

EFFICIENT LANGUAGE MODEL TRAINING THROUGH CROSS-LINGUAL AND PROGRESSIVE TRANSFER LEARNING

Malte Ostendorff & Georg Rehm

German Research Center for Artificial Intelligence (DFKI GmbH)

Berlin, Germany

{first.lastname}@dfki.de

ABSTRACT

Most language models are primarily pretrained on English text, limiting their use for other languages. As the model sizes grow, the performance gap between English and other languages with fewer compute and data resources increases even further. Consequently, more resource-efficient training methods are needed. To address this problem, we introduce a cross-lingual and progressive transfer learning approach, called CLP-Transfer, that transfers models from a source language, for which pretrained models are available to a new target language. As opposed to prior work, which focused on the cross-lingual transfer, we extend the transfer to the model size. Given a pretrained model in a source language, we aim for a same-sized model in a target language. Instead of training a model from scratch, we exploit a smaller model that is in the target language but requires much fewer resources. Both small and source models are then used to initialize the token embeddings of the larger model based on the overlapping vocabulary of the source and target language. All remaining weights are reused from the model in the source language. This approach outperforms the sole cross-lingual transfer and can save up to 80% of the training steps compared to the random initialization.

1 INTRODUCTION

Language models (LMs) based on the Transformer architecture (Vaswani et al., 2017) dominate today’s NLP. These models are typically pretrained on primarily English text (Zhang et al., 2022; Black et al., 2022), except for a few multilingual models (Scao et al., 2022; Lin et al., 2021; Shliazhko et al., 2022). Given that multilingual models have been shown to perform suboptimal compared to monolingual ones (Conneau et al., 2020; Nozza et al., 2020), other languages than English benefit less from the recent progress in NLP. As the model sizes grow, the performance gap between the models for English and other languages with fewer resources increases even further. This gap is emphasized by Hoffmann et al. (2022), as they show that model performance is not only bound by computing resources but mainly by data. Consequently, more resource-efficient training methods are needed to bridge the gap for languages with fewer resources available.

Transfer learning is generally known to improve the training efficiency of various machine learning problems (Houlsby et al., 2019). Regarding LMs, efficient methods for task, language, or domain adaption have been proposed (Pfeiffer et al., 2020; Guo et al., 2022). To obtain monolingual LMs for low-resource languages, Minixhofer et al. (2022) and de Vries & Nissim (2021) have shown that available pretrained models can be recycled. These cross-lingual transfer learning approaches reduce the training effort. However, they only transfer across languages and neglect the sizes of the LMs. While training a large model may not be feasible in a low-resource setting, training a small or medium model is likely possible, as demonstrated by AraGPT2 (Antoun et al., 2021), CamemBERT (Martin et al., 2020), or Finnish BERT (Virtanen et al., 2019).

This paper presents CLP-Transfer, which is a cross-lingual and progressive transfer learning approach. As opposed to prior work, which focused on the cross-lingual transfer between two languages, we extend the transfer to the dimension of the model size. Given a large and pretrained model in a source

language, we aim for a same-sized model in a target language. Instead of training the large model in the target language from scratch, we first train a smaller model that requires much fewer resources (or reuse public models). Both small and source models are then used to initialize the token embeddings of the large target model based on the overlapping vocabulary of the source and target language. All remaining Transformer weights are reused from the large model in the source language.

We evaluate CLP-Transfer for decoder-only LMs based on GPT2 (Radford et al., 2019) and BLOOM (Scao et al., 2022). We use Ukrainian or German as target languages. The source models are either in English or multilingual. We find that CLP-Transfer outperforms the sole cross-lingual transfer and can save up to 80% of the training steps compared to the random initialization.

2 RELATED WORK

Cross-lingual Transfer. Exploiting pretrained models or data across languages is a common approach in NLP research (Zoph et al., 2016; Lin et al., 2019; Nguyen & Chiang, 2017). For instance, Artetxe et al. (2020) proposed to replace the tokenizer and only train the token embeddings while freezing other Transformer layers of a multilingual BERT model. Such a transfer approach produces monolingual models that can be independently fine-tuned to specific languages. de Vries & Nissim (2021) followed a similar approach to transfer a GPT2 model to a new language. Specifically, they transfer English GPT2 to Dutch and Italian by exclusively relearning the token embeddings and not the other model weights. This forces the LM to learn token embeddings that are aligned between English and the target language. However, freezing most parameters also limits the model’s ability to learn about the new language. More recently, Minixhofer et al. (2022) introduced the WECHSEL method that uses bilingual dictionaries to map the token embeddings from the source to the target language. It reuses the Transformer weights from the source model and continues training them.

