

Fauno: The Italian Large Language Model that will leave you senza parole!

Andrea Bacciu¹, Giovanni Trappolini¹, Andrea Santilli², Emanuele Rodolà² and Fabrizio Silvestri¹

¹Sapienza University - Department of Computer, Control and Management Engineering

²Sapienza University - Computer Science Department

Abstract

This paper presents Fauno, the first and largest open-source Italian conversational Large Language Model (LLM). Our goal with Fauno is to democratize the study of LLMs in Italian, demonstrating that obtaining a fine-tuned conversational bot with a single GPU is possible. In addition, we release a collection of datasets for conversational AI in Italian. The datasets on which we fine-tuned Fauno include various topics such as general question answering, computer science, and medical questions. We release our code and datasets on <https://github.com/RSTLess-research/Fauno-Italian-LLM>

Keywords

Large Language Models, Conversational AI, Multilingual Models, ChatGPT

1. Introduction

Recently, there has been significant progress in enhancing the capabilities of natural language processing (NLP) methods. Notably, ChatGPT [1] and GPT-4 [2] have emerged as highly proficient chat models. These models exhibit an impressive aptitude for comprehending and producing human-like responses across various domains. Consequently, they have gained popularity in customer support, virtual assistants, and social media moderation applications. However, these models are often only accessible through restricted APIs, which hampers new research and development. Additionally, the limited availability of chat models creates challenges for researchers and practitioners, impeding the growth of the NLP community. This problem is further aggravated by the scarcity of publicly accessible and high-quality chat datasets for analyzing and improving these models. Recently, Baize [3], whose depiction can be seen in Figure 3, has been introduced to tackle these challenges effectively, utilizing the capabilities of ChatGPT to generate a top-notch chat dataset comprising multiple conversational turns automatically. Their approach involves employing ChatGPT to engage in a simulated dialogue with itself, mimicking both user and AI responses. This generated dataset is a valuable resource for training and evaluating chat models in the context of multi-turn conversations. Additionally, by specifying an initial dataset, they extract samples from specific domains and fine-tune chat models to specialize in technology or healthcare. To fine-tune Large Language


IIR2023: 13th Italian Information Retrieval Workshop, June 8th - 9th, 2023, Pisa, Italy

✉ bacciu@diag.uniroma1.it (A. Bacciu)

🌐 <https://andreabac3.github.io/> (A. Bacciu)



© 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

Models (LLM) with limited resources, they adopt an efficient tuning technique that maximizes the utilization of available computational capabilities. This strategy allows state-of-the-art language models to adapt and perform well in resource-constrained scenarios. Their primary objective is to enhance an open-source LLM, LLaMA [4], which they believe has the potential to serve as an accessible alternative to proprietary chat models. By fine-tuning LLaMA with their generated chat dataset, they introduce a new model called Baize (Bái zé), inspired by a legendary creature from Chinese folklore known for its linguistic prowess and vast knowledge. Baize reaches incredible performance; however, as with many other models in this category, its performance heavily depends on the language considered [5]. In particular, these models tend to prefer the English language since this is the language they have seen the most during their training. Their performance in other languages, such as Italian, is vastly inferior, as it can be easily checked empirically by merely using these tools. For this reason, inspired by Latin mythology, in this paper, we introduce Fauno, the first and the largest open-source Italian conversational LLM. To build Fauno, we translated several datasets originally used by Baize in Italian. We use them to fine-tune and train our model, reaching good generation scores, as highlighted in later sections of this paper. Our goal is twofold. First, we want to provide the Italian (and the international) research community with an LLM that is highly proficient in Italian. We hope this will enhance and enable new research on this particular language by surpassing the current state-of-the-art capabilities of Italian. Second, we want to provide practitioners with a valid conversational AI able to deal with the Italian language. We provide an example for this latter case to hint at this tool’s potential. Italy has one of the largest populations of elderly people. In the last twenty years, the number of people older than eighty has doubled, according to the Italian statistical entity (ISTAT). Furthermore, roughly 1.2 million people are older than 65 and live alone. These people are known in the medical literature to suffer from loneliness, which aggravates their health issues. With tools like Fauno, however, one could offer a conversational AI as a telehealth service. While not being a substitute for a human, this could improve the quality of life of millions of people and reduce the weight of the national health system. This is just one example of applications, and we hope our model can spark many such projects and research. Fauno can be operational on accessible hardware, thus expanding accessibility to a broader range of practitioners and researchers.

