# LLM-Adapters: An Adapter Family for Parameter-Efficient Fine-Tuning of Large Language Models

**Zhiqiang Hu**[1]    **Yihuai Lan**    **Lei Wang**[2*]    **Wanyu Xu**[3]    **Ee-Peng Lim**[2]

**Roy Ka-Wei Lee**[1]    **Lidong Bing**[4]    **Xing Xu**[5]    **Soujanya Poria**[1]

[1]**Singapore University of Technology and Design**
[2]**Singapore Management University**
[3]**Southwest Jiaotong University**
[4]**DAMO Academy, Alibaba Group**
[5]**University of Electronic Science and Technology of China**

## Abstract

The success of large language models (LLMs), like GPT-3 and ChatGPT, has led to the development of numerous cost-effective and accessible alternatives that are created by fine-tuning open-access LLMs with task-specific data (e.g., ChatDoctor) or instruction data (e.g., Alpaca). Among the various fine-tuning methods, adapter-based parameter-efficient fine-tuning (PEFT) is undoubtedly one of the most attractive topics, as it only requires fine-tuning a few external parameters instead of the entire LLMs while achieving comparable or even better performance. To enable further research on PEFT methods of LLMs, this paper presents **LLM-Adapters**, an easy-to-use framework that integrates various adapters into LLMs and can execute these adapter-based PEFT methods of LLMs for different tasks. The framework includes state-of-the-art open-access LLMs such as LLaMA, BLOOM, OPT, and GPT-J, as well as widely used adapters such as Series adapter, Parallel adapter, and LoRA. The framework is designed to be research-friendly, efficient, modular, and extendable, allowing the integration of new adapters and the evaluation of them with new and larger-scale LLMs. Furthermore, to evaluate the effectiveness of adapters in LLMs-Adapters, we conduct experiments on six math reasoning datasets. The results demonstrate that using adapter-based PEFT in smaller-scale LLMs (7B) with few extra trainable parameters yields comparable, and in some cases superior, performance to that of powerful LLMs (175B) in zero-shot inference on simple math reasoning datasets. Overall, we provide a promising framework for fine-tuning large LLMs on downstream tasks. We believe the proposed LLMs-Adapters will advance adapter-based PEFT research, facilitate the deployment of research pipelines, and enable practical applications to real-world systems. We will keep all of the code open-source and continue to update the framework with new adapters, LLMs, and tasks. The code of our framework can be found at https://github.com/AGI-Edgerunners/LLM-Adapters.

## 1   Introduction

Large language models (LLMs) such as GPT-3 (Brown et al., 2020), BLOOM (Scao et al., 2022), and LLaMA (Touvron et al., 2023) have shown impressive performance in various natural language
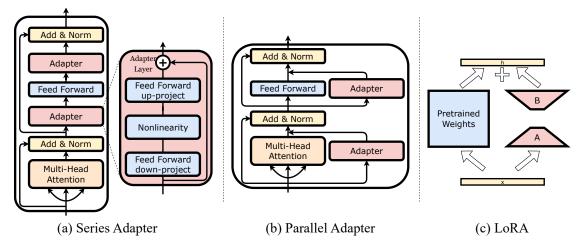
---

*Corresponding Author

Figure 1: A detailed illustration of the model architectures of three different adapters: (a) Series Adapter (Houlsby et al., 2019), (b) Parallel Adapter (He et al., 2021), and (c) LoRA (Hu et al., 2021).

processing (NLP) tasks. Fine-tuning (Brown et al., 2020; Raffel et al., 2020; Wei et al., 2021; Taori et al., 2023) is a popular technique that involves adapting LLMs to specific downstream tasks by training LLMs on task-specific datasets. However, the most powerful LLMs, specifically instruction-following LLMs, such as ChatGPT[1] and GPT-4[2], are currently in a closed-source state. Instead, these models only offer user interfaces or APIs, and their source code is not accessible to researchers or developers. Therefore, researchers or developers are unable to utilize these LLMs as backbone models for developing fine-tuning methods for specific downstream tasks. The lack of access to the source code of these LLMs hinders the innovation and advance of the state-of-the-art.

