

ReTok: Replacing Tokenizer to Enhance Representation Efficiency in Large Language Model

Shuhao Gu, Mengdi Zhao, Bowen Zhang, Liangdong Wang,
Jijie Li, Guang Liu*

Data Research Team, BAAI
{shgu, liuguang}@baai.ac.cn

Abstract

Tokenizer is an essential component for large language models (LLMs), and a tokenizer with a high compression rate can improve the model’s representation and processing efficiency. However, the tokenizer cannot ensure high compression rate in all scenarios, and an increase in the average input and output lengths will increase the training and inference costs of the model. Therefore, it is crucial to find ways to improve the model’s efficiency with minimal cost while maintaining the model’s performance. In this work, we propose a method to improve model representation and processing efficiency by replacing the tokenizers of LLMs. We propose replacing and reinitializing the parameters of the model’s input and output layers with the parameters of the original model, and training these parameters while keeping other parameters fixed. We conducted experiments on different LLMs, and the results show that our method can maintain the performance of the model after replacing the tokenizer, while significantly improving the decoding speed for long texts.

1 Introduction

Tokenizer is a basic component of large language models (LLMs) (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023), which is used to preprocess text by converting it into a sequence of tokens. This process allows the model to handle and analyze text data more efficiently by breaking it down into manageable pieces. Currently, different LLMs typically use different tokenizers, which are generally trained on their own training corpora with various methods such as BPE (Sennrich et al., 2016), Word-Piece (Song et al., 2021), Unigram (Kudo, 2018), etc. Generally, these methods determine the vocabulary of the tokenizer based on the frequency of different tokens in the training data. This allows the

model to obtain an efficient representation of the training data, achieving a relatively high compression rate, i.e., shorter sequence lengths. However, this also results in a lower compression rate for the tokenizer when there is a significant distribution difference between the test data and the training data (Tran, 2020; Liu et al., 2023; Dobler and de Melo, 2023). For example, a tokenizer trained on general domain corpora will have a lower compression rate for specific tasks like code and math. On the one hand, using a low-compression-rate tokenizer will cause the model to output longer sequences during inference, which increases the computational cost and inference time (Ahia et al., 2023). On the other hand, it will also cause the model to consume more computation during supervised fine-tuning (SFT) on specific languages or tasks, affecting training efficiency and potentially impacting the final performance (Dagan et al., 2024). Besides, if we need to utilize different LLMs in practical applications, e.g., using different LLMs for collaboration (Xiong et al., 2023), we must optimize and adapt to different tokenizers, which increases development and maintenance costs. Therefore, substituting pretrained model’s original tokenizer with a arbitrary one that provides higher compression rates could substantially reduce the cost of using LLM.

In this work, we propose to **Replace Toknizer (ReTok)** of LLM to improve the model’s representation and processing efficiency, while maintaining the model’s performance as much as possible. First, using the Llama3 tokenizer (Meta-AI, 2024) as an example, which demonstrates high compression rates for English and code, we quickly improve its compression rate for Chinese by expanding the vocabulary. Then, we propose replacing the model’s input and output layers based on the new tokenizer. Meanwhile, we initialize the new parameters according to the correspondence between the tokens in the new tokenizer’s vocabulary and

** corresponding author

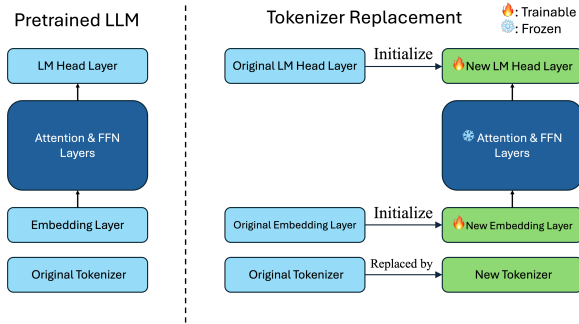


Figure 1: Illustration of the proposed method.

the original tokenizer’s vocabulary. Lastly, these new parameters of the model will be updated, and other parameters will be fixed during model training. Experiments were conducted on the Qwen1.5-0.5B (Bai et al., 2023), Aquila2-7B (BAAI, 2023), and Llama3-8B (Meta-AI, 2024) models. The results indicate that our method can maintain the performance of the original LLM after replacing the tokenizer, while significantly reducing the inference time of the model in specific domains.

