
PHOENIX: OPEN-SOURCE LANGUAGE ADAPTION FOR DIRECT PREFERENCE OPTIMIZATION

Matthias Uhlig
Hochschule Ansbach
uhlig19498@hs-ansbach.de

Sigurd Schacht
Hochschule Ansbach
sigurd.schacht@hs-ansbach.de

Sudarshan Kamath Barkur
Hochschule Ansbach
s.kamath-barkur@hs-ansbach.de

ABSTRACT

Large language models have gained immense importance in recent years and have demonstrated outstanding results in solving various tasks. However, despite these achievements, many questions remain unanswered in the context of large language models. Besides the optimal use of the models for inference and the alignment of the results to the desired specifications, the transfer of models to other languages is still an underdeveloped area of research. The recent publication of models such as Llama-2 and Zephyr has provided new insights into architectural improvements and the use of human feedback. However, insights into adapting these techniques to other languages remain scarce. In this paper, we build on latest improvements and apply the Direct Preference Optimization (DPO) approach to the German language. The model is available at <https://huggingface.co/DRXD1000/Phoenix>.

Keywords Large Language Models · Finetuning · Direct Preference Optimization · Mistral

1 Introduction

Large Language Models like GPT-3 [1], Llama-2 [2] or Mistral [3] demonstrate very good results in solving a wide range of tasks and form good foundation models for adaptation to chat models [4], [5] or other domain-specific applications. These and other open-source models benefit from a wide range of technological improvements. Structural improvements such as Grouped-Query Attention [6] or Flash-Attention [7], [8] for example were able to drastically reduce both the calculation time and the memory overhead of the attention block. Furthermore, high-quality closed-source models such as GPT-4[9], ChatGPT [10], Claude 2 [11] and Crowdsourcing efforts [12] allowed the creation of high-quality datasets. These models act as teachers, distilling the information and allowing the creation of student models trained with this distilled information.

However, despite this highly diverse field of research, which, in addition to the points already mentioned, also deals with inference optimization i.a. [13]–[16] and improving the availability of those models to the general public and low-resource researchers [17], [18], there still remains one element, the language, which has continued to be highly homogeneous. All previously mentioned enhancements to data, models, or pipelines aim to improve the results on English-based benchmarks, even though there are only around 300 million native speakers around the globe. The focus on those standard benchmarks through which models are compared results in a drastic under-representation of languages other than English. Furthermore, it also influences the selection process regarding the pretraining corpus. Although there are around 100 million German native speakers, the Llama-2 model family corpus only contains 0.17% of German texts [2]. As a result, the models are under-trained for languages other than English.

2 Related Work

Training the model in multiple languages is possible. This can be achieved using pre-training the model first with the corpus consisting of multiple languages, followed by fine-tuning for downstream tasks. For example, Google trained the mT5 model from scratch over 101 Languages [19], Meta trained the XGLM model with specific tasks for more than 20 languages [20] and the BigScience Workshop trained BLOOM Model with over 46 languages [21]. However, the large number of languages in these models means that the performance in a specific language declines sharply and does not

come close to the quality of monolingual models [22]. To solve this problem, approaches have already been developed. For example, the WECHSEL approach [22] replaces the tokenizer used with a new tokenizer(for the target language) such that the new tokens are semantically similar. The CLP-Transfer approach [23] uses a smaller, resource-efficient model in the target language. This smaller model, along with the source model, helps initialize the token embeddings of the larger model by leveraging the shared vocabulary of both languages. The rest of the weights are retained from the source language model. However, these methods require an extremely high computing effort that most organizations cannot cope with. Concerning the variants, only a GPT-2 model of [22] and a Bloom-6.7B model of [23] exist. Due to the now obsolete architectures of these models, neither Flash Attention nor Grouped Query Attention can be used with them.

Since models learn almost their entire knowledge in pretraining [24], the LAION team [25] pursued another approach by further training current models using a large language-specific text corpus in the style of casual language modeling. With an additional 65 billion training tokens, these models can outperform their original English counterparts in German tasks and partially outperform the original models in English tasks [25]. Besides this, the LAION team was also able to translate known instruction data sets into German with the help of the OpenAI API to train German chat models.