Progressive Transfer. Going from a small to a larger model is also known as progressive growing and was originally proposed to improve training stability. Simonyan & Zisserman (2014) showed that starting from an efficient and small model and gradually increasing the model capacity yields more stable training. The paradigm of progressive growth can also be used to accelerate model training which has been shown for various model architectures. Karras et al. (2017) demonstrate this for GANs, Graves (2016) for RNNs, and Gu et al. (2021) for BERT models. Furthermore, Gong et al. (2019) grow a BERT model in terms of its depth, i.e., they use trained weights of a shallow model to initialize a deeper model and achieve 25% shorter training time.

3 METHODOLOGY

Our objective is to obtain a large LMs $M_t^{(\text{large})}$ with $p^{(\text{large})}$ parameters for a target language t . To increase the efficiency, we omit the standard from-scratch training approach, i.e., random initialization of $M_t^{(\text{large})}$. Instead, our goal is to find a good initialization of the parameter weights of $M_t^{(\text{large})}$ such that training effort is reduced. To achieve this, we exploit an already pretrained large LM $M_s^{(\text{large})}$, also with $p^{(\text{large})}$ parameters and the same model architecture but in a source language s , and a small pretrained LM $M_t^{(\text{small})}$, with significantly fewer parameters $p^{(\text{small})} \ll p^{(\text{large})}$ in the target language t . The Transformer layer weights \mathbf{W}_t from the large target model are initialized with the weights of $M_s^{(\text{large})}$. Similarly, token embedding weights V_t for that the tokens that exist in both the target and source language vocabulary are initialized with V_s . For the remaining token embeddings weights, a combination of $M_s^{(\text{large})}$ and $M_t^{(\text{small})}$ is used. To get our approach to work, we rely on two assumptions about the vocabularies and token embedding spaces of the source and target LMs.

3.1 ASSUMPTIONS

Shared vocabulary. Our approach relies on the tokenizers of source and target languages sharing a substantial fraction of their vocabulary. Given the tokenizer in the source and target language with their vocabularies V_s and V_t , we assume that the number of tokens occurring in both vocabularies $V_s \cap V_t$ is significantly larger than zero, i.e., $|V_s \cap V_t| \gg 0$. Languages with the same script and from the same language family typically share more tokens, e.g., the overlap between German and English is higher compared to Ukrainian and English. Notably, there will be always a certain overlap since

	125M	350M	1.3B	2.7B	6.7B	13B
125M	1.00	0.38	0.58	0.56	0.52	0.54
350M	0.38	1.00	0.37	0.37	0.35	0.36
1.3B	0.58	0.37	1.00	0.69	0.66	0.65
2.7B	0.56	0.37	0.69	1.00	0.70	0.74
6.7B	0.52	0.35	0.66	0.70	1.00	0.73
13B	0.54	0.36	0.65	0.74	0.73	1.00

Figure 1: Similarity of OPT’s token embeddings measured as overlapping k -nearest neighbors.

Vocabulary s	Vocabulary t	$ V_s \cap V_t $
English GPT2	German (ours)	24.04%
Multil. BLOOM	German (ours)	5.55%
Multil. XGLM	German (ours)	2.62%
English GPT2	Ukrainian (ours)	7.51%
English GPT2	Arabic GPT2	6.95%
English GPT2	Finnish GPT2	13.71%
Multil. BLOOM	Arabic GPT2	6.52%
Multil. BLOOM	Finnish GPT2	3.54%

Table 1: Normalized number of overlapping vocabulary tokens between different tokenizers.

Byte-Pair Encoding (Sennrich et al., 2016) is the tokenization algorithm. As shown in Tab. 1, the assumption holds for our source and target combinations. While the English and German tokenizers share 24% of their vocabulary, the multilingual BLOOM also shares 5% of the German vocabulary despite its much larger vocabulary size. Ukrainian and English have 7% overlapping tokens.

Token embeddings. An LM has the token embeddings $V \in \mathbb{R}^{|V| \times h}$ that map each token v in the vocabulary V to its vector representation $v \in \mathbb{R}^h$ with the hidden size of h . For larger models, the hidden size h of the token embedding is typically also larger compared to one of smaller models, i.e., $h^{(large)} > h^{(small)}$. Despite varying in terms of h , we assume that relative positioning in the token embedding space remains comparable across model sizes. The embeddings of a small model $V^{(small)}$ would share spacial properties with the embeddings $V^{(large)}$ of a large model.

