

Vikhr: The Family of Open-Source Instruction-Tuned Large Language Models for Russian

Anonymous ACL submission

Abstract

There has been a surge in the development of various Large Language Models (LLMs). However, text generation for languages other than English often faces significant challenges, including poor generation quality and the reduced computational performance due to the disproportionate representation of tokens in model’s vocabulary. In this work, we address these issues and introduce Vikhr, a new state-of-the-art open-source instruction-tuned LLM designed specifically for the Russian language. “Vikhr” refers to the name of the Mistral LLM series and means “strong gust of wind.” Unlike previous efforts for Russian that utilize computationally inexpensive LoRA adapters on top of English-oriented models, Vikhr features an adapted tokenizer vocabulary and undergoes the continued pre-training and instruction tuning of all weights. This approach not only enhances the model’s performance but also significantly improves its computational and contextual efficiency. The remarkable performance of Vikhr across various Russian-language benchmarks can also be attributed to our efforts in expanding instruction datasets and corpora for continued pre-training. Vikhr not only sets the new state of the art among open-source LLMs for Russian, but even outperforms some proprietary closed-source models on certain benchmarks. The model weights, instruction sets, and code will be publicly available¹.

1 Introduction

Instruction tuning has unlocked in Large Language Models (LLMs) vast zero-shot capabilities without the need of careful prompt engineering (Ouyang et al., 2022). The most rapid research and development efforts are currently devoted to English LLMs. There has been a surge in English open-source models: Llama series (Touvron et al., 2023a,b), Mistral series (Jiang et al., 2023), Vicuna series (Chiang et al., 2023), etc. This growth is driven

by the abundance of raw training data in English and dedicated efforts to create comprehensive sets of instruction-output pairs. Despite the fact that LLMs oriented on English have some multilingual capabilities (Zhao et al., 2024) due to small portions of texts in various languages leaked into their training datasets (Touvron et al., 2023a), their overall performance in these languages remains relatively low. Although they can usually generate portions of coherent texts, these models struggle with reasoning in non-English languages, lack culture-specific knowledge, and are highly inefficient in terms of tokenization. This inefficiency arises due to the way bite-pair tokenization algorithms work: they split the infrequent words into multiple tokens. Since multilingual data typically represents a small portion of the training dataset, non-English words are often split in many pieces. This leads to more steps during prompt processing and text generation, shorter effective context windows, and ultimately lower quality (Tikhomirov and Chernyshev, 2023; Petrov et al., 2024). This disparity places non-English languages at a disadvantage.

There is a research direction focused on developing multilingual LLMs that work well for multiple popular languages: BLOOMz (Muennighoff et al., 2023), mGPT (Shliazhko et al., 2022), Bactrian-X (Li et al., 2023), PALO (Maaz et al., 2024), Aya101 from CohereAI (Üstün et al., 2024), etc. These models are typically trained on rich multilingual datasets and are less skewed towards English. However, when aiming to perform well across multiple languages simultaneously, these models must still share their vocabulary and parameters. This often hinders their performance for each particular language in isolation, especially for the popular smaller model sizes, such as 7B and 13B.

The goal of maximizing the LLM performance for a specific language within a certain number of parameters has led researchers to develop bi-lingual LLMs. For example, Jais (Sengupta et al., 2023)

¹<http://anonymous.repo>

083 focus only on English and Arabic. The inclusion of
084 English data in pre-training alongside Arabic data
085 is motivated by the significantly larger volume of
086 English data available. This helps LLMs substan-
087 tially enhance skills such as logical and common
088 sense reasoning, which are also applied when gen-
089 erating text in Arabic.

090 Russian is one of the high-resource languages
091 and is typically represented in multilingual LLMs.
092 Additionally, there are several proprietary closed-
093 source LLMs, such as MTS AI, GigaChat, and Yan-
094 dexGPT, that meet or even surpass their English-
095 oriented flagship competitors when it comes to text
096 processing and generation in Russian. However,
097 controllable research often requires white-box ac-
098 cess to LLM logits and layer outputs, the ability to
099 modify weights and a model architecture, and con-
100 sistent answers for reproducibility, which is often
101 impossible in closed-source LLMs due to their con-
102 stant development and retirement. There are only a
103 few open-source LLMs designed for Russian (such
104 as Saiga (Gusev, 2023), ruGPT (AI Forever, 2022),
105 ruadapt (Tikhomirov and Chernyshev, 2023)). Of
106 these, only Saiga and ruadapt are instruction-tuned.