While Supervised Finetuning(SFT) provides a good approach to teach a model the conversation or response pattern of humans or other, mostly larger, models, they tend to give false or toxic answers [26]. To improve the alignment of models and increase the safety and helpfulness of LLMs, new training methods using reinforcement learning have been applied. This approach also known as Reinforcement Learning with Human Feedback(RLHF) has been used among others by the OpenAI Team to train InstructGPT [27] and Meta to improve Llama-2 [2]. InstructGPT thereby outperformed GPT-3 175B significantly in Human Evaluation although the model only has 1.3B parameters. The use of the Proximal Policy Optimization(PPO) [28] has been the State of the Art Pipeline to to train a RLHF Model [29], but the need for a separate reward model added additional complexity to the process. A newer approach, Direct Preference Optimization(DPO) [30], made this step dispensable and allowed a direct optimization of the SFT Model by comparing the scores for a chosen and rejected answer and adapting the weights of the model to maximize the rewards.

Since the release of the LAION models [25] further improvements have become available that could not be included in their training. Neftune [31], for example, has significantly improved the quality of instruction tuning. Furthermore, the work of HuggingFace for the creation of the Zephyr model [5] using Direct Preference Optimization(DPO) [30] enabled significant progress in the alignment of language models and thus made it possible to compete with Zephyr 7B against models one order of magnitude bigger. To improve the work of the LAION team [25] and to further promote the use of LLMs for the German language, we extend the approach pursued so far and optimize the Mistral [3] adaptation of the LAION team[25] by extended Supervised Finetuning and the use of DPO.

3 Method

While the HuggingFace Hub already offers a variety of translated German instruction datasets, like variants of ShareGPT, Alpaca, and Evol-Instruct provided by FreedomAI[32], Open-Platypus from the LAION team[25] and Dolly 15K [33] translated to German by ourselves [34], most of the entries consist of single turn conversations which limit their suitability for multi-turn chatbots.

3.1 Translation with open-source tools

The HuggingFace team behind the Zephyr model [5] has chosen a different dataset and used a cleaned and filtered versions of data from the Tsinghua University [35] for instruction data and UltraFeedback dataset[36] for DPO. While the adaptations of these datasets offer a high diversity of multi-turn conversations and feedback data, translating them with traditional API services would be a huge financial investment. Together, the Ultrachat version from HuggingFace[5] and the post-processed version of the Ultrafeedback dataset from Argilla[37] hold over 1.5 billion characters. Translating this amount of data for example with a commercial Translation API, like the DeepL Pro API¹, would cost over 30.000 Euros.

Due to advances in the quality of open-source translation models, such costly APIs are not the only option. That's why, for this project, the ALMA model from Microsoft[38], which has shown outstanding quality in translation tasks, was used to translate the data.

To get a high throughput while maintaining quality, the vLLM library for inference of the ALMA Model [39] has been used on a Nvidia A100 80GB GPU. The texts have been broken into chunks separated by the newline character to maintain the original text structure as best as possible. After the translation, the chunks were concatenated again. With

¹<https://www.deepl.com/en/pro/>

this approach, the translation cost of these two datasets could be reduced to approximately 30 Euros, which is 100 times cheaper. This low price has been achieved because the University of Applied Science Ansbach has the necessary infrastructure. But even with a rented GPU, the costs would sum up to about 300 Euros, which still offers huge savings.

Concerning the datasets used, special care was taken to ensure that only data whose license also permits commercial use is included. This approach made it necessary to filter the Open-Platypus dataset from the LAION Team [25], as it contains data only approved for research.

3.2 Supervised Finetuning

The process of supervised fine-tuning and alignment training is closely related to the process described in the paper for Zephyr model [5] and the Alignment Handbook by HuggingFace [40]. The base model has been swapped for the German adaption of the Mistral model². The hyperparameters used for training on an 8x A100 80GB Instance can be seen in Table 1. For the SFT training data conversations with less than 2049 tokens have been selected from the previously mentioned instruction datasets.

