answering dataset ViTruthfulQA (Nguyen et al., 2023), comprising 199 questions.¹⁹ Each of the 199 questions is fed into 11 experimental models to generate responses, which are then anonymously shuffled. Here, we utilize the greedy search decoding method, which is more suitable for LLM comparison (Lin & Chen, 2023). Two annotators independently assess each generated response on whether the response is correct or not. A response is annotated as “correct” only if it contains the accurate answer for the corresponding question without any hallucinated information. We host a discussion session with the annotators to resolve annotation conflicts.

Table 1 presents the accuracy results obtained, showing that overall, our 4B-parameter PhoGPT-4B-Chat is highly competitive compared to the closed-source GPT-4 model, yielding better accuracy scores than the remaining closed-source models GPT-3.5-turbo and Gemini Pro 1.0, as well as all open-source baselines. Furthermore, when it comes to 147 out of the 199 questions that specifically ask for information related to Vietnam, PhoGPT-4B-Chat achieves the highest accuracy.

4 CONCLUSION

We have introduced state-of-the-art open-source 4B-parameter LLMs for Vietnamese, including the base pre-trained PhoGPT-4B and its chat variant, PhoGPT-4B-Chat. We hope that these models will foster future research and applications of Vietnamese LLMs.

¹⁴<https://huggingface.co/sail/Sailor-7B-Chat>

¹⁵<https://huggingface.co/sail/Sailor-4B-Chat>

¹⁶<https://huggingface.co/LR-AI-Labs/vbd-llama2-7B-50b-chat>

¹⁷<https://huggingface.co/vilm/vinallama-7b-chat>

¹⁸<https://huggingface.co/google/gemma-7b-it>

¹⁹The original dataset consists of 213 questions; however, 14 of them are not questions or are unclear, e.g. “Tên gọi nào không phải là một loại trái cây Việt Nam?” (Which name is not a Vietnamese fruit?). Therefore, we remove them, resulting in a final evaluation set of 199 questions.

LIMITATIONS

PhoGPT has certain limitations. For example, it is not good at tasks involving reasoning, coding or mathematics. PhoGPT may generate harmful, hate speech, biased responses, or answer unsafe questions. Users should be cautious when interacting with PhoGPT that can produce factually incorrect output.

ACKNOWLEDGMENTS

We extend our thanks to Nhung Nguyen (v.nhungnt89@vinai.io) for crawling and pre-processing health data and to Thien Huu Nguyen (v.thienh4@vinai.io) for the initial discussions.

REFERENCES

- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. Falcon-40B: an open large language model with state-of-the-art performance. 2023.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen Technical Report. *arXiv preprint*, arXiv:2309.16609, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. In *Proceedings of NeurIPS*, 2020.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness. In *Proceedings of NeurIPS*, 2022.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations, 2023.
- Gemma-Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Pier Giuseppe Sessa, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimentko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel,

- Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open Models Based on Gemini Research and Technology. *arXiv preprint*, arXiv:2403.08295, 2024.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. Mistral 7B. *arXiv preprint*, arXiv:2310.06825, 2023.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of SIGOPS*, 2023.
- Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji, and Timothy Baldwin. Bactrian-X : A Multilingual Replicable Instruction-Following Model with Low-Rank Adaptation. *arXiv preprint*, arXiv:2305.15011, 2023.
- Yen-Ting Lin and Yun-Nung Chen. LLM-Eval: Unified Multi-Dimensional Automatic Evaluation for Open-Domain Conversations with Large Language Models. In Yun-Nung Chen and Abhinav Rastogi (eds.), *Proceedings of NLP4ConvAI*, pp. 47–58, 2023.
- Minh Thuan Nguyen, Khanh Tung Tran, Nhu Van Nguyen, and Xuan-Son Vu. ViGPTQA - state-of-the-art LLMs for Vietnamese question answering: System overview, core models training, and evaluations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 754–764, 2023.
- Ofir Press, Noah Smith, and Mike Lewis. Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation. In *Proceedings of ICLR*, 2022.
- MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023. URL www.mosaicml.com/blog/mpt-7b. Accessed: 2023-05-05.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timoth e Lacroix, Baptiste Rozi re, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and Efficient Foundation Language Models. *arXiv preprint*, arXiv:2302.13971, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models. *arXiv preprint*, arXiv:2307.09288, 2023b.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *Proceedings of NIPS*, pp. 5998–6008, 2017.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R mi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of EMNLP: System Demonstrations*, pp. 38–45, 2020.

This figure "HumanResults.png" is available in "png" format from:

<http://arxiv.org/ps/2311.02945v3>