107 This work aims to build an efficient and effective
108 open-source instruction-following LLM for Rus-
109 sian facilitating multilingual natural language pro-
110 cessing research. Building even a small LLM that
111 targets a particular language from scratch requires
112 a lot of computational resources. Consequently,
113 many researchers simply fine-tune LoRA adapters
114 (Hu et al., 2021) for English-oriented LLMs on
115 some language-specific data. While this approach
116 can improve model generation quality, it does not
117 address computational inefficiency because the tok-
118 enizer and model vocabulary remain unchanged.
119 In contrast, our approach not only fine-tunes a
120 base LLM on Russian language data but also re-
121 constructs its underlying tokenizer and vocabulary,
122 alongside suggesting an improved method for con-
123 tinued pre-training. Additionally, we have signifi-
124 cantly expanded the available Russian datasets for
125 instruction tuning. The developed LLM achieves
126 state-of-the-art results for the Russian language
127 among other open-source counterparts across a
128 wide range of benchmarks.

129 Contributions of the paper are the following:

- 130 • We have constructed Vikhr – a state-of-the-
131 art open-source instruction-following LLM
132 oriented on the Russian language. In addition
133 to its high generation quality, Vikhr features

an efficient tokenizer that enables rapid text
generation and good context utilization.

- We have developed a pipeline for adapting
English-oriented LLMs to the Russian lan-
guage. The pipeline implements vocabulary
adaptation, continued pre-training with regu-
larization to prevent “catastrophic forgetting”,
and instruction tuning.
- We have expanded the datasets for continued
pre-training of Russian language models and
previously available instruction datasets.
- We conducted an extensive evaluation of sev-
eral open-source LLMs on evaluation bench-
marks for Russian, demonstrating that Vikhr
achieves new state-of-the-art results.

2 Related Work

One of the first notable series of generative LLMs
for Russian is ruGPT (AI Forever, 2022; Zmitro-
vich et al., 2023). The authors created several mod-
els trained for the vanilla language modelling task
with the sizes of up to 13b. The models were cre-
ated from the scratch and trained on large Russian
corpora. They are able to handle the linguistic
nuances of Russian more effectively than multilin-
gual models (Muennighoff et al., 2022). Since the
training data was mostly in Russian, these models
have efficient tokenization, but the lack of multilin-
gual data (e.g. in English) limits their performance.
ruGPT models are not instruction tuned.

Gusev (2023) suggests to leverage reasoning ca-
pabilities of existing English-oriented LLMs and
adapt them to the Russian language by training
LoRA adapters. They also create an Alpaca-like set
of Russian instruction-output pairs and performed
instruction tuning. They have established the Saiga
model series, which has a competitive performance
and used to be a reasonable choice for off-the-shelf
open-source Russian LLM for the past year. How-
ever, the tokenizer in these models is not adapted,
so they experience issues with context and compu-
tational efficiency.

Tikhomirov and Chernyshev (2023) address
these issues in Saiga. In addition to model tun-
ing on Russian data, they also adapt the model
tokenizer. They note that improving tokenization
helps to both improve the efficiency of the model
and its performance while reducing memory con-
sumption. However, during continued pre-training,
the authors freeze the model weights except LM
heads and token embeddings, which probably re-

Content	Length	Tokenization Result
Original Sentence	31	Машинное обучение изменяет мир
Mistral Tokenizer	13	['Ма', 'шин', 'ное', 'об', 'у', 'чение', 'из', 'мен', 'я', 'ет', 'ми', 'р']
Vikhr Tokenizer	7	['Ма', 'шин', 'ное', 'обучение', 'изменяет', 'мир']

Table 1: Tokenizer comparisons between the original Mistral model and Vikhr

sults in the suboptimal performance.