Parameter	Value
Batch Size	512
Number of Steps	228
Packing	True
Learning Rate Scheduler	cosine
Optimizer	Adam
Train Loss	0.8767
Eval Loss	0.8753

Table 1: Hyperparameters for SFT-Training

Parameter	Value
Batch Size	64
Number of Steps	941
Learning Rate Scheduler	Linear
Optimizer	Adam
Train Accuracy	78.75%
Eval Accuracy	82.5%

Table 2: Hyperparameters for DPO-Training

3.3 Direct Preference Optimization

As with the SFT step the DPO training and parameters in Table 1 are aligned with the Alignment Handbook [40]. The main difference can be seen in the number of training epochs. While the Zephyr model was trained for three epochs Phoenix was only trained for one. As the paper from the Zephyr model [5] shows, the improvement from additional DPO epochs in mt-bench score is only marginal. The DPO training data consists of the translated version of the Argilla Ultrafeedback dataset [37].

Metric	Value
First Turn	6.39375
Second Turn	5.1625
Categories	
Writing	7.45
Roleplay	7.9
Reasoning	4.3
Math	3.25
Coding	2.5
Extraction	5.9
STEM	7.125
Humanities	7.8
Average	5.778125

Table 3: MT-Bench results of Phoenix

²<https://huggingface.co/LeoLM/leo-mistral-hessianai-7b>

4 Evaluation

As with training, the evaluation of German models is still an underdeveloped research area. Nevertheless, the LAION team [25] has adapted the mt-bench and parts of the lm-evaluation harness [41]. The results shown in Table 3 present the answer to the German mt-bench³ of the Phoenix model. While the model still has some shortfalls in math and coding our model is able to outperform the Llama-2-70b-chat version of the LAION team[25] in Reasoning and Roleplay and performs quite well compared to other models of the same size.

While this test alone is not sufficient for evaluating the full diversity of LLMs, it is one of the only available standardized sources for the German language. This displays the necessity of further research in multilingual LLM evaluation. To extend the comparison between the models, we compared the Mistral model of the LAION team (LeoLM-Mistral-7B-Chat)[25] against Phoenix in the German versions of the HellaSwag, ARC and MMLU Benchmarks⁴.

Model	HellaSwag-de	ARC-challenge-de	MMLU-DE
LeoLM-Mistral-7B-Chat	0.47639912	0.38310580	0.38944
Phoenix	0.5281	0.4428	0.3952

Table 4: LM Evaluation Harness-DE results

5 Conclusion

In this paper we present, one of the first German DPO-aligned LLM. As the evaluation shows it performs on par with some of the best available models of its size and even manages to compete with 10 times larger models. Furthermore, we showcase how to create high-quality translated data at low cost to improve the LLM training process for different languages and task.