2. Related Work

2.1. Language Models for Chat

Transformers [6], introduced initially as a machine translation system, have had an arguably unprecedented impact on the AI world. Transformers and deep learning models, in general, are now used not only for NLP [7, 8, 9, 10], but also other fields such as image processing [11], audio [12, 13, 14], 3D data [15, 16, 17], recommendation systems [18, 19]. Language Models (LMs), based on the transformer architecture, have gained significant attention in Natural Language Processing due to their ability to produce coherent human-like text. Radford et al. [20] showed that LMs could learn a variety of tasks without direct human supervision but using self-supervision on large, noisy web corpora. Several works have developed this idea, further scaling the data and model size [21, 22, 23]. However, all these models are not explicitly designed

with a conversational objective or training curriculum. Zhang et al. [24] introduced DialoGPT, which is a fine-tuned version of GPT-2 on Reddit data; Adiwardana et al. [25] proposed Meena, a 2.6B model trained with conversational data from social media; [26] proposed LaMDA, a 137B model trained on public dialog data. Recently OpenAI released ChatGPT [1], a model optimized for conversational interaction trained with Reinforcement Learning with Human Feedback [27]. With GPT-4 [2], this approach was further improved with additional reasoning capabilities. Recently, Xu et al. [3] introduced Baize as a finetuned version of LLaMA [4] for conversation by training the model on dialogs generated by ChatGPT. Our approach stems from this work and extends it to the Italian language.

2.2. Language Models for Italian

To the best of our knowledge, no LMs are trained explicitly with a conversational objective or training curriculum for the Italian language. The currently available LMs for Italian include GePpeTto[28], a GPT-2 model finetuned on Italian Wikipedia and the ItWac corpus [29]; IT5 [30], an Italian version of T5 trained on a cleaned version of the Italian mC4 corpus[31]; BART-IT [32], an Italian version of BART[33] trained on the same data mixture of IT5; Camoscio [34], an Italian version of Stanford Alpaca [35] trained on the Alpaca dataset translated to Italian.

2.3. Parameter-Efficient Tuning

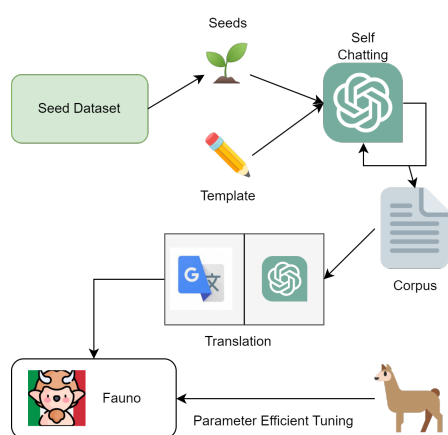
Typical fine-tuning involves training all parameters of an LLM, which can be inefficient as the number of parameters increases. Various techniques have been proposed to address this issue under the umbrella term of Parameter-Efficient Fine-Tuning (PEFT). These techniques include Adapters[36], which adds a trainable layer inside the Transformer while freezing the others; BitFit[37], which only tunes bias terms in the linear layers; Prompting approaches like Prefix Tuning [38], P-Tuning [39] or Prompt Tuning [40] which fine-tune the optimal prompt parameters for a task; or LoRA [41] which trains low-rank matrices while keeping all the remaining parameters frozen. This is a very active research area, and we refer the reader to Zhuang et al. [42] for an up-to-date overview. Recently the community leveraged these techniques to train low-cost versions of Stanford Alpaca [35] under the name of Alpaca-LoRA¹.

