

FuLG: 150B Romanian Corpus for Language Model Pretraining

Vlad-Andrei Bădoi
University Politehnica of
Bucharest
Bucharest, Romania

Mihai-Valentin Dumitru
University Politehnica of
Bucharest
Bucharest, Romania

Alexandru M.
Gherghescu
University Politehnica of
Bucharest
Bucharest, Romania

Alexandru Agache
University Politehnica of
Bucharest
Bucharest, Romania

Costin Raiciu
University Politehnica of
Bucharest
Broadcom Inc.
Bucharest, Romania

ABSTRACT

Research in the field of language models is rapidly evolving, with many open models being released to the public. Openly available pretraining corpora usually focus on only a handful of languages, with many others either missing completely or extremely underrepresented. In this report, we introduce FuLG¹, a hundred-fifty-billion-token Romanian corpus extracted from CommonCrawl. We present our methodology for filtering FuLG and compare it via ablation studies against existing Romanian corpora.

1 INTRODUCTION

In a world where Large Language Models (LLMs) have shown great potential for natural language tasks, the corpora used for pretraining play a central role in their overall capabilities. However, most state-of-the-art models are not accompanied by their training datasets and authors only sparsely discuss their composition.

In response to this, researchers and practitioners have produced numerous data corpora for widely spoken languages on the Internet. These are often derived from CommonCrawl², a public repository of crawled web pages, which has indexed more than 250 billion pages. Other datasets go further by including curated data such as books, social media discussions, and research papers. Trillion-token datasets such as Dolma [19], FineWeb [14], or RedPajama [1] have enabled the development of billion-parameter models with loss-optimal training, but they cover only a small fraction of languages.

This limitation is evident in the language understanding capabilities of open models, as they are often trained on publicly available data. For instance, OLMo models [5] support only six languages, while other initiatives such as LLM360

K2 [9] are English-centric. Even models developed by tech giants, such as Llama [22] or Gemma [20], have limited output capabilities in less commonly spoken languages. Only recently have we seen some developments in this area, with Llama 3 supporting 30 languages, albeit with lower performance levels compared to English and limited information about the training dataset.

To improve the standing of the Romanian language in future models, we delve into filtering Romanian content from CommonCrawl, documenting the process and releasing an open dataset. We built a pipeline that employs deduplication techniques, common signal filters, and FastText for language detection. We document the steps taken to obtain the final text corpus. As a result, we release FuLG, an openly available corpus that is three times larger than existing Romanian corpora.

In this work we make two main contributions. We release FuLG, a 156B-token (589GB tokenized), or 220B tokens with the Llama 3 tokenizer, corpus for LLM pretraining and fine-tuning in the Romanian language. Second, we document the process of filtering data from CommonCrawl, which may be employed for other underrepresented languages.

2 RELATED WORK

A crucial factor in model performance is the size, quality and diversity of the pretraining corpus. These datasets typically comprise content from various sources, primarily web pages, due to their sheer number and availability, followed by social media discussions, academic papers, books, code repositories, and encyclopedias, covering a wide range of topics. The main sources of data are either private crawls of the Internet or the publicly available CommonCrawl repository.

¹<https://hf.co/datasets/faur-ai/fulg>

²<https://commoncrawl.org>

While information about closed models’ training data is scarce, even open models like Llama or Gemma lack transparency regarding their training datasets, as it is a central part of the competitive advantage. We are thus in a peculiar scenario, where the best open-source models aren’t truly open-source when it comes to pretraining data, which deeply affects any smaller initiative.

On the other hand, open datasets are essential components for truly open models and play a vital role in democratizing future LLMs. Currently, numerous open datasets exist, many derived from CommonCrawl. Notable examples include: Dolma [19] (3T tokens), C4 [2] (175B tokens), The Pile [4] (387B tokens), ROOTS [8] (400B tokens), Refined-Web [15] (600B tokens), RedPajama v2 [1] (30T tokens), FineWeb [14] (15T tokens), Zyda [21] (1.3T tokens), LLM360 Amber [9] (1.2T tokens), MAP-Neo [26] (4.5T tokens), OSCAR [13]. While these datasets provide sufficient quality and size for English-based models, they often lack adequate representation of less commonly spoken languages. For example, after training a Romanian GPT-NeoX tokenizer on the Romanian part of OSCAR, we obtain a number of only 10B tokens on OSCAR, and 41B tokens on mC4 [17] (slightly higher number of tokens if using an English tokenizer, similar to Llama or OLMo, due to worse compression ratio). However, considering the scaling laws of dataset size relative to model size, these quantities are insufficient for optimal LLM training. Beyond these, several smaller datasets exist for various Natural Language Processing (NLP) tasks in Romanian [6, 10–12, 18, 23] but they generally contain less than 500M tokens, which is far from adequate for LLM development.