In this work, we take advantage of pre-trained English-oriented LLMs, adapt LLM tokenizer for better computational efficiency, leverage continued pre-training on vast Russian-language corpora with regularization for preventing “catastrophic forgetting”, construct a novel extended set of Russian instruction-output pairs, and perform instruction tuning. The created LLM adaptation pipeline along with the data for continued pre-training and instruction tuning enables Vikhr to achieve new state-of-the-art results for Russian, maintain high performance for English, and demonstrate high computational efficiency.

3 LLM Construction Pipeline

The construction of Vikhr starts from one of English-oriented LLMs. In this work, we discuss the Vikhr model based on Mistral 7B. The strong logical and common reasoning capabilities, as well as the extensive world knowledge present in these LLMs provide an excellent starting point for our model. These features partially transfer to Vikhr, enhancing its performance in generating text in Russian. The process of LLM adaptation to Russian starts with the vocabulary adaptation. Then we perform continued pre-training of the LLM on large Russian datasets to mitigate the vocabulary shift and introduce culture specific knowledge. Finally, we perform fine-tuning of Vikhr on a set of instruction-output pairs in Russian.

3.1 Vocabulary Adaptation

The big drawback of English-oriented LLMs is that each Russian word would be split into multiple tokens: a common case is when symbols in the word become an individual tokens (see example in Table 1). This slows down the generation by multiple times, reduces the amount of information

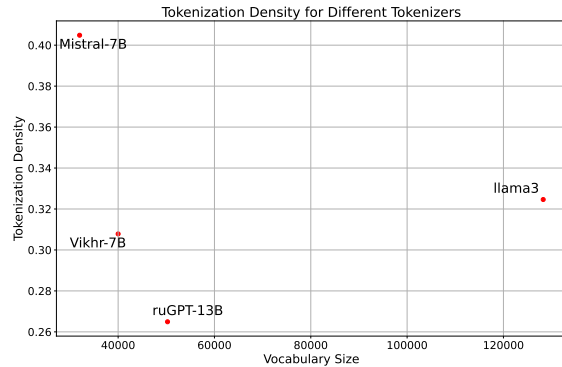


Figure 1: The Vikhr tokenizer efficiency in comparison to tokenizers of other models.

Data Source	Approx. size (GB)	Tokens (Billion)
Scientific papers	20	2.5
News articles	4	1
Wikipedia	25	4
Habr	6	1
Other sources	20	2.5

Table 2: The statistics of the Russian-language datasets for continued pre-training.

that could be stored in the context, and drastically hurts the generation quality.

To mitigate this problem in Vikhr, we adopt the approach suggested in (Cui et al., 2023; Tikhomirov and Chernyshev, 2023), where authors rebuild the tokenizer using a language-specific corpus. In particular, we trained a SentencePiece tokenizer (Kudo and Richardson, 2018) with a 40k vocabulary on the RuLM dataset (Gusev, 2023). As can be seen from Figure 1, the resulting tokenizer for Russian is much more efficient than the tokenizer of the original English-oriented model.

3.2 Continued Pre-training

The new vocabulary requires also new embedding matrices and LM heads. The tokens that were present in the original vocabulary are initialized with the old embeddings, the new tokens are initialized by averaging the embeddings of their pieces in the original embedding matrix (Hewitt, 2021). The similar approach is also applied to LM heads. Training model with these modifications requires much more computational resources than the mainstream technique for adaptation of LLMs to new languages based on LoRA adapters (Hu et al., 2021), as it requires to perform continued pre-training of the whole model and on much more

Hyperparam.	Value
LR	1×10^{-3}
AdamW eps	1×10^{-8}
Num warmup steps	10
AdamW betas	0.99, 0.95
Accumulation steps	128
Batch size	3
Epochs	1
Sequence length	1024

Table 3: The hyperparameters for continued pre-training.

language-specific data to mitigate the shift in the vocabulary.

The dataset for continued pre-training is constructed from Russian Wikipedia, news articles, scientific papers, top 100k up-voted posts on Habr, and some other sources. The statistics of these datasets is presented in Table 2. The total number of tokens used for this step is 11 billion.

We note that the continued pre-training of a LLM might partially eliminate the reasoning capabilities present in the original English-oriented model. This drastically affects the model performance. In our preliminary experiments, continued pre-training may result even in worse performance on Russian benchmarks compared to the original model. To alleviate the “catastrophic forgetting”, we use the loss regularization with KL penalty between the probability distribution of Vikhr and the reference English-oriented original LLM:

$$L_{\text{Vikhr}} = L_{\text{CE}} + KL(P_{\text{Vikhr}} || P_{\text{Ref}}). \quad (1)$$

In practice, we implement this approach using the SLERP interpolation of model losses (Goddard et al., 2024).