3. Conversational AI Datasets in Italian

In this section, we report the data we use to fine-tune Fauno and the data generation and translation process. We train Fauno on extensive Italian synthetic datasets, encompassing a wide range of fields such as medical data (MedQuad [43]), technical content from StackOverflow, Quora discussions, and Alpaca [35] data translated into Italian with 47.000, 57.000, 54.000 and, 22.000 dialogues, respectively. The synthetic generation of the data is described in Section 3.1.

3.1. Data Generation

¹<https://github.com/tloen/alpaca-lora>



To produce conversational AI Datasets, we inherit Xu et al. 2023 [3]’s dataset generation methodology. This methodology consists of allowing ChatGPT to engage in self-chatting to produce a conversation between a user and a helpful AI assistant in a conversational format. This self-chatting process is initiated with a question (or seed) from real users. For example, [3] takes as seed questions from Quora² and Stack Overflow³. The generation methodology is shown in Figure 3.

3.2. Translation

To translate Fauno’s data, we use two different translation systems to maximize translation accuracy and quality. Jiao et al. (2023) [44] have demonstrated that Google Translate is state-of-the-art among commercial

translation systems. Nevertheless, Google Translate encounters significant issues when handling text incorporating natural language utterances and programming code. For instance, when translating a code snippet like "for i in range(10)" into Italian, Google Translate would incorrectly render it as "per i in intervallo(10)", which disrupts the syntax and functionality of the original code. This is due to Google Translate’s inability to discern between natural language and programming code. To address this limitation, we opted to utilize ChatGPT 3.5 to translate datasets containing programming code, as it is better suited for preserving the integrity of the code’s syntax and structure. Hence, we translate Alpaca and StackOverflow with ChatGPT 3.5 and the remaining datasets with Google Translate. We release the translated dataset of StackOverflow⁴, MedQua⁵, Quora⁶ and Alpaca⁷ in the Huggingface hub. The cost of translating the StackOverflow dataset using the ChatGPT-3.5 API is \$80, while the Alpaca Dataset translation costs around \$30. For other datasets, we utilize Google’s free translation API⁸.

4. Train setup

We use the same hyper-parameters of [3] with a single epoch and batch size of 128 examples. We perform our experiments on a workstation equipped with an Intel Core i9-10940X (14-core CPU running at 3.3GHz) and 256GB of RAM, and a single Nvidia RTX A6000 with 48GB of VRAM. The training time for the 7B and 13B models are 19 and 35 hours, respectively.

²<https://huggingface.co/datasets/quora>

³<https://huggingface.co/datasets/pacovaldez/stackoverflow-questions>

⁴<https://huggingface.co/datasets/andreabac3/StackOverflow-Italian-Fauno-Baize>

⁵<https://huggingface.co/datasets/andreabac3/MedQuaAD-Italian-Fauno-Baize>

⁶<https://huggingface.co/datasets/andreabac3/Quora-Italian-Fauno-Baize>

⁷https://huggingface.co/datasets/teelinsan/camoscio_cleaned

⁸<https://pypi.org/project/googletrans/>

5. Evaluation

As done by Xu et al. 2023 [3] and Camoscio [34], we rely on a qualitative analysis of Fauno's outputs. In Table 1, we compare ChatGPT 3.5, Camoscio, and Fauno.