3 DATA ACQUISITION AND FILTERING

FuLG is derived from CommonCrawl, a vast collection of web page crawls dating back to 2007. CommonCrawl releases snapshots of portions of the Internet several times annually. These snapshots exhibit low similarity rates between releases, enabling us to consistently extract new data from each snapshot. Our work encompasses snapshots from 2013 to May 2024.

We developed a pipeline based on existing software, which we elaborate on shortly. Given the petabytes of data requiring filtration, we utilized a distributed environment with multiple nodes for data acquisition. For deduplication and quality filtering, we employed a single large-memory node.

Data Acquisition. To process CommonCrawl snapshots, we leveraged the CCNet pipeline [24]. CCNet facilitates distributed processing of snapshots, including downloading in WET format, language identification via the FastText algorithm [7], and deduplication of common paragraphs. We encountered two primary challenges with CCNet:

- (1) Due to the fact that CCNet development has stalled, we quickly hit roadblocks related to package versions and environment-specific problems with our SLURM setup, which needed modifications to the source code.
- (2) Since we ran the pipeline on a cluster shared with other users, we did not have all the hardware for ourselves. Strict limits imposed by the SLURM system for fair sharing meant we needed to adapt the source code, which was not built with this in mind. We faced limitations with the maximum size of job arrays supported, often capped at 100. To address this, we extended CCNet to support job batching. Initially, CCNet would submit a job array equal in size to the number of shards; our modifications enabled it to break the job array into multiple smaller submissions to the SLURM scheduler.

We kept only documents with a language score for Romanian exceeding 0.5, meaning that there are documents that may include other languages alongside Romanian.

Deduplication. For filtering and deduplication, we employed code from RedPajama [1]. Exact deduplication reduced the dataset size by 37%, while fuzzy deduplication with a 0.8 threshold further reduced it by 50%. We updated the deduplication pipeline to not only identify duplicates but also remove them from the dataset.

Content Filtering. Next, we introduced a content filtering step to our pipeline. Using a regex-based approach, we filtered HTML text extraction artifacts, such as navbar text, from documents. To filter potentially controversial content, we utilized a dictionary, removing documents containing specific words in either the content or URL. For Personally Identifiable Information (PII), we replaced phone numbers, email addresses, and links with special tokens.

Quality Filtering. As with deduplication, we utilized existing code from RedPajama to compute a set of quality signals. We then extended the code to filter documents based on the quality signals introduced by Gopher [16], reducing the dataset size by 50% to 156B tokens using the OLMo tokenizer trained on Romanian. (589GB). With a different tokenizer such as Llama 3, which was trained on many other languages, we obtain around 220B tokens; we expect much worse compression for English-only tokenizers, such as Llama 2. The thresholds and rules we employed for filtering are as follows:

- Fraction of characters in the most common n-grams exceeded: 0.2 for 2-grams, 0.18 for 3-grams, 0.16 for 4-grams, 0.15 for 5-grams, 0.14 for 6-grams, 0.13 for 7-grams, 0.12 for 8-grams, 0.11 for 9-grams, 0.10 for 10-grams
- Fewer than 50 words or more than 100,000 words

- Median word length less than 3 or greater than 10 characters
- More than 90% of lines start with bullet points
- More than 30% of lines end with ellipses
- Less than 30% of lines end with punctuation

The data acquisition process spanned several months, while deduplication and filtering were completed in less than 18 hours on a single large-memory node.

4 FuLG IN ACTION

To evaluate FuLG we conducted ablation studies using a 1B decoder-only model based on OLMo [5], with a sequence length of 2048 and the OLMo tokenizer trained on the OSCAR-Ro corpus. We used a global batch size of 256. We made several changes to the hyperparameters to adjust to the 1B model size: weight decay of 0.001 and a max learning rate of 4e-4 annealed by a cosine schedule, with an initial warmup of 1000 steps.