³<https://github.com/bjoernpl/FastEval>

⁴https://github.com/bjoernpl/lm-evaluation-harness-de/tree/mmlu_de

References

- [1] T. B. Brown, B. Mann, N. Ryder, *et al.*, “Language models are few-shot learners,”
- [2] H. Touvron, L. Martin, and K. Stone, “Llama 2: Open foundation and fine-tuned chat models,”
- [3] A. Q. Jiang, A. Sablayrolles, A. Mensch, *et al.*, “Mistral 7B,” Oct. 2023. arXiv: 2310.06825 [cs.CL].
- [4] R. Taori, I. Gulrajani, Y. Dubois, X. Li, P. Liang, and T. B. Hashimoto, “Alpaca: A strong, replicable instruction-following model,”
- [5] L. Tunstall, E. Beeching, N. Lambert, *et al.*, *Zephyr: Direct distillation of LM alignment*, Oct. 25, 2023. arXiv: 2310.16944 [cs]. [Online]. Available: <http://arxiv.org/abs/2310.16944> (visited on 10/30/2023).
- [6] J. Ainslie, J. Lee-Thorp, M. de Jong, Y. Zemlyanskiy, F. Lebrón, and S. Sanghai, *GQA: Training generalized multi-query transformer models from multi-head checkpoints*, Oct. 23, 2023. arXiv: 2305.13245 [cs]. [Online]. Available: <http://arxiv.org/abs/2305.13245> (visited on 10/30/2023).
- [7] T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré, *FlashAttention: Fast and memory-efficient exact attention with IO-awareness*, Jun. 23, 2022. arXiv: 2205.14135 [cs]. [Online]. Available: <http://arxiv.org/abs/2205.14135> (visited on 10/30/2023).
- [8] T. Dao, *FlashAttention-2: Faster attention with better parallelism and work partitioning*, Jul. 17, 2023. arXiv: 2307.08691 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.08691> (visited on 11/23/2023).
- [9] OpenAI, *GPT-4 technical report*, Mar. 27, 2023. arXiv: 2303.08774 [cs]. [Online]. Available: <http://arxiv.org/abs/2303.08774> (visited on 11/15/2023).
- [10] “Introducing ChatGPT.” (), [Online]. Available: <https://openai.com/blog/chatgpt> (visited on 11/15/2023).
- [11] “Introducing claude 2.1,” Anthropic. (), [Online]. Available: <https://www.anthropic.com/index/claude-2-1> (visited on 11/23/2023).
- [12] A. Köpf, Y. Kilcher, D. von Rütte, *et al.*, *OpenAssistant conversations – democratizing large language model alignment*, Apr. 14, 2023. arXiv: 2304.07327 [cs]. [Online]. Available: <http://arxiv.org/abs/2304.07327> (visited on 10/30/2023).
- [13] T. Dettmers, M. Lewis, Y. Belkada, and L. Zettlemoyer, *LLM.int8(): 8-bit matrix multiplication for transformers at scale*, Nov. 10, 2022. arXiv: 2208.07339 [cs]. [Online]. Available: <http://arxiv.org/abs/2208.07339> (visited on 10/30/2023).
- [14] E. Frantar, S. Ashkboos, T. Hoefler, and D. Alistarh, *GPTQ: Accurate post-training quantization for generative pre-trained transformers*, Mar. 22, 2023. arXiv: 2210.17323 [cs]. [Online]. Available: <http://arxiv.org/abs/2210.17323> (visited on 10/30/2023).
- [15] E. Frantar and D. Alistarh, *SparseGPT: Massive language models can be accurately pruned in one-shot*, Mar. 22, 2023. arXiv: 2301.00774 [cs]. [Online]. Available: <http://arxiv.org/abs/2301.00774> (visited on 10/30/2023).
- [16] J. Lin, J. Tang, H. Tang, *et al.*, *AWQ: Activation-aware weight quantization for LLM compression and acceleration*, Oct. 3, 2023. arXiv: 2306.00978 [cs]. [Online]. Available: <http://arxiv.org/abs/2306.00978> (visited on 10/30/2023).
- [17] E. J. Hu, Y. Shen, P. Wallis, *et al.*, *LoRA: Low-rank adaptation of large language models*, Oct. 16, 2021. arXiv: 2106.09685 [cs]. [Online]. Available: <http://arxiv.org/abs/2106.09685> (visited on 10/30/2023).
- [18] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, *QLoRA: Efficient finetuning of quantized LLMs*, May 23, 2023. arXiv: 2305.14314 [cs]. [Online]. Available: <http://arxiv.org/abs/2305.14314> (visited on 10/30/2023).
- [19] L. Xue, N. Constant, A. Roberts, *et al.*, “Mt5: A massively multilingual pre-trained text-to-text transformer,” Oct. 2020. arXiv: 2010.11934 [cs.CL].
- [20] X. V. Lin, T. Mihaylov, M. Artetxe, *et al.*, “Few-shot learning with multilingual language models,” Dec. 2021. arXiv: 2112.10668 [cs.CL].
- [21] BigScience Workshop, T. L. Scao, A. Fan, *et al.*, “BLOOM: A 176b-parameter open-access multilingual language model,” Nov. 2022. arXiv: 2211.05100 [cs.CL].
- [22] B. Minixhofer, F. Paischer, and N. Rekabsaz, “WECHSEL: Effective initialization of subword embeddings for cross-lingual transfer of monolingual language models,” in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022, pp. 3992–4006. DOI: 10.18653/v1/2022.naacl-main.293. arXiv: 2112.06598 [cs]. [Online]. Available: <http://arxiv.org/abs/2112.06598> (visited on 10/30/2023).