To test this assumption, we compare token embeddings with different sizes from English OPT models (Zhang et al., 2022). Specifically, we retrieve the set of k -nearest neighbours N_v with $k = 10$ for each token v and measure the overlapping neighbors for different model sizes, e.g., $N_v^{(large)} \cap N_v^{(small)}$. This measure is normalized and computed for all available tokens. As shown in Fig. 1, OPT token embeddings are comparable across model sizes. The similarity between embedding spaces increases when the model size is comparable. We find even between the smallest and the largest model (125M and 13B parameters) a 54% overlap. It is unclear why the 350M model has the lowest embedding similarity compared to all other models, independent of their size difference.

3.2 CROSS-LINGUAL & PROGRESSIVE TRANSFER

The weights of a model in a language i are comprised of token embeddings V_i and the Transformer weights W_i . We want to initialize $V_t^{(large)}$ and $W_t^{(large)}$ for our target language t and the large model size. The Transformer weights are simply initialized with the ones from the source language s , i.e., $W_t^{(large)} = W_s^{(large)}$. To initialize $V_t^{(large)}$, we exploit $V_s^{(large)}$ and $V_t^{(small)}$, which are the token embeddings of a smaller model in the target language. The embeddings of overlapping tokens that simultaneously exist in the source and target vocabulary $v \in V_s \cap V_t$ are directly initialized with the source embeddings: When a token is not part of the overlapping vocabulary $v \notin V_s \cap V_t$, we initialize its embedding v_t as the weighted average over the embeddings \hat{v} of the overlapping token:

$$v_t = \begin{cases} v_s, & \text{if } v \in V_s \cap V_t \\ \sum_{\hat{v} \in V_s \cap V_t} \frac{\delta(v_t^{(small)}, \hat{v}_t^{(small)})}{\delta(v_t^{(small)}, \hat{v}_t^{(small)})}, & \text{otherwise} \end{cases} \tag{1}$$

The weight δ aims to transfer the spacial properties from the small model to the large model and is the normalized cosine similarity of the small embeddings of overlapping v and missing \hat{v} tokens:

$$\delta(v, \hat{v}) = \frac{\cos(v_t^{(small)}, \hat{v}_t^{(small)})}{\sum_{\substack{\hat{v}' \in V_s \cap V_t, \\ v' \in V_s \cup V_t}} \cos(v_t^{(small)}, \hat{v}'_t^{(small)})} \tag{2}$$

The intuition is those embeddings that are more similar in the $E_t^{(small)}$ should contribute more to the construction of their corresponding token in the large model. This approach allows us to recycle the pretrained weights of a source large model while preserving the spacial properties of the embedding space of the target language and simultaneously adjusting it to the vocabulary of our target language.

4 EXPERIMENT DESIGN

We evaluate the CLP-Transfer approach by transferring the English GPT2 (Radford et al., 2019) and multilingual BLOOM (Scao et al., 2022) to monolingual Ukrainian or German models. Both model types are evaluated at different scales. More specifically, we grow the GPT2 model from 124M to 774M and 117M to 1.5B parameters and the BLOOM model from 1.5B to 6.4B parameters.

4.1 MODELS

Model Architectures. Both models (GPT2 and BLOOM) are decoder-only Transformer models (Vaswani et al., 2017) trained with causal language modeling. GPT2 uses learned positional embeddings, whereas BLOOM uses ALiBi (Press et al., 2022). Another difference is that BLOOM applies normalization on the token embedding layer to improve training stability.

In our experiments with GPT2, we aim for a monolingual Ukrainian model with 774M parameters (GPT2-Large) and a German model with 1.5B parameters (GPT2-XL). The source models are the English GPT2 models as provided by Radford et al. (2019). The small Ukrainian model is trained from scratch by us. As the small German model, we use a GPT2-base model with 117M parameters trained with WECHSEL (Minixhofer et al., 2022). Additionally, we conduct experiments with BLOOM. For this experiment, our objective is the training of a German model based on BLOOM with 7.1B parameters as the source model. The BLOOM 7.1B model has 30 layers, 32 attention heads, and a hidden size of 4096. Our German BLOOM target model uses a different tokenizer with a smaller vocabulary size (see below). Therefore, its token embedding layer contains fewer parameters than the multilingual BLOOM model. As a result, the target model has only 6.4B parameters. The small German model is a BLOOM model with 1.5B parameters trained with our method (24 layers, 16 attention heads, and a hidden size of 2048).

Tokenizers. All tokenizers use Byte-Pair Encoding (Sennrich et al., 2016). The vocabulary size of English GPT2 is 50,257 tokens. BLOOM covers 46 natural languages and 13 programming languages. Therefore, BLOOM’s vocabulary has 250,880 tokens, 5x larger than the one from English GPT2. For our Ukrainian and German tokenizers, we opt for the same vocabulary sizes as the English one (50,257 tokens) and train them on subsets of the respective training sets.