2 Related Work

The work on tokenizers can generally be divided into two categories based on the objectives. The first category of work studies how to design the optimal tokenizer to help the model learn more effectively. Singh and Strouse (2024) explored the impact of different tokenization methods for numbers on mathematical problems. Clark et al. (2022) explored having the model discard the tokenizer to process text, using characters to directly represent the input sequence. Zheng et al. (2023) proposed the Syntax-Aware Tokenization method, enhancing the model’s capability for code tasks. Rajaraman et al. (2024) presented a theoretical framework to compare and analyze different tokenization algorithms, and provided several suggestions on how to improve tokenization effectiveness. Besides, there are also studies focused on improving tokenization effectiveness for some specific languages, such as Arabic (Alyafeai et al., 2023) and Japanese (Tolmachev et al., 2018).

The second category of work focuses on how to modify the existing tokenizers of the pre-trained model. Xue et al. (2022) studied how to let the pre-trained model process byte sequences with minimal modifications. Minixhofer et al. (2024) proposed to train a hypernetwork taking a tokenizer as input and predicting the corresponding embeddings. Gee

Tokenizer	Size	Compression Ratio				Average
		Zh	En	Code	Math	
Aquila2	100k	4.70	4.42	3.20	3.77	4.02
Qwen1.5	151k	4.27	4.51	3.62	3.35	3.94
Llama3	128k	3.45	4.61	3.77	3.88	3.93
ReTok	143k	4.60	4.61	3.78	3.88	4.22

Table 1: Compression rates for Chinese, English, code, and math with different tokenizers. ReTok is the tokenizer we developed based on the Llama3 tokenizer, which has a higher overall compression rate.

et al. (2024) and Dagan et al. (2024) studied how to adapt the tokenizers of pretrained models to improve efficiency and performance for downstream tasks.

In these works, our work is more similar to the second category of work, where we also improve the tokenizers based on pretrained models. Meanwhile, our work focuses more on how to improve the representation efficiency of the model with as little cost as possible while maintaining the overall performance of the model.

3 Methods

The proposed method comprises three distinct steps: obtaining a high-compression tokenizer, replacing and initializing the model’s input and output layers, and training the model. An illustration of our method is shown in Figure 1.

3.1 Vocabulary Expansion

First, we need to obtain a high-compression tokenizer, which can be obtained either by training from scratch or by extending an existing tokenizer. In this work, we exemplify with the Llama3 tokenizer, distinguished for its commendable compression capabilities in English and code, as a paradigm. By integrating Chinese lexicons into its vocabulary, we enhance its compression efficacy for Chinese text, thereby obtaining a more efficient representation. Initially, we employ the same pre-tokenization approach as the original tokenizer to segment the Chinese training corpus into distinct chunks. Subsequently, on this corpus, we utilize byte-level BPE method to learn a new vocabulary and merge it with the original vocabulary to obtain the new vocabulary as well as the tokenizer. We compared the compression rates of different tokenizers across Chinese, English, code, and math corpora, where

Model	PIQA (5-shot)	ARC-C (25-shot)	HellaSwag (10-shot)	MMLU (5-shot)	CMMLU (5-shot)	AGIEval (0-shot)	BBH (3-shot)	Human-Eval (0-shot)	Gsm8k (5-shot)	Average
Qwen1.5-0.5B	69.8	31.6	49.1	39.94	45.57	22.46	21.24	9.15	13.04	33.54
+ ReTok	70	30.6	48.4	38.74	41.32	24.38	24.83	11.59	8.64	33.17
Aquila2-7B	76.3	42.4	68.7	45.66	49.94	26.17	29.73	6.1	8.04	39.23
+ ReTok	76.3	42.8	68.4	43.72	41.45	23.84	22.73	14.02	6.44	37.74
LLama3-8b	82.1	59.6	82.1	66.66	50.66	26.37	61.25	26.22	50.04	56.11
+ ReTok	81.9	56.4	82.1	65.65	46.35	25.53	59.92	29.27	51.4	55.39

Table 2: Results of the main experiments.

the compression rate is computed as:

$$r = \frac{\text{byte}(s)}{\text{tokenize}(s)}, \quad (1)$$

where $\text{byte}(s)$ represents the sequence length when a sentence is represented with bytes, while $\text{tokenize}(s)$ represents the sequence length after tokenization. The results are shown in Table 1, which indicate that the our newly expanded tokenizer (ReTok) achieves the highest overall compression rate.