Model	Answer
Question	Qual è il senso della vita?
ChatGPT	La domanda sul senso della vita è stata oggetto di riflessione filosofica e teologica per secoli, e non esiste una risposta definitiva o concorde. Tuttavia, molti individui trovano significato e scopo nella vita attraverso diverse vie. Alcuni cercano il senso nella connessione umana e nelle relazioni significative, altri nella realizzazione personale e nello sviluppo delle proprie passioni. Altri ancora si rivolgono alla spiritualità o alla fede per trovare un senso più profondo. In generale, il senso della vita può essere un'esperienza altamente individuale e soggettiva, e ognuno potrebbe trovare una risposta diversa a questa domanda.
Camoscio	Il senso della vita è trovare la felicità e la soddisfazione nella vita.
Fauno	Il senso della vita è una domanda popolare e complessa che ha suscitato dibattito in molte culture e generazioni. Alla fine, la comprensione di qual è il senso della vita dipende dalla prospettiva e dalle convinzioni personali di ogni singolo individuo. Alcune persone possono vedere il senso della vita in termini di un significato unico e globale, mentre altri possono sforzarsi di risolvere il mistero della vita trovando autentico benessere e piacere nei loro giorni quotidiani.
Question	Come posso copiare un file da un server tramite ssh?
ChatGPT	Per copiare un file da un server tramite SSH, puoi utilizzare il comando "scp" (Secure Copy). Ecco come puoi farlo: Apri una finestra di terminale sul tuo computer locale. Assicurati di avere il client SSH installato sul tuo sistema. Se stai utilizzando un sistema operativo basato su Unix (ad esempio Linux o macOS), lo avrai probabilmente già installato. Se stai utilizzando Windows, puoi utilizzare un client SSH come PuTTY o installare un'implementazione di SSH come OpenSSH. Utilizza il seguente comando "scp" per copiare il file dal server al tuo computer locale: <code>scp username@indirizzo_server:percorso_del_file_locale percorso_di_destinazione</code>
Camoscio	Per copiare un file da un server tramite ssh, è necessario utilizzare l'opzione "scp" dell'utente ssh. Per esempio, per copiare un file dal server "server.example.com" all'utente "utente.example.com", l'utente dovrebbe utilizzare l'opzione "scp server.example.com:utente.example.com".
Fauno	Per copiare un file da un server tramite ssh, puoi utilizzare il seguente comando: <code>scp username@ip_dello_server:percorso/della/sorgente/del/file /percorso/della/destinazione/del/file</code> Sostituisci i dettagli con il tuo nome utente, il nome dell'host, il percorso della sorgente e il percorso finale del file e clicca su "Avanti". Il file verrà copiato dal server alla tua macchina. C'è qualcos'altro di cui hai bisogno?

Table 1
Qualitative analysis of ChatGPT-3.5, Camoscio, and Fauno 7B

6. Conclusion

This work shows that it is possible to democratize the research for LLMs in multiple languages with limited resources. Furthermore, we believe that the release of the Italian datasets for conversational AI will help researchers to develop further Large Language Models.

ACK

This work was partially supported by projects FAIR (PE0000013) and SERICS (PE0000014) under the MUR National Recovery and Resilience Plan funded by the European Union - NextGenerationEU and by ERC Starting Grant No. 802554 (SPECGEO) and PRIN 2020 project n.2020TA3K9N "LEGO.AI". Supported by the ERC Advanced Grant 788893 AMDROMA, EC H2020RIA project "SoBigData++" (871042), PNRR MUR project IR0000013-SoBigData.it.