We trained three identical models on different pretraining datasets: Oscar, mC4 and FuLG. We leveraged the Hugging Face transformers [25] library for training. We trained to completion on each dataset (290k steps for FuLG, 75k for mC4 and 20k for OSCAR). Since FuLG was much bigger than both OSCAR and mC4 (about 4 times bigger than mC4), we also included an earlier checkpoint from FuLG at 70k steps. We do note that it slightly underperforms, but we believe this naturally happens because of the cosine decay learning rate schedule.

Dataset	Perplexity
FuLG (290k)	16.06
FuLG (70k)	21
mC4	15.38
OSCAR	23.57

Table 1: Perplexities for datasets.

Perplexity. A commonly employed method to assess the quality of a dataset is to fix a model, train it with different datasets, and measure perplexity against a curated evaluation set. While perplexity is not a definitive measure of dataset quality [19], we use it as a sanity check, with values similar to existing datasets confirming our approach. We constructed a perplexity dataset covering multiple domains, sourced from Wikipedia, news articles, textbooks, research papers, and books, totaling 74M tokens. Although we did not specifically perform a separate decontamination step, as this is only a preliminary evaluation, we do note that the collection methods of the evaluation dataset were manual, with care in picking high-quality sources. We do plan to apply a decontamination step in further experiments, to make

sure our results are fair. Before calculating perplexity, the evaluation dataset was cleaned using the *clean-text* Python package³. The results are shown in Table 1. Both ablations of FuLG show perplexity values similar to those of existing datasets.

Dataset	Grammar	Creativity	Complexity
FuLG (290k)	8.5	7.1	4.6
FuLG (70k)	7.25	5.875	4
mC4	7.5	5.2	3.1
OSCAR	6.3	4.2	2.8

Table 2: Results for the story generation task.

Story generation We further qualitatively evaluated the trained models. Inspired by the TinyStories [3] work, we used the models to generate stories from given prompts and asked GPT-4 to rate the creativity, grammar, and overall complexity of each story. We only considered responses that were coherent stories, discarding anomalous outputs. Both OSCAR and FuLG (70k) generated significantly more anomalies than the other ablations. In Table 2, we present our findings, which suggest that FuLG may enable better performance for specific tasks.

5 DISCUSSION AND FUTURE WORK

The availability of high-quality datasets for less commonly spoken languages is crucial for the democratization of LLMs. While proprietary models often demonstrate proficiency across a wide range of languages, open models frequently underperform in this aspect. By developing large, high-quality open datasets for diverse languages, we can foster the creation of superior open-source models. These improved models could potentially be utilized by governments worldwide to meet their digitalization needs.

FuLG represents an initial effort to improve dataset size and quality for the Romanian language. We’re already underway to implementing a few other improvements which will further increase dataset quality:

- Data processing: We currently work with Common Crawl snapshots in WET format. However, as previous research [citation] indicates, employing a more sophisticated HTML parser could significantly improve both the quality of extracted text and the overall volume of usable data.
- Language-specific optimization: Romanian language particularities could be leveraged to adapt and refine quality filters. Currently, we use thresholds designed for the English language, but developing language-specific criteria could yield better results.

³<https://github.com/jfilter/clean-text>

- Novel quality filters: Identifying and implementing new quality filters tailored to the Romanian language could further enhance the dataset's overall quality and relevance.

ACKNOWLEDGEMENTS

We acknowledge the EuroHPC Joint Undertaking for awarding this project access to the EuroHPC supercomputer LUMI, hosted by CSC (Finland) and the LUMI consortium through a EuroHPC Regular Access call. Some of the experiments were performed on the Luxembourg national supercomputer MeluXina. The authors gratefully acknowledge the LuxProvide teams for their expert support. We acknowledge the computational resources provided by the PRACE award granting access to Discoverer in SofiaTech, Bulgaria. The authors gratefully acknowledge the HPC RIVR consortium and EuroHPC JU for funding this research by providing computing resources of the HPC system Vega at the Institute of Information Science. We acknowledge VSB – Technical University of Ostrava, IT4Innovations National Supercomputing Center, Czech Republic, for awarding this project access to the LUMI supercomputer, owned by the EuroHPC Joint Undertaking, hosted by CSC (Finland) and the LUMI consortium through the Ministry of Education, Youth and Sports of the Czech Republic through the e-INFRA CZ (grant ID: 90254). This work was funded by a VMWare gift.