To speed up the process of continued pre-training, we use an optimized Flash attention implementation². As an optimization algorithm, we leverage AdamW as it trades some memory efficiency in favor of robustness to the hyperparameter choice. The hyperparameters used for continued pre-training are presented in Table 3.

3.3 Instruction Tuning

Instruction tuning is an essential step in reaching high zero-shot performance with LLMs. It also allows to obtain more natural communication with the model without complex prompting. Further fine-tuning techniques such as RLHF (Ouyang

²<https://huggingface.co/docs/optimum/bettertransformer/tutorials/convert>

et al., 2022), which require input from the assessors, are also crucial for such tasks as multicriteria alignment. However, the most significant performance gains are still achieved through instruction tuning (Jha et al., 2023).

Previously, Gusev (2023) constructed an open-source set of instruction-output pairs for the Russian language (Saiga). The core Saiga dataset was created similar to Alpaca by querying ChatGPT (gpt-3.5-turbo) (Taori et al., 2023). In this work, we extend this set by translating two English instruction datasets. First, we translated instructions for the FLAN model (Wei et al., 2021) and generated answers in Russian using ChatGPT. Originally, FLAN instructions were constructed automatically from annotated datasets using templates to facilitate multitask and zero-shot capabilities of seq2seq models. Later, it was shown that this data also helps to improve decoder-only chat-oriented models as well. Second, we construct Veles³ by translating the English OpenHermes (Teknium, 2023) instruction dataset. We also include without translation Nectar⁴ (Zhu et al., 2023) – the English instruction dataset. It helps to keep the performance of Vikhr high also for English. Since the majority of the outputs were machine generated there are many low quality outputs. To mitigate this problem, we filtered out low quality pairs using a reward model trained on human data. For the reward model, we selected the e5-large-multilingual model (Wang et al., 2024). This model was particularly suitable for our needs due to its ability to handle multilingual data efficiently, ensuring that the classifier could accurately assess the quality of responses in both Russian and English. We trained reward model on answer preference dataset⁵. These dataset collect from human prompts, and marked up with gpt4. By applying this reward model, we filtered out low-quality instruction-output pairs, significantly enhancing the overall performance and reliability of our Vikhr instruction datasets.

The statistics of the Vikhr instruction datasets is presented in Table

5.

Contrary to Saiga, we do not use LoRA adapters and just as in the phase of continued pre-training, we update all model parameters. The hyperparameters for the instruction tuning phase are presented

³<http://anonymous.repo>

⁴<https://huggingface.co/datasets/berkeley-nest/Nectar>

⁵<http://anonymous.repo>

	PPL	ruMMLU
With Data Filtration	4.3	0.8
Without Data Filtration	6.45	0.63

Table 4: Performance Metrics with and without Data Filtration

Instruction Set	Language	# instances
Veles	Russian	30k
Nectar	English	50k
Saiga	Russian	100k
ruFLAN	Russian	500k

Table 5: The statistics of instruction datasets.

in Table 6.

3.4 Hardware

Vikhr was trained on eight NVIDIA A100 GPUs 80GB. We spend approximately 1,000 GPU hours for the continued pre-training phase and 60 hours for instruction tuning.

4 Experiments

4.1 Experimental Setup

Benchmarks. The evaluation was performed on MMLU (Hendrycks et al., 2021), Ru-MMLU⁶, CheGeKa, Russian SuperGLUE (Shavrina et al., 2020), and MERA (Fenogenova et al., 2024). MMLU (En-MMLU) evaluates LLMs across 57 subjects with multiple-choice questions, assessing a model’s broad knowledge and reasoning abilities. We use this benchmark to verify that the model retains bi-lingual capabilities. In the results, we report the accuracy@1 score. RuMMLU is a translation of MMLU with GPT-3.5 to Russian. Just as for MMLU, we report the accuracy@1 score. CheGeKa is based on questions from the game “What? Where? When?”. This benchmark contains