4.2 BASELINES

We compare against from-scratch training and WECHSEL. (1) With from-scratch training, the LM is trained from scratch in the target language with randomly initialized weights. The from-scratch baseline for the BLOOM experiments (BLOOM 6.7B) was trained with minor changes to the transferred BLOOM-CLP 6.4B model. The baseline BLOOM 6.7B follows the model size proposed by Brown et al. (2020). It has 32 layers instead of 30 layers and was not trained on GC4 and Open Legal Data but on other German datasets. (2) The WECHSEL method as introduced by Minixhofer et al. (2022) applies cross-lingual transfer to monolingual LMs. WECHSEL uses bilingual dictionaries to map the token embeddings from a source language to a target language and reuses the Transformer weights from the source model.

5 RESULTS

We show the results for one monolingual Ukrainian and two German models, i.e., GPT2 774M, GPT2 1.5B, and BLOOM 6.4B. The models are evaluated with validation perplexity w.r.t. consumed tokens.

Transferring GPT2. The first experiment evaluates CLP-Transfer by training a Ukrainian and a German GPT2 model. For German, CLP-Transfer is compared against from-scratch training and

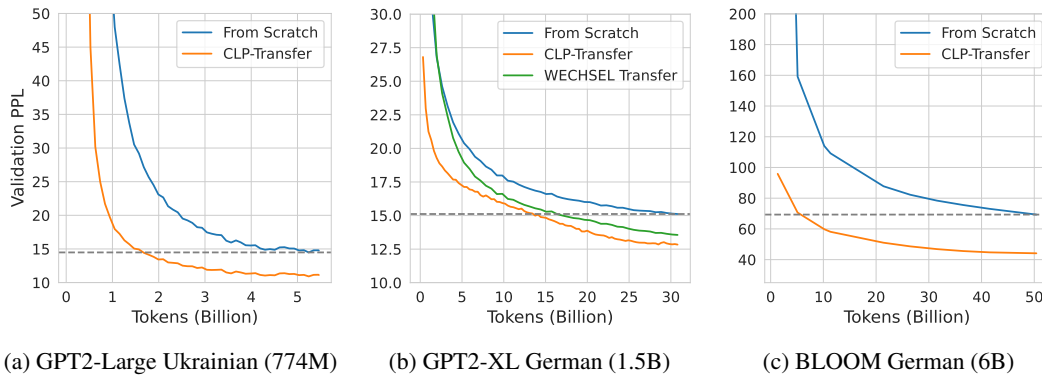


Figure 2: Validation PPL comparing CLP-Transfer, from-scratch, and WECHSEL (cross-lingual transfer). CLP-Transfer achieves the same PPL as from-scratch but already after 20% of tokens for BLOOM (32% of tokens for GPT2-Large; 50% of tokens for GPT2-XL); see the dashed line.

WECHSEL’s cross-lingual transfer method. Fig. 2b shows the validation perplexity (PPL) of each method in relation to the consumed training tokens. We find that CLP-Transfer outperforms the baselines. The validation PPL of CLP-Transfer is constantly the lowest of all three methods. At the end of the German training (after 30.8B tokens), CLP-Transfer yields a 12.8 PPL, followed by WECHSEL with 13.5 PPL. The worst result has the from-scratch training with 15.1 PPL. CLP-Transfer achieve the same PPL as from-scratch training but already have 50% of the consumed tokens. During the first phase of the training (0-5B tokens), the improvements of CLP-Transfer are most significant. These results demonstrate that our transfer learning approach is superior to from-scratch training even at the end of the training or can achieve the same results more efficiently. Moreover, using a small model in the target language yields further efficiency gains compared to WECHSEL’s sole cross-lingual transfer. Fig. 2a shows a similar outcome for the Ukrainian models, i.e., after 5.4B tokens CLP-Transfer yields 11.1 PPL compared to 14.7 PPL through from-scratch training. This demonstrates that CLP-Transfer can even transfer models across languages with different scripts (from Latin to Cyrillic).