To mitigate this issue, The team of Stanford Alpaca (Taori et al., 2023) utilizes the Self-instruct strategy (Wang et al., 2022) to generate an instruction-following dataset comprising 52k data examples and fine-tune the open-source LLaMA (7B) model using these instructions to obtain an instruction-following LLaMA model. However, the massive size of LLMs, often comprising billions of parameters (Brown et al., 2020; Zhang et al., 2022; Chowdhery et al., 2022; Touvron et al., 2023), makes fine-tuning the entire LLM for a downstream task a highly challenging and inefficient. To overcome this challenge, Alpaca-LoRA[3] proposes integrating a parameter-efficient fine-tuning (PEFT) (Houlsby et al., 2019; Lester et al., 2021; Mangrulkar et al., 2022; Fu et al., 2022) method, Low-Rank Adaptation (LoRA) (Hu et al., 2021), into Alpaca. This integration allows the parameter-efficiently fine-tuned model to achieve comparable performance to full fine-tuning but with few trainable parameters. Tthe success and potential of Alpaca and Alpaca-LoRA sparke a range of adaptations and applications, including Chinese Alpaca[4], Japanese Alpaca[5], Thai Alpaca[6], medical Alpaca (ChatDoctor) (Yunxiang et al., 2023), movie recommendation Alpaca (RecAlpaca) (Wang et al., 2023), multi-modal Alpaca (LLaMA-Adapter) (Zhang et al., 2023). Additionally, the toolbox of fine-tuning LLM, known as LMFlow (Diao et al., 2023), has been developed.

In this paper, we draw inspiration from the Adapter-Transformer toolkit (Pfeiffer et al., 2020a) and introduce **LLMs-Adapters**, a user-friendly framework that seamlessly integrates a range of adapters into LLMs and enables efficient and effective adapter-based PEFT of LLMs on different tasks. The propposed framework incorporates cutting-edge open-access LLMs including LLaMA-7B (Touvron et al., 2023), BLOOM-7.1B (Scao et al., 2022), and GPT-J (Wang and Komatsuzaki, 2021), as well as

---

[1] https://chat.openai.com/chat    [2] https://chat.openai.com/chat?model=gpt-4    [3] https://github.com/tloen/alpaca-lora

[4] https://github.com/LC1332/Chinese-alpaca-lora    [5] https://github.com/kunishou/Japanese-Alpaca-LoRA

[6] https://huggingface.co/Thaweewat/thai-buffala-lora-7b-v0-1

popular adapters (Figure 1) such as Series Adapter (Houlsby et al., 2019), Parallel Adapter (Pfeiffer et al., 2020b; He et al., 2021), and LoRA (Hu et al., 2021). Through the use of our proposed LLMs-Adapters, we can fine-tune a small number of external adapter parameters while leaving the parameters of the pretrained LLMs unchanged. We evaluate the performance of the introduced LLMs-Adapters by conducting experiments on six math reasoning datasets, including AQuA (Ling et al., 2017), GSM8K (Cobbe et al., 2021), MultiArith (Roy and Roth, 2016), AddSub (Hosseini et al., 2014), SingleEq (Koncel-Kedziorski et al., 2015), and SVAMP (Patel et al., 2021). The experimental results demonstrate that on simple math reasoning datasets (i.e., AddSub, SingleEq, and MultiArith), using adapter-based PEFT in smaller-scale LLMs (7B) with few extra trainable parameters can yield comparable or even superior performance to that of powerful LLMs (175B) in zero-shot inference.

Overall, we offer a promising framework for fine-tuning large LLMs on a range of downstream tasks. We envision that LLMs-Adapters will serve as a valuable resource for advancing research on adapter-based PEFT of LLMs, facilitating the deployment of such research pipelines, and enabling practical applications of this technique to real-world systems. Additionally, we will keep all of the code open source and will continue to update the framework with new adapters, LLMs, and tasks.

## 2   Adapter Family

Adapters for LLMs refer to neural modules integrated into LLMs, which contain a small number of extra trainable parameters, allowing for efficient fine-tuning of specific tasks without affecting the pre-trained parameters of the LLM. Here we use the notation $\Theta$ to represent the parameters of the LLM and $\Phi$ to represent the parameters of the adapter module. During training, the parameters of the LLM ($\Theta$) remain fixed, while the parameters of the adapter module ($\Phi$) are adjusted to perform a specific task. As a result, the representations generated by the LLM are not distorted due to task-specific tuning, while the adapter module acquires the capability to encode task-specific information.