3.2 Layer Replacement & Initialization

In this step, we replace and reinitialize the input layer (i.e., embedding layer) and the output layer (i.e., LM head layer) of the original model. To maximize the utilization of the original model’s knowledge and accelerate the convergence of the new parameters, we initialize the newly added parameters based on the model’s original layers. In detail, for each token in the new vocabulary, if it is also present in the original vocabulary, we initialize its corresponding new parameters in both the embedding layer and the LM head layer with the parameters from the original model’s layers. For tokens that do not exist in the original vocabulary, we process these tokens using the original tokenizer to obtain their corresponding ID lists from the original vocabulary. We then retrieve the associated parameters from the original embedding layer and LM head layer, compute their average, and use this average to initialize the new parameters.

3.3 Model Training

Then, we choose to train the reinitialized embedding layer and LM head layer of the model, while keeping the main parameters of the model, i.e. attention layers, FFN layers, and layer norm parameters, fixed to prevent catastrophic forgetting. Given the training data $\mathcal{T} = \{t_0, \dots, t_n\}$, we use the standard language modeling objective to training the

newly added parameters:

$$\mathcal{L}(\mathcal{T}) = \sum_i \log P(t_i | t_{i-1}, \dots, t_0; \Theta), \quad (2)$$

where the conditional probability P is modeled using a neural network with parameters Θ . Once the model has converged, we can unfreeze the main parameters of the model and conduct joint training on all the parameters, which will further improve the model’s performance.

4 Experiments

4.1 Experimental Setups

Data To train the reinitialized parameters effectively, it is essential to construct a comprehensive training dataset. Ideally, using the original data of the pretrained model would be the best way to restore the model’s performance. However, in practice, it is often impossible to access the entire pretraining dataset of the pretrained model. Therefore, in this work, we use publicly available datasets as much as possible for our experiments. Our collected training dataset is divided into three parts: Chinese, English, and Code. The English data primarily comes from sources such as RedPajama (Together Computer, 2023), Pile (Gao et al., 2021a), C4 (Raffel et al., 2020), and Falcon (Penedo et al., 2023). The Chinese data is mainly sourced from Wudao (Yuan et al., 2021) and ChineseWebText (Chen et al., 2023). The code data is derived from sources like StarCoder (Li et al., 2023b).

Training Details We conducted experiments based on Qwen1.5-0.5B (Bai et al., 2023), Aquila2-7B (BAAI, 2023), and Llama3-8B (Meta-AI, 2024) models. The global batch size for training was set to 384K tokens. We used a cosine learning rate scheduler with a maximum learning rate of 1e-4 and a warm-up period of 1000 steps. The training was conducted on 32 NVIDIA A-100 GPUs.

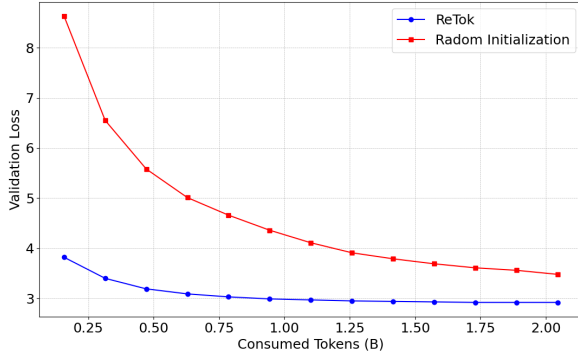


Figure 2: Comparison of validation loss changes between ReTok and random initialization based on Llama3-8B.

Evaluation The test tasks can be divided into two categories: comprehension tasks and generative tasks. The comprehension tasks include PIQA (5-shot) (Bisk et al., 2020), ARC-Challenge (25-shot) (Clark et al., 2018), HellaSwag (10-shot) (Zellers et al., 2019), MMLU (5-shot) (Hendrycks et al., 2021), and CMMLU (5-shot) (Li et al., 2023a), while the generative tasks include AGIEval (0-shot) (Zhong et al., 2023), Big-Bench Hard (BBH) (3-shot) (Suzgun et al., 2023), Human-Eval (0-shot) (Chen et al., 2021), and Gsm8k (5-shot) (Cobbe et al., 2021). These tasks cover the capabilities in Chinese, English, code, and math. All tasks were tested using frameworks lm-evaluation-harness (Gao et al., 2021b) and OpenCompass (OpenCompass Team, 2023). All the results were obtained from a single run.

4.2 Main Results

The overall results are shown in Table 2. On most test sets, our method ReTok achieved performance similar to the original models, demonstrating the effectiveness of our approach. However, there were notable changes in performance for certain test sets, such as HumanEval and CMMLU. This may be because that our training data included a significant amount of code content and relatively less Chinese text. This phenomenon suggests that the distribution of pre-training data can influence model performance, and our future work will focus on the impact of different data on model performance.