- [23] M. Ostendorff and G. Rehm, *Efficient language model training through cross-lingual and progressive transfer learning*, Jan. 23, 2023. arXiv: 2301.09626 [cs]. [Online]. Available: <http://arxiv.org/abs/2301.09626> (visited on 10/30/2023).
- [24] C. Zhou, P. Liu, P. Xu, *et al.*, *LIMA: Less is more for alignment*, May 18, 2023. arXiv: 2305.11206 [cs]. [Online]. Available: <http://arxiv.org/abs/2305.11206> (visited on 10/30/2023).
- [25] B. Plüster, *Leolm: Igniting german-language llm research*, Sep. 2023. [Online]. Available: <https://laion.ai/blog/leo-llm/>.
- [26] D. Ganguli, L. Lovitt, J. Kernion, *et al.*, *Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned*, Nov. 22, 2022. arXiv: 2209.07858 [cs]. [Online]. Available: <http://arxiv.org/abs/2209.07858> (visited on 10/30/2023).
- [27] L. Ouyang, J. Wu, X. Jiang, *et al.*, *Training language models to follow instructions with human feedback*, Mar. 4, 2022. arXiv: 2203.02155 [cs]. [Online]. Available: <http://arxiv.org/abs/2203.02155> (visited on 10/30/2023).
- [28] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” Jul. 2017. arXiv: 1707.06347 [cs.LG].
- [29] D. M. Ziegler, N. Stiennon, J. Wu, *et al.*, “Fine-tuning language models from human preferences,” Sep. 2019. arXiv: 1909.08593 [cs.CL].
- [30] R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn, *Direct preference optimization: Your language model is secretly a reward model*, May 29, 2023. arXiv: 2305.18290 [cs]. [Online]. Available: <http://arxiv.org/abs/2305.18290> (visited on 10/30/2023).
- [31] N. Jain, P.-Y. Chiang, Y. Wen, *et al.*, “NEFTune: Noisy embeddings improve instruction finetuning,” Oct. 2023. arXiv: 2310.05914 [cs.CL].
- [32] Z. Chen, S. Yan, J. Liang, *et al.*, *MultilingualSIFT: Multilingual Supervised Instruction Fine-tuning*, version 0.1, Jul. 2023. [Online]. Available: <https://github.com/FreedomIntelligence/MultilingualSIFT.git>.
- [33] M. Conover, M. Hayes, A. Mathur, *et al.* “Free dolly: Introducing the world’s first truly open instruction-tuned llm.” (2023), [Online]. Available: <https://www.databricks.com/blog/2023/04/12/dolly-first-open-commercially-viable-instruction-tuned-llm> (visited on 06/30/2023).
- [34] M. Uhlig, *Dolly-15k-german*, <https://huggingface.co/datasets/DRXD1000/Dolly-15k-German>, 2023.
- [35] N. Ding, Y. Chen, B. Xu, *et al.*, “Enhancing chat language models by scaling high-quality instructional conversations,” *arXiv preprint arXiv:2305.14233*, 2023.
- [36] G. Cui, L. Yuan, N. Ding, *et al.*, *Ultrafeedback: Boosting language models with high-quality feedback*, 2023. arXiv: 2310.01377 [cs.CL].
- [37] Argilla, *Argilla/ultrafeedback-binarized-preferences*, <https://huggingface.co/datasets/argilla/ultrafeedback-binarized-preferences>, 2023.
- [38] H. Xu, Y. J. Kim, A. Sharaf, and H. H. Awadalla, “A paradigm shift in machine translation: Boosting translation performance of large language models,” Sep. 2023. arXiv: 2309.11674 [cs.CL].
- [39] W. Kwon, Z. Li, S. Zhuang, *et al.*, “Efficient memory management for large language model serving with pagedattention,” in *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [40] L. Tunstall, E. Beeching, N. Lambert, N. Rajani, A. M. Rush, and T. Wolf, *The alignment handbook*, <https://github.com/huggingface/alignment-handbook>, 2023.
- [41] L. Gao, J. Tow, S. Biderman, *et al.*, *A framework for few-shot language model evaluation*, version v0.0.1, Sep. 2021. DOI: 10.5281/zenodo.5371628. [Online]. Available: <https://doi.org/10.5281/zenodo.5371628>.