⁶https://github.com/NLP-Core-Team/mmlu_ru

Hyperparam.	Value
LR	1×10^{-5}
AdamW, eps	1×10^{-8}
Num warmup steps	10
AdamW, betas	0.99, 0.95
Accumulation steps	64
Batch size	3
Num epochs	3
Sequence length	1024

Table 6: The hyperparameters for instruction tuning.

challenging open-ended questions, requiring logical reasoning and world knowledge. It includes 29,376 training and 416 test instances. The reported evaluation metric is the F1 score. Russian SuperGLUE is a benchmark similar to well-known English SuperGLUE (Wang et al., 2019). It tests LLMs on various natural language understanding tasks like reading comprehension and textual entailment. The metric reported in the results is accuracy@1. The MERA benchmark encompasses 21 evaluation tasks for generative LLMs in 11 skill domains. Note that among other tasks MERA also includes CheGeKa, RuMMLU, and one of the sub-tasks of SuperGLUE (RWSD). The reported evaluation metric is the total score, which is the average of scores across all non-diagnostic tasks.

Baselines. We compare Vikhr to six open-source and two proprietary closed-source competitors of the similar size. Open-source models: aya101 – a massively multilingual LLM from CohereAI that follows instructions in 101 languages⁷, it shows state-of-the-art results among massively multilingual LLMs; Mistral-7B-0.2-instruct – an English-oriented LLM that was used as the base model for Vikhr; rccmsu/ruadapt_mistral_saiga_7b_v0.1 – a Russian-oriented LLM that was constructed from the Mistral model using similar adaptations of the tokenizer, token embeddings, and the LM head (Tikhomirov and Chernyshev, 2023); saiga-mistral-7b-lora and saiga-llama3-8b – two versions of the Saiga models based on English-oriented LLMs and obtained by fine-tuning LoRA adapters on the Saiga instruction dataset⁸. Closed-source proprietary models for Russian: MTS AI Chat⁹ and GigaChat-7b. The access to GigaChat weights is closed, so the reported results are taken from the leaderboards¹⁰. The results of MTS AI Chat are also taken from the leaderboard¹¹.

4.2 Results

The evaluation results are presented in Table 7. As we can see, Vikhr outperforms all open-source models, including the ones that were built specifically for Russian. It also slightly outperforms its parent model Mistral on the En-MMLU benchmark, which might be the result of longer pre-training.

⁷<https://huggingface.co/CohereForAI/aya-101>

⁸<https://huggingface.co/collections/IlyaGusev>

⁹https://huggingface.co/MTSAIR/multi_verse_model

¹⁰<https://mera.a-ai.ru/ru/submits/10257>

¹¹<https://mera.a-ai.ru/ru/submits/10290>

LLM	Pre-train on Russian	Training Method	En-MMLU	Ru-MMLU	CheGeKa	Russian SuperGLUE	MERA
MTS AI Chat 7B (closed-source) \diamond	false	sft+dpo	-	0.689	0.083	0.56	0.479
GigaChat-7B (closed-source) \diamond	true	sft+dpo	-	0.67	0.451*	0.71*	0.479
aya101	false	pt+sft	0.41	0.37	0.005	0.36	0.320
Mistral-7B-Instruct-v0.2	false	none	0.60	<u>0.78</u>	0.005	0.57	0.400
rccmsu/ruadapt-mistral-7b-v0.1	false	pt+sft	0.61	0.72	0.005	0.64	0.421
rugpt13b	true	none	0.25	0.25	0.132	0.52	0.208
saiga-mistral-7b-lora	false	sft	0.60	0.76	0.223	0.64	0.442
saiga-llama3-8b	false	sft	0.59	<u>0.78</u>	<u>0.225</u>	<u>0.66</u>	<u>0.476</u>
Vikhr-7B-instruct_0.2	true	pt+sft	0.62	0.80	0.231	0.67	0.485

Table 7: Evaluation results for Russian and multilingual LLMs. Pre-train on Russian means that the model underwent (continued) pre-training on Russian data. The following abbreviations are used: sft – instruction tuning, pt – (continued) pre-training; dpo – direct preference optimization. \diamond The results for GigaChat and MTS AI are taken from the leaderboards. The best result among open-source models is highlighted with bold, the second best is underscored. The best result among closed-source proprietary models is marked with *.