REFERENCES

- [1] Together Computer. 2023. *RedPajama: an Open Dataset for Training Large Language Models*. <https://github.com/togethercomputer/RedPajama-Data>
- [2] Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. 2021. Documenting large webtext corpora: A case study on the colossal clean crawled corpus. *arXiv preprint arXiv:2104.08758* (2021).
- [3] Ronen Eldan and Yuanzhi Li. 2023. Tinystories: How small can language models be and still speak coherent english? *arXiv preprint arXiv:2305.07759* (2023).
- [4] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027* (2020).
- [5] Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. 2024. Olmo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838* (2024).
- [6] Radu Ion, Elena Irimia, Dan Stefanescu, and Dan Tufis. 2012. ROM-BAC: The Romanian Balanced Annotated Corpus. In *LREC*. 339–344.
- [7] Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Hervé Jégou, and Tomas Mikolov. 2016. Fasttext. zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651* (2016).
- [8] Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, et al. 2022. The bigscience roots corpus: A 1.6 tb composite multilingual dataset. *Advances in Neural Information Processing Systems* 35 (2022), 31809–31826.
- [9] Zhengzhong Liu, Aurick Qiao, Willie Neiswanger, Hongyi Wang, Bowen Tan, Tianhua Tao, Junbo Li, Yuqi Wang, Suqi Sun, Omkar Pangarkar, et al. 2023. Llm360: Towards fully transparent open-source llms. *arXiv preprint arXiv:2312.06550* (2023).
- [10] Mihai Manolescu and Çağrı Çöltekin. 2021. Roff-a romanian twitter dataset for offensive language. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*. 895–900.
- [11] Ludmila Midrigan-Ciochina, Victoria Boyd, Lucila Sanchez-Ortega, Diana Malancea_Malac, Doina Midrigan, and David P Corina. 2020. Resources in Underrepresented Languages: Building a Representative Romanian Corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. 3291–3296.
- [12] Verginica Barbu Mititelu, Dan Tuflă, and Elena Irimia. 2018. The reference corpus of the contemporary Romanian language (CoRoLa). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- [13] Pedro Javier Ortiz Suárez, Benoît Sagot, and Laurent Romary. 2019. 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*), Piotr Bański, Adrien Barbaresi, Hanno Biber, Evelyn Breiteneder, Simon Clematide, Marc Kupietz, Harald Lüngen, and Caroline Iliadi (Eds.). Leibniz-Institut für Deutsche Sprache, Mannheim, 9 – 16. <https://doi.org/10.14618/ids-pub-9021>
- [14] Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. 2024. The FineWeb Datasets: Decanting the Web for the Finest Text Data at Scale. *arXiv:2406.17557* <https://arxiv.org/abs/2406.17557>
- [15] Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidi, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116* (2023).
- [16] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446* (2021).
- [17] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *arXiv e-prints* (2019). [arXiv:1910.10683](https://arxiv.org/abs/1910.10683)
- [18] Christof Schöch, Tomaž Erjavec, Roxana Patras, and Diana Santos. 2021. Creating the european literary text collection (elitec): Challenges and perspectives. *Modern Languages Open* (2021).
- [19] Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Author, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. 2024. Dolma: An open corpus of three trillion tokens for language model pretraining research. *arXiv preprint arXiv:2402.00159* (2024).
- [20] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295* (2024).
- [21] Yury Tokpanov, Beren Millidge, Paolo Glorioso, Jonathan Pilault, Adam Ibrahim, James Whittington, and Quentin Anthony. 2024. Zyda: A 1.3 T Dataset for Open Language Modeling. *arXiv preprint arXiv:2406.01981* (2024).

- [22] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- [23] Tamás Váradi, Svetla Koeva, Martin Yamalov, Marko Tadić, Bálint Sass, Bartłomiej Nitóń, Maciej Ogródniczuk, Piotr Pezik, Verginica Barbu Mititelu, Radu Ion, et al. 2020. The MARCELL legislative corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. 3761–3768.
- [24] Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2019. CCNet: Extracting high quality monolingual datasets from web crawl data. *arXiv preprint arXiv:1911.00359* (2019).
- [25] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Online, 38–45. <https://www.aclweb.org/anthology/2020.emnlp-demos.6>
- [26] Ge Zhang, Scott Qu, Jiaheng Liu, Chenchen Zhang, Chenghua Lin, Chou Leuang Yu, Danny Pan, Esther Cheng, Jie Liu, Qunshu Lin, et al. 2024. Map-neo: Highly capable and transparent bilingual large language model series. *arXiv preprint arXiv:2405.19327* (2024).