6 Appendix

This section shows the comparison of the outputs for the same prompts between the Mixtral-8x7B Mixture of Experts(MoE) model and the Phoenix model with 7B parameters. For the comparison of the prompts side by side, we use the playground from COAI.

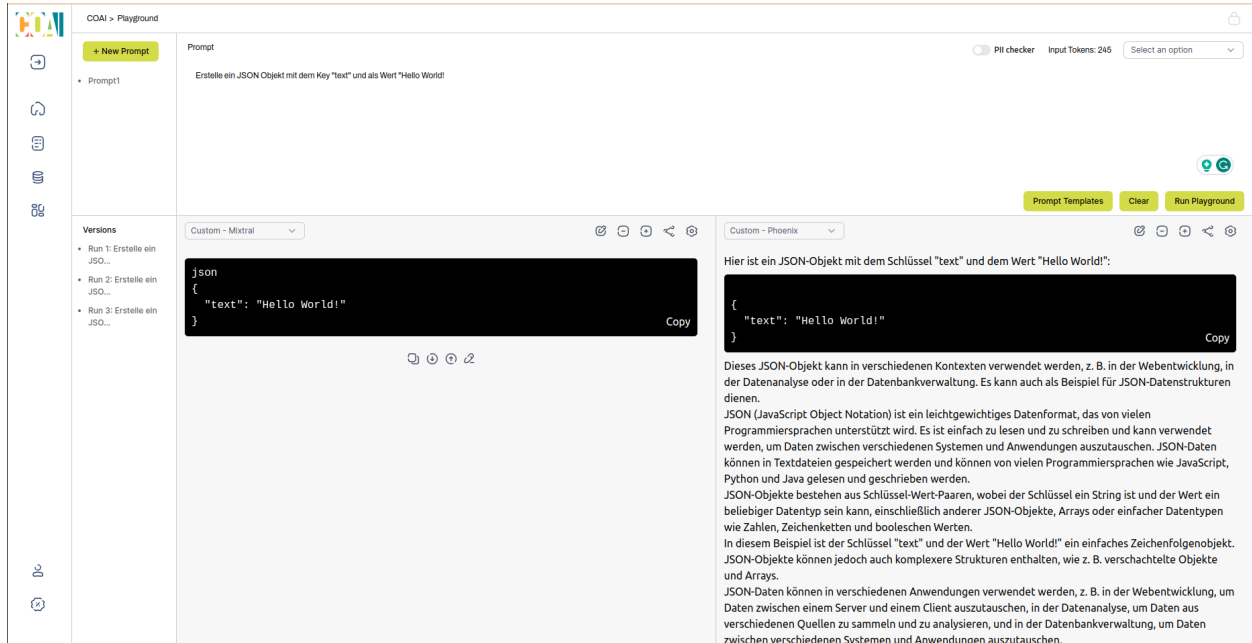


Figure 1: The COAI Playground with the outputs for the two models side by side. The Mixtral MoE model is on the left and the Phoenix model is on the right.

Table 5: Comparison of Responses from Two Different LLMs: Mixtral MoE and the Phoenix

Prompt	Mixtral MoE	Phoenix
Erstelle ein JSON Objekt mit dem Key "text" und als Wert "Hello World!"	Generates only the JSON object	Generates the JSON object, and talks about the use of JSON objects, the concept of Key and Value.
Wie hoch muss man einen Fisch an einer Angel hängen, damit eine Katze ihn nicht erreichen kann?	Talks about how the fishing line should be kept at a height. Talks about the distance of 1.5 to 2 meters over the floor.	Talks about how it is not possible to keep the fish safe from a cat, as the cats are skilled climbers. Talks about the life of the fish and its safety.