Transferring BLOOM. The second experiment applies CLP-Transfer on a multilingual BLOOM model to train a monolingual German model with 6.4B parameters. In this experiment, we compare only against from-scratch training. As shown in Fig. 2c, CLP-Transfer again outperforms the from-scratch training. After complete training on 50.4B tokens, CLP-Transfer yields a 44.1 PPL, whereas from-scratch training is significantly worse with 69.3 PPL. 20% of training tokens are sufficient for CLP-Transfer to be on par with from-scratch training. This suggests that the efficiency gains from CLP-Transfer are even more prevalent at 6B compared to 1.5B parameters. We attribute this outcome to the training data containing too few tokens for 6B models. The validation PPL is still decreasing at the end of the training suggesting that the model is not fully trained yet. According to Hoffmann et al. (2022), a compute-optimal LM at the 6B scale would require approx. 142B tokens which our BLOOM model training did not consume. The German GPT2-XL training is much closer to being compute-optimal (33B tokens).

6 CONCLUSION

CLP-Transfer is a cross-lingual and progressive transfer learning approach for the efficient training of large LMs. Our experiments demonstrate that monolingual Ukrainian or German models initialized with CLP-Transfer reduce the training effort. The CLP-Transfer models achieve better results when trained on the same number of tokens than from-scratch training or WECHSEL transfer. To obtain the same perplexity as from-scratch training, CLP-Transfer needs only 20-50% of the original token count, depending on model type and language. This yields up to an 80% reduction in training effort. Such a reduction lowers the barriers to the training of large LMs in low-resource settings. We make code, model checkpoints, and a Web-based demo available.¹

¹Demo: <https://opengptx.dfki.de>; Repository: <https://github.com/malteos/clp-transfer>

ACKNOWLEDGMENTS

We thank the anonymous reviewers for their insightful comments and suggestions. The work presented in this paper has received funding from the German Federal Ministry for Economic Affairs and Climate Action (BMWK) through the project OpenGPT-X (project no. 68GX21007D). The authors gratefully acknowledge the Gauss Centre for Supercomputing e.V. and the GWK support for funding this project by providing computing time through the John von Neumann Institute for Computing (NIC) on the GCSSupercomputer JUWELS at Jülich Supercomputing Centre (JSC) and through the Center for Information Services and HPC (ZIH) at TU Dresden.

REFERENCES

- Wissam Antoun, Fady Baly, and Hazem Hajj. AraGPT2: Pre-trained transformer for Arabic language generation. In *Proceedings of the Sixth Arabic Natural Language Processing Workshop*, pp. 196–207, Kyiv, Ukraine (Virtual), April 2021. Association for Computational Linguistics. URL <https://aclanthology.org/2021.wanlp-1.21>.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of monolingual representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4623–4637, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.421. URL <https://aclanthology.org/2020.acl-main.421>.
- Sidney Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, Usvsn Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. GPT-NeoX-20B: An open-source autoregressive language model. In *Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models*, pp. 95–136, virtual+Dublin, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.bigscience-1.9. URL <https://aclanthology.org/2022.bigscience-1.9>.
- 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. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 2020-Decem, 2020. ISSN 10495258.
- Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. Xnli: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2018.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised Cross-lingual Representation Learning at Scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 8440–8451, Online, 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.747.
- Wietse de Vries and Malvina Nissim. As Good as New. How to Successfully Recycle English GPT-2 to Make Models for Other Languages. *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 836–846, 2021. doi: 10.18653/v1/2021.findings-acl.74.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, September 2021. URL <https://doi.org/10.5281/zenodo.5371628>.
- Linyuan Gong, Di He, Zhuohan Li, Tao Qin, Liwei Wang, and Tiejun Liu. Efficient training of BERT by progressively stacking. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2337–2346. PMLR, June 2019.

- Alex Graves. Adaptive Computation Time for Recurrent Neural Networks. 2016. doi: 10.48550/ARXIV.1603.08983.
- Xiaotao Gu, Liyuan Liu, Hongkun Yu, Jing Li, Chen Chen, and Jiawei Han. On the Transformer Growth for Progressive BERT Training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5174–5180, Online, 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.406.
- Xu Guo, Boyang Albert Li, and Han Yu. Improving the sample efficiency of prompt tuning with domain adaptation. *ArXiv*, abs/2210.02952, 2022.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training Compute-Optimal Large Language Models. 3(2020):1–36, 2022.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, 2019.
- Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive Growing of GANs for Improved Quality, Stability, and Variation. 2017. doi: 10.48550/ARXIV.1710.10196.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating Training Data Makes Language Models Better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8424–8445, Dublin, Ireland, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.577.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shrutu Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian Li. Few-shot Learning with Multilingual Language Models. 2021.
- Yu-Hsiang Lin, Chian-Yu Chen, Jean Lee, Zirui Li, Yuyan Zhang, Mengzhou Xia, Shrutu Rijhwani, Junxian He, Zhisong Zhang, Xuezhe Ma, Antonios Anastasopoulos, Patrick Littell, and Graham Neubig. Choosing transfer languages for cross-lingual learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3125–3135, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1301. URL <https://aclanthology.org/P19-1301>.
- Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. CamemBERT: a tasty French language model. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7203–7219, Online, July 2020. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.acl-main.645>.
- Benjamin Minixhofer, Fabian Paischer, and Navid 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*, pp. 3992–4006, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.293. URL <https://aclanthology.org/2022.naacl-main.293>.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M. Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. Crosslingual Generalization through Multitask Finetuning, November 2022.