This framework presents three types of adapters within LLMs: Series Adapter, Parallel Adapter, and LoRA. Our plan for future work involves updating the framework with new adapters.

**Series Adapter.** Inspired by Houlsby et al. (2019), our framework adds the bottleneck feed-forward layer in series to each multi-head and feed forward layer of a Transformer (Vaswani et al., 2017) block, which serves as the basis for most LLMs (Brown et al., 2020; Touvron et al., 2023; Wang and Komatsuzaki, 2021). Figure 1(a) shows that bottleneck adapters consist of a two-layer feed-forward neural network that includes a down-projection matrix, a non-linearity function, an up-projection, and a residual connection between input and output.

**Parallel Adapter.** The Parallel Adapter implements Pfeiffer et al. (2020b)'s architecture, which integrates bottleneck feed-forward layers in Parallel with the multi-head and feed-forward layers of a Transformer block in LLMs. As shown in Figure 1(b), the adapters are incorporated alongside each transformer layer.

**LoRA.** Hu et al. (2021) propose LoRA aimed at efficiently fine-tuning pre-trained models with fewer trainable parameters. LoRA introduces trainable low-rank decomposition matrices into LLMs' existing layers, enabling the model to adapt to new data while keeping the original LLMs fixed to retain the previous knowledge. Specifically, LoRA performs a reparameterization of each model layer expressed as a matrix multiplication by injecting low-rank decomposition matrices, as illustrated in Figure 1(c). This reparameterization enables the model to be fine-tuned without the need to compute the full dense matrix multiplication. By reducing the rank of the matrices, LoRA also helps to reduce the number of parameters in fine-tuning LLMs.

3

LLMs-Adapters provides a configuration file that allows users to customize the architecture settings to facilitate flexibility and extensibility. Therefore, users can select the block and layer for inserting any adapters or use standard adapter architectures from the literature.

# 3 Evaluation

## 3.1 Datasets

Adapter methods in LLM-Adapters are evaluated on six math reasoning datasets. These are: (1) the GSM8K (Cobbe et al., 2021) dataset, consisting of linguistically diverse math word problems written by human problem writers for grade school students; (2) the SVAMP (Patel et al., 2021) benchmark, which includes one-unknown arithmetic word problem for students up to 4th grade and is derived from an existing dataset with slight modifications; (3) the MultiArith (Roy and Roth, 2016) dataset, featuring math word problems that require multiple reasoning steps and operations; (4) the AddSub (Hosseini et al., 2014) dataset, comprising of addition and subtraction arithmetic word problems; (5) the AQUA (Ling et al., 2017) dataset, which contains algebraic word problems with natural language rationales; and (6) the SingleEq (Koncel-Kedziorski et al., 2015) dataset, consisting of single-equation algebra word problems with multiple math operations over non-negative rational numbers and one variable.

In order to prepare the training data for the adapters in our experiment, we perform a random shuffling of each math reasoning dataset, followed by dividing each dataset into two subsets: a training set comprising 80% of the data and a test set containing 20% of the data. Next, we combine all the training datasets into a single, unified dataset and use it to train all the adapters. To obtain the necessary supervision signals for training the adapters, we extracted rationales and answers for training data examples from the log files of Zero-shot-COT Kojima et al. (2022) experimented with GPT-3 text-Davince-003[7]. We have acquired a total of 3260 and 816 mathematical word problems, accompanied by corresponding rationales and answers, for the purpose of fine-tuning and testing, respectively. It is noteworthy that we have not implemented any filtering procedure to exclude samples with erroneous rationales or answers.