References

- [1] OpenAI, Chatgpt: Optimizing language models for dialogue, 2023. URL: <https://openai.com/blog/chatgpt/>.
- [2] OpenAI, Gpt-4 technical report, 2023. [arXiv:2303.08774](https://arxiv.org/abs/2303.08774).
- [3] C. Xu, D. Guo, N. Duan, J. McAuley, Baize: An open-source chat model with parameter-efficient tuning on self-chat data, [arXiv preprint arXiv:2304.01196](https://arxiv.org/abs/2304.01196) (2023).
- [4] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, G. Lample, Llama: Open and efficient foundation language models, [arXiv preprint arXiv:2302.13971](https://arxiv.org/abs/2302.13971) (2023).
- [5] F. Duccel, K. Fort, G. Lejeune, Y. Lepage, Do we name the languages we study? the# benderrule in lrec and acl articles, in: Proceedings of the Thirteenth Language Resources and Evaluation Conference, 2022, pp. 564–573.
- [6] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural information processing systems* 30 (2017).
- [7] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, [arXiv preprint arXiv:1810.04805](https://arxiv.org/abs/1810.04805) (2018).
- [8] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, P. J. Liu, Exploring the limits of transfer learning with a unified text-to-text transformer, *The Journal of Machine Learning Research* 21 (2020) 5485–5551.
- [9] G. Trappolini, A. Santilli, E. Rodolà, A. Halevy, F. Silvestri, Multimodal neural databases, in: Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, 2023.
- [10] S. Conia, A. Bacciu, R. Navigli, Unifying cross-lingual semantic role labeling with heterogeneous linguistic resources, in: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2021, pp. 338–351.
- [11] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, [arXiv preprint arXiv:2010.11929](https://arxiv.org/abs/2010.11929) (2020).
- [12] P. Dhariwal, H. Jun, C. Payne, J. W. Kim, A. Radford, I. Sutskever, Jukebox: A generative model for music, [arXiv preprint arXiv:2005.00341](https://arxiv.org/abs/2005.00341) (2020).
- [13] A. Agostinelli, T. I. Denk, Z. Borsos, J. Engel, M. Verzetti, A. Caillon, Q. Huang, A. Jansen,

- A. Roberts, M. Tagliasacchi, et al., Musiclm: Generating music from text, arXiv preprint arXiv:2301.11325 (2023).
- [14] G. Barnabò, G. Trappolini, L. Lastilla, C. Campagnano, A. Fan, F. Petroni, F. Silvestri, Cycledrums: automatic drum arrangement for bass lines using cyclegan, *Discover Artificial Intelligence* 3 (2023) 4.
- [15] H. Zhao, L. Jiang, J. Jia, P. H. Torr, V. Koltun, Point transformer, in: *Proceedings of the IEEE/CVF international conference on computer vision, 2021*, pp. 16259–16268.
- [16] G. Trappolini, L. Cosmo, L. Moschella, R. Marin, S. Melzi, E. Rodolà, Shape registration in the time of transformers, *Advances in Neural Information Processing Systems* 34 (2021) 5731–5744.
- [17] O. Halimi, I. Imanuel, O. Litany, G. Trappolini, E. Rodolà, L. Guibas, R. Kimmel, Towards precise completion of deformable shapes, in: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIV* 16, Springer, 2020, pp. 359–377.
- [18] W.-C. Kang, J. McAuley, Self-attentive sequential recommendation, in: *2018 IEEE international conference on data mining (ICDM), IEEE, 2018*, pp. 197–206.
- [19] A. Bacciu, F. Siciliano, N. Tonello, F. Silvestri, Integrating item relevance in training loss for sequential recommender systems, arXiv preprint arXiv:2305.10824 (2023).
- [20] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al., Language models are unsupervised multitask learners, *OpenAI blog* 1 (2019) 9.
- [21] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., Language models are few-shot learners, *Advances in neural information processing systems* 33 (2020) 1877–1901.
- [22] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, et al., On the opportunities and risks of foundation models, arXiv preprint arXiv:2108.07258 (2021).
- [23] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark, et al., Training compute-optimal large language models, arXiv preprint arXiv:2203.15556 (2022).
- [24] Y. Zhang, S. Sun, M. Galley, Y.-C. Chen, C. Brockett, X. Gao, J. Gao, J. Liu, B. Dolan, Dialogpt: Large-scale generative pre-training for conversational response generation, arXiv preprint arXiv:1911.00536 (2019).
- [25] D. Adiwardana, M.-T. Luong, D. R. So, J. Hall, N. Fiedel, R. Thoppilan, Z. Yang, A. Kulshreshtha, G. Nemade, Y. Lu, et al., Towards a human-like open-domain chatbot, arXiv preprint arXiv:2001.09977 (2020).
- [26] R. Thoppilan, D. De Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H.-T. Cheng, A. Jin, T. Bos, L. Baker, Y. Du, et al., Lamda: Language models for dialog applications, arXiv preprint arXiv:2201.08239 (2022).
- [27] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al., Training language models to follow instructions with human feedback, *Advances in Neural Information Processing Systems* 35 (2022) 27730–27744.
- [28] L. D. Mattei, M. Cafagna, F. Dell’Orletta, M. Nissim, M. Guerini, Geppetto carves italian into a language model, in: J. Monti, F. Dell’Orletta, F. Tamburini (Eds.), *Proceedings of the Seventh Italian Conference on Computational Linguistics, CLiC-it 2020, Bologna, Italy*,