- Toan Q. Nguyen and David Chiang. Transfer learning across low-resource, related languages for neural machine translation. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pp. 296–301, Taipei, Taiwan, November 2017. Asian Federation of Natural Language Processing. URL <https://aclanthology.org/I17-2050>.
- Debora Nozza, Federico Bianchi, and Dirk Hovy. What the [MASK]? Making Sense of Language-Specific BERT Models. 2020. doi: 10.48550/ARXIV.2003.02912.
- Malte Ostendorff, Till Blume, and Saskia Ostendorff. Towards an Open Platform for Legal Information. In *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020*, pp. 385–388, August 2020. doi: 10.1145/3383583.3398616.
- Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. AdapterHub: A Framework for Adapting Transformers. pp. 46–54, 2020. doi: 10.18653/v1/2020.emnlp-demos.7.
- Ofir Press, Noah A. Smith, and Mike Lewis. Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation, April 2022.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language Models are Unsupervised Multitask Learners. *arXiv*, 2019.
- Georg Rehm (ed.). *European Language Grid: A Language Technology Platform for Multilingual Europe*. Cognitive Technologies. Springer, Cham, Switzerland, January 2023.
- Georg Rehm and Andy Way (eds.). *European Language Equality: A Strategic Agenda for Digital Language Equality*. Cognitive Technologies. Springer, June 2023. In print.
- Teven Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, M Saiful Bari, Stella Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, Ofir Press, Colin Raffel, Victor Sanh, Sheng Shen, Lintang Sutawika, Jaesung Tae, Zheng Xin Yong, Julien Launay, and Iz Beltagy. What Language Model to Train if You Have One Million GPU Hours? 2022. doi: 10.48550/ARXIV.2210.15424.
- Dietmar Schabus, Marcin Skowron, and Martin Trapp. One million posts: A data set of german online discussions. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pp. 1241–1244, Tokyo, Japan, August 2017. doi: 10.1145/3077136.3080711.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1162. URL <https://aclanthology.org/P16-1162>.
- Oleh Shliashko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. mGPT: Few-Shot Learners Go Multilingual, April 2022.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. Asynchronous pipelines for processing huge corpora on medium to low resource infrastructures. *Proceedings of the Workshop on Challenges in the Management of Large Corpora (CMLC-7) 2019*. Cardiff, 22nd July 2019, pp. 9 – 16, Mannheim, 2019. Leibniz-Institut für Deutsche Sprache. doi: 10.14618/ids-pub-9021. URL <https://nbn-resolving.org/urn:nbn:de:bsz:mh39-90215>.
- Jannis Vamvas and Rico Sennrich. X-Stance: A multilingual multi-target dataset for stance detection. In *Proceedings of the 5th Swiss Text Analytics Conference (SwissText) 16th Conference on Natural Language Processing (KONVENS)*, Zurich, Switzerland, jun 2020. URL <http://ceur-ws.org/Vol-2624/paper9.pdf>.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. *Proceedings of the 31st International Conference on Neural Information Processing Systems*, (Nips):6000–6010, June 2017. doi: 10.1017/CBO9780511809071.
- Antti Virtanen, Jenna Kanerva, Rami Ilo, Jouni Luoma, Juhani Luotolahti, Tapio Salakoski, Filip Ginter, and Sampo Pyysalo. Multilingual is not enough: Bert for finnish, 2019. URL <https://arxiv.org/abs/1912.07076>.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=gEZrGCozdqR>.
- Michael Wiegand and Melanie Siegel. Overview of the germeval 2018 shared task on the identification of offensive language. In *Proceedings of the GermEval*, 2018.
- Michael Wojatzki, Eugen Ruppert, Sarah Holschneider, Torsten Zesch, and Chris Biemann. Germeval 2017: Shared task on aspect-based sentiment in social media customer feedback. *Proceedings of the GermEval*, pp. 1–12, 2017.
- Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. PAWS-X: A cross-lingual adversarial dataset for paraphrase identification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3687–3692, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1382. URL <https://aclanthology.org/D19-1382>.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona 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. OPT: Open Pre-trained Transformer Language Models. 2022.
- Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. Transfer learning for low-resource neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 1568–1575, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1163. URL <https://aclanthology.org/D16-1163>.