## 3.2 Methods for Comparison

The study by Houlsby et al. (2019) introduces a novel method called **S-Adapter$^h$**, which involves integrating adapter layers after both multi-attention modules and MLP modules for each transformer layer. By solely training the parameters of adapter layers, **S-Adapter$^h$** can achieve near state-of-the-art performance on text classification tasks.In contrast, Pfeiffer et al. (2020b) propose a modified method, called **S-Adapter$^p$**, which only integrates task adapters after the MLP modules in each transformer layer. Furthermore, Pfeiffer et al. (2020b) also introduce an invertible adapter architecture for adapting a pre-trained multilingual model to a new language. Another approach, called Parallel Adapters (**P-Adapter**), integrates adapter layers in parallel with multi-head attention layers or MLP layers. Hu et al. (2021) propose **LoRA**, which freezes the pret-rained model weights and injects trainable rank decomposition matrices into each layer of the Transformer architecture, greatly reducing the number of trainable parameters for downstream tasks.

---

[7] https://platform.openai.com/docs/models/gpt-3-5

Table 1: Accuracy comparison of different LLMs parameter-efficiently fine-tuned with different adapters on six math reasoning datasets. Teacher here means Zero-shot CoT (Kojima et al., 2022). We use GPT-3.5 `text-Davinci-003` for Zero-shot CoT and instruction data generation. Params represents the number of trainable parameters.

| LLM | Method | Params | MultiArith | GSM8K | AddSub | AQuA | SingleEq | SVAMP | Average |
|---|---|---|---|---|---|---|---|---|---|
| GPT-3.5 | Teacher | - | 83.8 | 56.4 | 85.3 | 38.9 | 88.1 | 69.9 | 70.4 |
| LLaMA$_{7B}$ | LoRA | 4.2M | 88.3 | 21.9 | 78.5 | 27.5 | 83.3 | 54.5 | 59.0 |
| | S-Adapter$^h$ | 200M | 88.3 | 18.5 | 69.6 | 27.4 | 85.2 | 52.5 | 56.9 |
| | S-Adapter$^p$ | 200M | 88.3 | 18.5 | 69.6 | 15.6 | 79.4 | 52.0 | 53.9 |
| | P-Adapter | 200M | 83.3 | 22.7 | 77.2 | 9.8 | 81.3 | 57.0 | 55.2 |
| BLOOM$_{7.1B}$ | LoRA | 4M | 46.7 | 4.2 | 32.9 | 11.7 | 41.2 | 22.5 | 26.5 |
| | S-Adapter$^h$ | 125M | 60.8 | 6.4 | 43.0 | 23.5 | 52.0 | 37.5 | 37.2 |
| | S-Adapter$^p$ | 188M | 70.6 | 8.3 | 50.6 | 13.7 | 50.0 | 35.5 | 38.1 |
| | P-Adapter | 125M | 55.0 | 5.7 | 35.4 | 27.5 | 49.0 | 28.0 | 33.4 |
| GPT-J$_{6B}$ | LoRA | 3.7M | 79.2 | 10.6 | 69.6 | 2.0 | 71.6 | 45.0 | 46.3 |
| | S-Adapter$^h$ | 117M | 82.5 | 4.5 | 55.7 | 3.9 | 67.6 | 39.5 | 42.3 |
| | S-Adapter$^p$ | 176M | 79.2 | 9.8 | 54.4 | 19.6 | 63.7 | 37.5 | 44.0 |
| | P-Adapter | 176M | 79.2 | 11.0 | 65.8 | 11.8 | 69.6 | 44.5 | 47.0 |

## 3.3 Implementation Details

In order to facilitate INT8 training of large models with parallel adapters, we have adopted a technique whereby the parallel adapter layers are incorporated into multi-head attention layers and MLP layers, in parallel with Linear layers. This approach offers flexibility as the parallel layers can be integrated with any Linear layer within the model. To ensure a fair comparison among different adapter methods, we have endeavored to approximately match the fine-tuning compute budget for each adapter. To this end, we have set the rank of the LoRA A and B matrix to 8, while for **S-Adapter$^h$** and **S-Adapter$^p$**, the bottleneck size has been set to 256 and 768, respectively. For **P-Adapter**, the adapter layers have been placed in multi-head attention modules with a bottleneck size of 256.

Our evaluation of adapter methods involves the utilization of LLaMA-7B, BLOOM-7.1B, and GPT-J-6B as the backbone models. We fine-tune and infer all of the adapter methods on these models using two 3090 GPUs, each with a memory capacity of 24 GB. The process of fine-tuning each model is accomplished in a brief duration of approximately 20 minutes.