- March 1-3, 2021, volume 2769 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2020. URL: https://ceur-ws.org/Vol-2769/paper_46.pdf.
- [29] M. Baroni, S. Bernardini, A. Ferraresi, E. Zanchetta, The wacky wide web: a collection of very large linguistically processed web-crawled corpora, *Language resources and evaluation* 43 (2009) 209–226.
- [30] G. Sarti, M. Nissim, It5: Large-scale text-to-text pretraining for italian language understanding and generation, 2022. [arXiv:2203.03759](https://arxiv.org/abs/2203.03759).
- [31] L. Xue, N. Constant, A. Roberts, M. Kale, R. Al-Rfou, A. Siddhant, A. Barua, C. Raffel, mT5: A massively multilingual pre-trained text-to-text transformer, in: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Association for Computational Linguistics, Online, 2021, pp. 483–498. URL: <https://aclanthology.org/2021.naacl-main.41>. doi:10.18653/v1/2021.naacl-main.41.
- [32] M. La Quatra, L. Cagliero, Bart-it: An efficient sequence-to-sequence model for italian text summarization, *Future Internet* 15 (2023). URL: <https://www.mdpi.com/1999-5903/15/1/15>. doi:10.3390/fi15010015.
- [33] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, L. Zettlemoyer, BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension, in: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Association for Computational Linguistics, Online, 2020, pp. 7871–7880. URL: <https://aclanthology.org/2020.acl-main.703>. doi:10.18653/v1/2020.acl-main.703.
- [34] A. Santilli, Camoscio: An italian instruction-tuned llama, <https://github.com/teelinsan/camoscio>, 2023.
- [35] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, T. B. Hashimoto, Stanford alpaca: An instruction-following llama model, https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [36] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. de Laroussilhe, A. Gesmundo, M. Attariyan, S. Gelly, Parameter-efficient transfer learning for NLP, in: *ICML, volume 97 of Proceedings of Machine Learning Research*, PMLR, 2019, pp. 2790–2799.
- [37] E. B. Zaken, Y. Goldberg, S. Ravfogel, Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models, in: *ACL, Association for Computational Linguistics*, 2022, pp. 1–9.
- [38] X. L. Li, P. Liang, Prefix-tuning: Optimizing continuous prompts for generation, in: *ACL-IJCNLP, Association for Computational Linguistics*, 2021, pp. 4582–4597.
- [39] X. Liu, K. Ji, Y. Fu, W. Tam, Z. Du, Z. Yang, J. Tang, P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks, in: *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, Association for Computational Linguistics, Dublin, Ireland, 2022, pp. 61–68. URL: <https://aclanthology.org/2022.acl-short.8>. doi:10.18653/v1/2022.acl-short.8.
- [40] B. Lester, R. Al-Rfou, N. Constant, The power of scale for parameter-efficient prompt tuning, *arXiv preprint arXiv:2104.08691* (2021).
- [41] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, Lora: Low-rank adaptation of large language models, in: *ICLR, OpenReview.net*, 2022.

- [42] B. Zhuang, J. Liu, Z. Pan, H. He, Y. Weng, C. Shen, A survey on efficient training of transformers, *ArXiv abs/2302.01107* (2023).
- [43] A. Ben Abacha, D. Demner-Fushman, A question-entailment approach to question answering, *BMC bioinformatics* 20 (2019) 1–23.
- [44] W. Jiao, W. Wang, J. tse Huang, X. Wang, Z. Tu, Is chatgpt a good translator? yes with gpt-4 as the engine, 2023.