A DATASETS

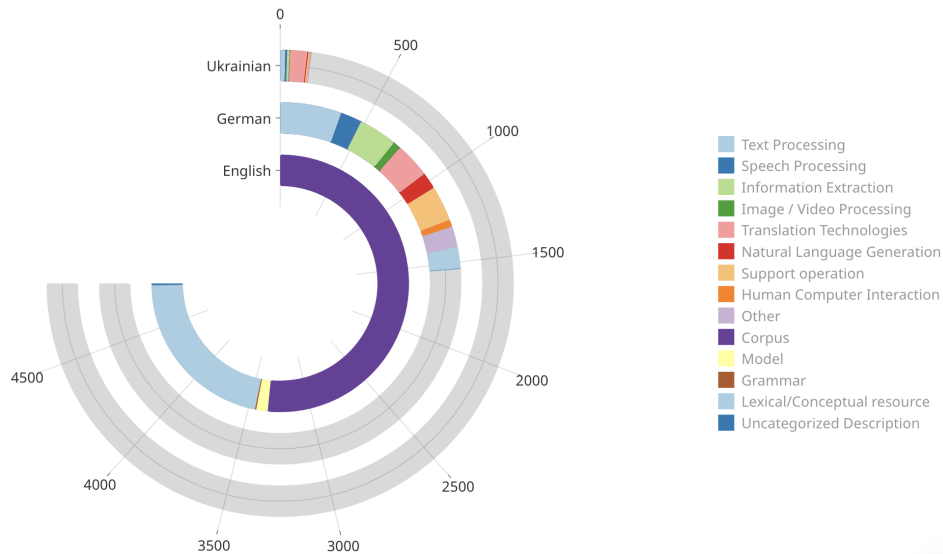


Figure 3: Comparison between the available language resources for English, Ukrainian, and German.

Depending on language and model type different training datasets are used. The datasets represent the Ukrainian and German data that is easily accessible. The small size of the datasets illustrates the lack of resources for Ukrainian and German. Fig. 3 shows the general gaps between these three languages w.r.t. the available resources in the European Language Grid (Rehm, 2023) as measured by the dashboard² implemented by the European Language Equality project (Rehm & Way, 2023).

Ukrainian GPT2 Training. The Ukrainian GPT2-CLP training relies on Web-crawled data from the Ukrainian subset of OSCAR v2109, v2201, and v2301 (Suárez et al., 2019). Wikimedia dumps are also used (date: 2023-03-01; Wikipedia, Wikibooks, Wikinews, Wikiquote, Wikisource, and Wikivoyage). We split the data into a 99:1 train-validation set. The Ukrainian training dataset comprises approximately 2.7B tokens. The data was used twice (two epochs).

German GPT2 Training. The German GPT2-CLP training relies exclusively on Web-crawled data from the German subset of OSCAR v2019 (Suárez et al., 2019).³ We follow the methodology from Minixhofer et al. (2022) to construct a separate training and validation dataset. Specifically, we used the first 4 GB of OSCAR as the training dataset, then the next 0.4GB as the validation dataset. The GPT2 training dataset comprises approximately 30.8B tokens.

BLOOM Training. To train the German BLOOM-CLP 6.4B model, we construct another dataset. We use again Web-crawled content from the German subset OSCAR but the more recent version of v2201 (excluding content tagged as *header*, *footer*, *noisy*, or *adult*) and from the GC4 Corpus⁴ (including only the *head* and *middle* parts). As both data sources originate from CommonCrawl and potentially have duplicated content, we deduplicate the Web-crawled content using the approach from Lee et al. (2022). We complement the Web data with German court decisions from Open Legal Data (Ostendorff et al., 2020). The BLOOM training dataset comprises approximately 50.4B tokens.