The second place with close scores for all 4 Russian language benchmarks is obtained by the Saiga model based on recently released Llama-3. The high scores of this model probably are the result of the transfer of the outstanding performance of Llama-3. Since Saiga based on Llama-3 outperforms Saiga based on Mistral, we expect that applying our adaptation pipeline to Llama-3 would also help further improving the state of the art.

We note that the original Mistral-7B-0.2-instruct, despite being an English-oriented model, demonstrates competitive performance in 3 out of 4 Russian benchmarks. This indicates demonstrates that such models could be viable alternatives. The only dataset, where its performance is very low is CheGeKa, which is related to open-ended question-answering. This may be due to the lack of culture-specific knowledge, as the English-oriented model has not seen much Russian texts. Note that the MTS AI Chat also shows very low results on CheGeKa, which might also indicate the lack of culture-specific knowledge.

The proprietary model GigaChat substantially outperforms Vikhr on CheGeKa and notably on Russian SuperGLUE. We assume this is due to the use of much larger Russian datasets for pre-training. However, surprisingly, it falls behind Vikhr on Ru-MMLU. On all benchmarks, Vikhr outperforms the the proprietary competitor from MTS AI.

5 Conclusion

We have presented Vikhr – a new state-of-the-art open-source instruction-following LLM oriented on the Russian language. To create Vikhr, we developed a comprehensive pipeline for adapting English-oriented LLMs to Russian. The pipeline

includes the adaptation of the tokenizer vocabulary, continued pre-training of the entire model, and instruction tuning. We have also constructed a new dataset for instruction tuning by expanding the Saiga dataset with automatically translated and cleaned English instruction datasets. Our extensive work enabled Vikhr to outperform the known baselines, while maintaining computational efficiency.

We hope that the published models will foster the research on LLMs and enhance the diversity of languages incorporated into research agendas.

Limitations

We do not introduce additional restrictions to the usage of our models. However, the users must comply with the license of the base model and instruction datasets.

We do not implement RLHF / DPO fine-tuning of Vikhr due to the lack of the resources for human annotation. We expect further performance improvements from these techniques.

We do not introduce additional instruction-output pairs to facilitate LLM alignment. However, we note that the majority of the data for supervised fine-tuning of Vikhr are obtained from the ChatGPT model series, so our model partially inherits its alignment.

Ethical Considerations

The development and deployment of Vikhr raise several ethical considerations that must be addressed to ensure its responsible use:

- Bias and Fairness: For developing Vikhr, we use publicly available data. Despite efforts to train Vikhr on diverse datasets, there is a risk