4.3 Convergence Speed Analysis

In this section, we analyzed the convergence speed of our method. We randomly selected a portion of data from the training set and, after a rigorous deduplication, used it as our validation set. We

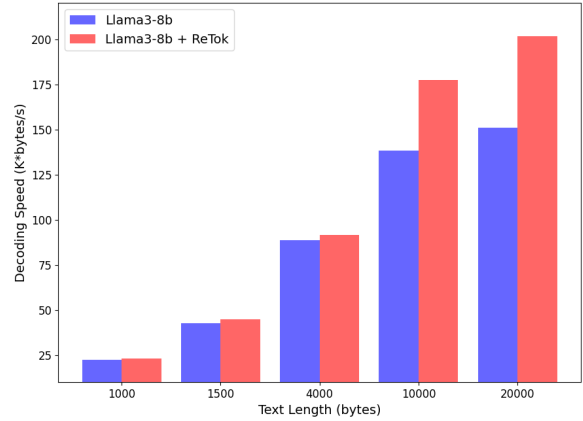


Figure 3: Comparison of decoding speeds based on the Llama3-8B model on the Chinese validation set for texts of different lengths.

tested the change of validation loss based on the Llama3-8B model with our method during training. Besides, we compared the validation loss changes when the input and output layers of the model were randomly initialized. The results are presented in Figure 2. The results show that, compared to the commonly used random initialization methods, our method quickly achieves convergence and delivers better performance.

4.4 Inference Speed Comparison

In this experiment, we compared the change in decoding speed of the model before and after tokenizer replacement based on the Llama3-8B model. We measured the decoding speed of the model on a Chinese validation set, because our proposed method ReTok primarily improves the compression rate of Chinese. We grouped the Chinese data based on byte-based sequence lengths and measured the model’s speed at each length group. Then, we used force decoding to process these data, ensuring that the lengths of input and output texts for different models were consistent (but the sequence lengths after tokenization were not the same). We compared the decoding speeds of different models, and the results are shown in Figure 3. The results indicate that our method can improve decoding speed, especially for longer texts.

5 Conclusion

In this work, we explored how to improve the model’s representation and processing efficiency by replacing the tokenizer of a LLM. We proposed using the original parameters to initialize the model’s input and output layers to accelerate convergence.

We conducted experiments on multiple test sets, and the results demonstrated that our method can maintain the original performance of the model while increasing decoding speed significantly.

Limitations

This work has two main limitations. First, the data we used comes from multiple sources, and we did not analyze the impact of each data source and different data ratios on different models. Second, our test set mainly focuses on the model’s capabilities in Chinese, English, code, and math, without testing the model’s other capabilities, such as multilingual abilities. We will leave these in future work.

References

- Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo Kasai, David R. Mortensen, Noah A. Smith, and Yulia Tsvetkov. 2023. [Do all languages cost the same? tokenization in the era of commercial language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 9904–9923. Association for Computational Linguistics.
- Zaid Alyafeai, Maged Saeed AlShaibani, Mustafa Ghaleb, and Irfan Ahmad. 2023. [Evaluating various tokenizers for arabic text classification](#). *Neural Process. Lett.*, 55(3):2911–2933.
- BAAI. 2023. Aquila2: Flagai’s open source language model. <https://github.com/FlagAI-Open/Aquila2>.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. [Qwen technical report](#). *ArXiv preprint*, abs/2309.16609.
- Yonatan Bisk, Rowan Zellers, Ronan LeBras, Jianfeng Gao, and Yejin Choi. 2020. [PIQA: reasoning about physical commonsense in natural language](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7432–7439. AAAI Press.
- 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*.
- Jianghao Chen, Pu Jian, Tengxiao Xi, Yidong Yi, Qianlong Du, Chenglin Ding, Guibo Zhu, Chengqing Zong, Jinqiao Wang, and Jiajun Zhang. 2023. [Chinesewebtext: Large-scale high-quality chinese web text extracted with effective evaluation model](#). *CoRR*, abs/2311.01149.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. [Evaluating large language models trained on code](#). *ArXiv preprint*, abs/2107.03374.
- Jonathan H. Clark, Dan Garrette, Iulia Turc, and John Wieting. 2022. [Canine: Pre-training an efficient tokenization-free encoder for language representation](#). *Trans. Assoc. Comput. Linguistics*, 10:73–91.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#). *ArXiv preprint*, abs/1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *CoRR*, abs/2110.14168.
- Gautier Dagan, Gabriel Synnaeve, and Baptiste Rozière. 2024. [Getting the most out of your tokenizer for pre-training and domain adaptation](#). *CoRR*, abs/2402.01035.
- Konstantin Dobler and Gerard de Melo. 2023. [FOCUS: effective embedding initialization for monolingual specialization of multilingual models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 13440–13454. Association for Computational Linguistics.

- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2021a. [The Pile: An 800GB dataset of diverse text for language modeling](#). *ArXiv preprint*, abs/2101.00027.
- 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. 2021b. [A framework for few-shot language model evaluation](#).
- Leonidas Gee, Andrea Zugarini, Leonardo Rigutini, and Paolo Torrioni. 2024. [Fast vocabulary transfer for language model compression](#). *CoRR*, abs/2402.09977.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Taku Kudo. 2018. [Subword regularization: Improving neural network translation models with multiple subword candidates](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 66–75. Association for Computational Linguistics.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023a. [Cmmllu: Measuring massive multitask language understanding in chinese](#). *ArXiv preprint*, abs/2306.09212.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy V, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Moustafa-Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 2023b. [Starcoder: may the source be with you!](#) *CoRR*, abs/2305.06161.
- Yihong Liu, Peiqin Lin, Mingyang Wang, and Hinrich Schütze. 2023. [OFA: A framework of initializing unseen subword embeddings for efficient large-scale multilingual continued pretraining](#). *CoRR*, abs/2311.08849.
- Meta-AI. 2024. Llama 3. <https://github.com/meta-llama/llama3?tab=readme-ov-file>.
- Benjamin Minixhofer, Edoardo Maria Ponti, and Ivan Vulić. 2024. [Zero-shot tokenizer transfer](#). *Preprint*, arXiv:2405.07883.
- OpenAI. 2023. [GPT-4 technical report](#). *ArXiv preprint*, abs/2303.08774.
- OpenCompass Team. 2023. Opencompass: A universal evaluation platform for foundation models. <https://github.com/open-compass/opencompass>.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Hamza Alobeidli, Alessandro Cappelli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. [The refinedweb dataset for falcon LLM: outperforming curated corpora with web data only](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- 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](#). *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Nived Rajaraman, Jiantao Jiao, and Kannan Ramchandran. 2024. [Toward a theory of tokenization in llms](#). *CoRR*, abs/2404.08335.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [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)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Aaditya K. Singh and DJ Strouse. 2024. [Tokenization counts: the impact of tokenization on arithmetic in frontier llms](#). *CoRR*, abs/2402.14903.
- Xinying Song, Alex Salcianu, Yang Song, Dave Dopson, and Denny Zhou. 2021. [Fast wordpiece tokenization](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 2089–2103. Association for Computational Linguistics.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. 2023. [Challenging big-bench tasks and whether chain-of-thought can solve them](#). In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 13003–13051. Association for Computational Linguistics.

Together Computer. 2023. [Redpajama: An open source recipe to reproduce llama training dataset](#).

Arseny Tolmachev, Daisuke Kawahara, and Sadao Kurohashi. 2018. [Juman++: A morphological analysis toolkit for scriptio continua](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP 2018: System Demonstrations, Brussels, Belgium, October 31 - November 4, 2018*, pages 54–59. Association for Computational Linguistics.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#). *ArXiv preprint*, abs/2302.13971.

Ke M. Tran. 2020. [From english to foreign languages: Transferring pre-trained language models](#). *CoRR*, abs/2002.07306.

Kai Xiong, Xiao Ding, Yixin Cao, Ting Liu, and Bing Qin. 2023. [Examining inter-consistency of large language models collaboration: An in-depth analysis via debate](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7572–7590, Singapore. Association for Computational Linguistics.

Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. [Byt5: Towards a token-free future with pre-trained byte-to-byte models](#). *Trans. Assoc. Comput. Linguistics*, 10:291–306.

Sha Yuan, Hanyu Zhao, Zhengxiao Du, Ming Ding, Xiao Liu, Yukuo Cen, Xu Zou, Zhilin Yang, and Jie Tang. 2021. [Wudaocorpora: A super large-scale chinese corpora for pre-training language models](#). *AI Open*, 2:65–68.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.

Wenqing Zheng, S. P. Sharan, Ajay Kumar Jaiswal, Kevin Wang, Yihan Xi, Dejia Xu, and Zhangyang Wang. 2023. [Outline, then details: Syntactically guided coarse-to-fine code generation](#). In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 42403–42419. PMLR.

Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. [Agieval: A human-centric benchmark for evaluating foundation models](#). *CoRR*, abs/2304.06364.