Evaluation Datasets. We evaluate the German models for their language modeling ability using the OSCAR validation set from the GPT2 training⁵, and for zero-shot learning on German downstream tasks. The tasks are sentiment analysis from GermEval 2017 (Wojatzki et al., 2017), hate

²<https://live.european-language-grid.eu/catalogue/dashboard>

³[https://hf.co/datasets/oscar\(subset:unshuffled_deduplicated_de\)](https://hf.co/datasets/oscar(subset:unshuffled_deduplicated_de))

⁴<https://german-nlp-group.github.io/projects/gc4-corpus.html>

⁵https://hf.co/datasets/malteos/wechsel_de

speech classification from GermEval 2018 (Wiegand & Siegel, 2018), news topic classification from GNAD10 (Schabus et al., 2017), paraphrase identification from PAWSX (Yang et al., 2019), natural language inference from XNLI (Conneau et al., 2018), and stance detection from X-Stance (Vamvas & Sennrich, 2020). All evaluation tasks are implemented using the `lm-evaluation-harness` framework (Gao et al., 2021).⁶

Task → Model ↓ / Metric →	Oscar PPL (↓)	GEval17 F1 (↑)	GEval18 F1 (↑)	GNAD10 F1 (↑)	PAWSX F1 (↑)	XNLI Acc. (↑)	XStance F1 (↑)	Avg. (↑)
Random	-	0.33	0.50	0.11	0.50	0.33	0.50	0.38
<i>Multilingual models:</i>								
mGPT 1.3B	2274.80	0.36	0.51	0.08	0.49	0.37	0.49	0.38
XGLM 564M	179.59	0.05	0.40	0.05	0.46	0.44	0.50	0.32
XGLM 1.7B	105.10	0.04	0.35	0.19	0.58	0.45	0.40	0.34
XGLM 7.5B	66.74	0.51	0.51	0.06	0.50	0.39	0.41	0.40
<i>Monolingual German models:</i>								
GPT2-WECHSEL 117M	594.40	0.04	0.51	0.18	0.49	0.40	0.51	0.35
GPT2-XL-WECHSEL 1.5B	157.95	0.05	0.55	0.10	0.41	0.49	0.34	0.32
GPT2-XL-CLP 1.5B	46.33	0.05	0.02	0.07	0.46	0.49	0.34	0.24
GPT2-XL 1.5B f-s.	187.71	0.04	0.51	0.15	0.52	0.47	0.34	0.34
BLOOM-CLP 1.5B	49.80	0.04	0.14	0.11	0.44	0.48	0.38	0.26
BLOOM-CLP 6.4B (50B t.)	44.09	0.56	0.51	0.13	0.52	0.43	0.44	0.43
BLOOM 6.7B f-s. (50B t.)	69.32	0.51	0.52	0.13	0.41	0.38	0.42	0.39
BLOOM 6.7B f-s. (72B t.)	64.03	0.56	0.51	0.09	0.40	0.37	0.49	0.40

Table 2: Evaluation results of German downstream tasks in a zero-shot setting. The average score excludes the OSCAR validation perplexity (PPL). Smaller models are on par or worse than the random baseline. Our transfer model BLOOM-CLP 6.4B achieves the best results on average.

B DOWNSTREAM EVALUATION

Even though we trained the models exclusively with a causal language modeling objective, we want them to perform well on other downstream tasks, as shown by Brown et al. (2020). Hence, we compare the models and additional baselines on six German benchmarks in a zero-shot setting.

We compare the monolingual German models against multilingual models trained on German data. We evaluate XGLM (Lin et al., 2021) ranging from 564M to 7.5B parameters and mGPT (Shliakhko et al., 2022) with 1.3B parameters. XGLM was trained on approx. 5.4% German data and mGPT on 8.2% German data.

Given that the from-scratch trained BLOOM 6.7B model (50B tokens) is presumably under-trained, we add an additional variation that was trained on 22B more tokens, i.e., BLOOM 6.7B (72B tokens). The evaluated tasks are sentiment analysis (GermEval 2017), hate speech classification (GermEval 2018), news topic classification (GNAD10), paraphrase identification (PAWSX), natural language inference (XNLI), and stance detection (X-Stance). Tab. 2 reports the validation PPL on German OSCAR⁵, the results for the individual tasks, and the average over the tasks.

The zero-shot performance of all models is disappointing. Most models achieve results on par or worse than the random baseline. Only the largest models (more than 6B parameters) are better than the random baseline on average. The BLOOM-CLP 6.4B model has the best average score of 0.43, followed by the from-scratch trained BLOOM 6.7B (72B tokens) and XGLM 7.5B.

We hypothesize that this outcome is due to the model size and token count being still too small. Studies from Black et al. (2022) or Shliakhko et al. (2022) report similar near-random results for models with comparable sizes. Another reason might be the poorly translated test datasets that produce less meaningful results. For instance, PAWSX contains a large fraction of machine-translated samples. To improve the downstream task performance, promising approaches are prompt engineering, i.e., tailoring the prompts more to the German language, and multi-task fine-tuning, as demonstrated by BLOOMZ (Muennighoff et al., 2022) or FLAN (Wei et al., 2022).

⁶<https://github.com/OpenGPTX/lm-evaluation-harness>