466	of inherent biases in the data which may be	John Hewitt. 2021. Initializing new word embeddings for pretrained language models .	516
467	reflected in the model’s outputs. Continuous		517
468	monitoring and evaluation are required to mit-	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu,	518
469	igate any biases, ensuring fair and unbiased	Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,	519
470	performance.	et al. 2021. Lora: Low-rank adaptation of large lan-	520
		guage models. In <i>International Conference on Learn-</i>	521
471	• Misinformation: As with any LLM, Vikhr	ing Representations.	522
472	has the potential to generate misleading or		
473	incorrect information. It is crucial to es-	Aditi Jha, Sam Havens, Jeremy Dohmann, Alex Trott,	523
474	tablish guidelines and mechanisms for users	and Jacob Portes. 2023. Limit: Less is more for in-	524
475	to verify the information provided by the	struction tuning across evaluation paradigms. <i>arXiv</i>	525
476	model, promoting critical assessment and	<i>preprint arXiv:2311.13133</i> .	526
477	cross-referencing with reliable sources.		
		Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	527
478	• Misuse: Vikhr can be used for malicious pur-	sch, Chris Bamford, Devendra Singh Chaplot, Diego	528
479	poses, such as generating harmful content,	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	529
480	spam, or deepfakes. Implementing usage re-	laume Lample, Lucile Saulnier, et al. 2023. Mistral	530
481	strictions and monitoring mechanisms to de-	7b. <i>arXiv preprint arXiv:2310.06825</i> .	531
482	tect and prevent misuse is critical to safeguard		
483	against these risks.	Taku Kudo and John Richardson. 2018. Sentencepiece:	532
		A simple and language independent subword tok-	533
		enizer and detokenizer for neural text processing. In	534
		<i>Proceedings of the 2018 Conference on Empirical</i>	535
		<i>Methods in Natural Language Processing: System</i>	536
		<i>Demonstrations</i> , pages 66–71.	537
484	References	Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji,	538
485	AI Forever. 2022. ru-gpts: Generative pre-trained trans-	and Timothy Baldwin. 2023. Bactrian-x: A multilin-	539
486	former models for russian. https://github.com/	gual replicable instruction-following model with low-	540
487	ai-forever/ru-gpts .	rank adaptation. <i>arXiv preprint arXiv:2305.15011</i> .	541
		Muhammad Maaz, Hanoona Rasheed, Abdelrahman	542
488	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,	Shaker, Salman Khan, Hisham Cholakkal, Rao M	543
489	Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan	Anwer, Tim Baldwin, Michael Felsberg, and Fa-	544
490	Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion	had S Khan. 2024. Palo: A polyglot large multi-	545
491	Stoica, and Eric P. Xing. 2023. Vicuna: An open-	modal model for 5b people. <i>arXiv preprint</i>	546
492	source chatbot impressing gpt-4 with 90%* chatgpt	<i>arXiv:2402.14818</i> .	547
493	quality .		
		Niklas Muennighoff, Thomas Wang, Lintang Sutawika,	548
494	Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient	Adam Roberts, Stella Biderman, Teven Le Scao,	549
495	and effective text encoding for chinese llama and	M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey	550
496	alpaca . <i>arXiv preprint arXiv:2304.08177</i> .	Schoelkopf, et al. 2022. Crosslingual generaliza-	551
		tion through multitask finetuning. <i>arXiv preprint</i>	552
497	Alena Fenogenova, Artem Chervyakov, Nikita Mar-	<i>arXiv:2211.01786</i> .	553
498	tynov, Anastasia Kozlova, Maria Tikhonova, Albina		
499	Akhmetgareeva, Anton Emelyanov, Denis Shevelev,	Niklas Muennighoff, Thomas Wang, Lintang Sutawika,	554
500	Pavel Lebedev, Leonid Sinev, et al. 2024. Mera:	Adam Roberts, Stella Biderman, Teven Le Scao,	555
501	A comprehensive llm evaluation in russian. <i>arXiv</i>	M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey	556
502	<i>preprint arXiv:2401.04531</i> .	Schoelkopf, et al. 2023. Crosslingual generalization	557
		through multitask finetuning. In <i>The 61st Annual</i>	558
503	Charles Goddard, Shamane Siriwardhana, Malikeh	<i>Meeting Of The Association For Computational Lin-</i>	559
504	Ehghaghi, Luke Meyers, Vlad Karpukhin, Brian	<i>guistics</i> .	560
505	Benedict, Mark McQuade, and Jacob Solawetz. 2024.		
506	Arcee’s mergekit: A toolkit for merging large lan-	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	561
507	guage models. <i>arXiv preprint arXiv:2403.13257</i> .	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	562
		Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	563
508	Ilya Gusev. 2023. rulm: A toolkit for training neural lan-	2022. Training language models to follow instruc-	564
509	guage models. https://github.com/IlyaGusev/	tions with human feedback. <i>Advances in neural in-</i>	565
510	rulm .	<i>formation processing systems</i> , 35:27730–27744.	566
		Aleksandar Petrov, Emanuele La Malfa, Philip Torr,	567
511	Dan Hendrycks, Collin Burns, Steven Basart, Andy	and Adel Bibi. 2024. Language model tokenizers	568
512	Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-	introduce unfairness between languages. <i>Advances</i>	569
513	hardt. 2021. Measuring massive multitask language	<i>in Neural Information Processing Systems</i> , 36.	570
514	understanding. <i>Proceedings of the International Con-</i>		
515	<i>ference on Learning Representations (ICLR)</i> .		