## 3.4 Results

Table 1 reports the accuracy comparison of different LLMs parameter-efficiently fine-tuned with different adapters on six math reasoning datasets. The teacher model (175B) outperforms adapter-based parameter-efficiently fine-tuned LLMs across various tasks. However, in the case of simple math reasoning datasets such as MultiArith, AddSub, and SingleEq, adapter-based methods such as LLaMA-7B with LoRA can achieve comparable performance. It indicates that given sufficient task-specific training data, adapter-based parameter-efficient fine-tuning (PEFT) of smaller LLMs may have the potential to perform similarly to very large language models. Moreover, the performance of adapter-based PEFT methods varies among LLMs of similar size. For example, the latest open-source LLMs, LLaMA, outperform GPT-J and BLOOM in most cases. Comparing different adapters,

LoRA achieves excellent performance with significantly fewer trainable parameters, suggesting that task-specific fine-tuning may not require excessive learnable parameters. Overall, these findings demonstrate the potential for adapter-based PEFT of smaller LLMs to achieve high performance on specific tasks with few trainable parameters.

## 4    Conclusion and Future Work

This paper presents LLM-Adapters, a framework that includes adapter-based parameter-efficient fine-tuning (PEFT) methods of large language models (LLMs) for various tasks. The framework is research-friendly, efficient, modular, and extendable and contains state-of-the-art open-access LLMs and widely used adapters. The experiments on six math reasoning datasets demonstrate that using adapter-based PEFT with the proposed framework LLMs-Adapters can achieve comparable or even better performance than powerful LLMs in zero-shot inference on simple math reasoning datasets. Future work includes integrating new adapters and evaluating them with larger-scale language models on more tasks to enable further research on PEFT methods of LLMs. The code is available at the provided link for further exploration and development.

## References

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. 1, 2, 3

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*. 2

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*. 3, 4

Shizhe Diao, Rui Pan, Hanze Dong, KaShun Shum, Jipeng Zhang, Wei Xiong, and Tong Zhang. 2023. Lmflow: An extensible toolkit for finetuning and inference of large foundation models. https://optimalscale.github.io/LMFlow/. 2

Zihao Fu, Haoran Yang, Anthony Man-Cho So, Wai Lam, Lidong Bing, and Nigel Collier. 2022. On the effectiveness of parameter-efficient fine-tuning. *arXiv preprint arXiv:2211.15583*. 2

Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*. 2, 3

Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. Learning to solve arithmetic word problems with verb categorization. In *EMNLP*, pages 523–533. 3, 4

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR. 2, 3, 4

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*. 2, 3, 4

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*. 4, 5

Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. 2015. Parsing algebraic word problems into equations. *Transactions of the Association for Computational Linguistics*, 3:585–597. 3, 4

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*. 2

Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167. 3, 4

Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, and Sayak Paul. 2022. Peft: State-of-the-art parameter-efficient fine-tuning methods. https://github.com/huggingface/peft. 2

Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In *Proceedings of NAACL*, pages 2080–2094. 3, 4

Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. 2020a. Adapterhub: A framework for adapting transformers. *arXiv preprint arXiv:2007.07779*. 2

Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2020b. Mad-x: An adapter-based framework for multi-task cross-lingual transfer. *arXiv preprint arXiv:2005.00052*. 3, 4

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67. 2

Subhro Roy and Dan Roth. 2016. Solving general arithmetic word problems. *arXiv preprint arXiv:1608.01413*. 3, 4

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*. 1, 2

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca. 2

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*. 1, 2, 3

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30. 3

Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. `https://github.com/kingoflolz/mesh-transformer-jax`. 2, 3

Lei Wang, Zhiqiang Hu, Yihuai Lan, Wanyu Xu, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023. Recalpaca: Low-rank llama instruct-tuning for recommendation. `https://github.com/AGI-Edgerunners/RecAlpaca`. 2

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*. 2

Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*. 2

Li Yunxiang, Li Zihan, Zhang Kai, Dan Ruilong, and Zhang You. 2023. Chatdoctor: A medical chat model fine-tuned on llama model using medical domain knowledge. *arXiv preprint arXiv:2303.14070*. 2

Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. 2023. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*. 2

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. OPT: open pre-trained transformer language models. *CoRR*, abs/2205.01068. 2