571	Neha Sengupta, Sunil Kumar Sahu, Bokang Jia,	Jun Zhao, Zhihao Zhang, Qi Zhang, Tao Gui, and Xu-	625
572	Satheesh Katipomu, Haonan Li, Fajri Koto,	anjing Huang. 2024. Llama beyond english: An em-	626
573	Osama Mohammed Afzal, Samta Kamboj, Onkar	pirical study on language capability transfer. <i>arXiv</i>	627
574	Pandit, Rahul Pal, et al. 2023. Jais and jais-chat:	<i>preprint arXiv:2401.01055</i> .	628
575	Arabic-centric foundation and instruction-tuned open		
576	generative large language models. <i>arXiv preprint</i>	Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu,	629
577	<i>arXiv:2308.16149</i> .	and Jiantao Jiao. 2023. Starling-7b: Improving llm	630
		helpfulness & harmlessness with rlaif.	631
578	Tatiana Shavrina, Alena Fenogenova, Emelyanov Anton,	Dmitry Zmitrovich, Alexander Abramov, Andrey	632
579	Denis Shevelev, Ekaterina Artemova, Valentin Ma-	Kalmykov, Maria Tikhonova, Ekaterina Taktasheva,	633
580	lykh, Vladislav Mikhailov, Maria Tikhonova, Andrey	Danil Astafurov, Mark Baushenko, Artem Snegirev,	634
581	Chertok, and Andrey Evlampiev. 2020. Russiansu-	Tatiana Shavrina, Sergey Markov, et al. 2023. A	635
582	perglue: A russian language understanding evalua-	family of pretrained transformer language models for	636
583	tion benchmark. In <i>Proceedings of the 2020 Con-</i>	russian. <i>arXiv preprint arXiv:2309.10931</i> .	637
584	<i>ference on Empirical Methods in Natural Language</i>		
585	<i>Processing (EMNLP)</i> , pages 4717–4726.	Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-	638
		Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel	639
586	Oleh Shliashko, Alena Fenogenova, Maria Tikhonova,	Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid,	640
587	Vladislav Mikhailov, Anastasia Kozlova, and Tatiana	Freddie Vargus, Phil Blunsom, Shayne Longpre,	641
588	Shavrina. 2022. mgpt: Few-shot learners go multilin-	Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer,	642
589	gual. <i>arXiv preprint arXiv:2204.07580</i> .	and Sara Hooker. 2024. Aya model: An instruction	643
		finetuned open-access multilingual language model.	644
590	Rohan Taori, Ishaan Shum, Pieter Abbeel, Carlos	<i>arXiv preprint arXiv:2402.07827</i> .	645
591	Guestrin, and Percy Liang. 2023. Stanford alpaca:		
592	An instruction-following language model. <i>GitHub</i> .		
593	Teknum. 2023. Openhermes 2.5: An open dataset of		
594	synthetic data for generalist llm assistants .		
595	Mikhail Tikhomirov and Daniil Chernyshev. 2023. Im-		
596	impact of tokenization on llama russian adaptation.		
597	<i>arXiv preprint arXiv:2312.02598</i> .		
598	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier		
599	Martinet, Marie-Anne Lachaux, Timothée Lacroix,		
600	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal		
601	Azhar, et al. 2023a. Llama: Open and effi-		
602	cient foundation language models. <i>arXiv preprint</i>		
603	<i>arXiv:2302.13971</i> .		
604	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-		
605	bert, Amjad Almahairi, Yasmine Babaei, Nikolay		
606	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti		
607	Bhosale, et al. 2023b. Llama 2: Open founda-		
608	tion and fine-tuned chat models. <i>arXiv preprint</i>		
609	<i>arXiv:2307.09288</i> .		
610	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-		
611	preet Singh, Julian Michael, Felix Hill, Omer Levy,		
612	and Samuel Bowman. 2019. Superglue: A stick-		
613	ier benchmark for general-purpose language under-		
614	standing systems. <i>Advances in neural information</i>		
615	<i>processing systems</i> , 32.		
616	Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang,		
617	Rangan Majumder, and Furu Wei. 2024. Multilin-		
618	gual e5 text embeddings: A technical report. <i>arXiv</i>		
619	<i>preprint arXiv:2402.05672</i> .		
620	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,		
621	Adams Wei Yu, Brian Lester, Nan Du, Andrew M		
622	Dai, and Quoc V Le. 2021. Finetuned language mod-		
623	els are zero-shot learners. In <i>International Confer-</i>		
624	<i>ence on Learning Representations</i> .		