

# RuWikiBench: Evaluating Large Language Models through encyclopedia article replication

*D. A. Grigoriev*<sup>12</sup>, *D. I. Chernyshev*<sup>13</sup>

© The Authors 2025. This paper is published with open access at SuperFri.org

In light of the growing interest in using large language models (LLMs) as tools for generating scientific texts, the evaluation of their ability to produce encyclopedic content is becoming increasingly relevant. However, for Russian-language materials this issue has not been sufficiently studied, and existing benchmarks do not cover key aspects of analytical work with sources. This paper presents RuWikiBench - an open benchmark based on Ruwiki for evaluating the ability of large language models to reproduce Wikipedia-style articles, built around three tasks: selection of relevant sources, article structuring, and section generation. The results popular open-source LLM evaluation show that even under ideal conditions, the best models do not always follow the expert logic of composing encyclopedic content: even with a perfect source retrieval system, the models cannot reproduce the reference outline, and the quality of section generation shows almost no improvement with increase of model parameters.

*Keywords: benchmark, Wikipedia, Ruwiki, large language model.*

## Introduction

Modern large language models demonstrate impressive text generation results for wide range of topics and domains. However, their capabilities for working with scientific and encyclopedic materials remain understudied, particularly for Russian-language texts. Existing methods for model capability evaluation predominantly focus on standard linguistic tasks, without paying sufficient attention to analytical abilities when working with scientific texts. For the Russian language, this problem is especially relevant due to the limited availability of specialized evaluation tools.

There are many benchmarks covering various linguistic tasks for the Russian language. RussianSuperGlue [12] evaluates general language understanding and basic natural language processing tasks. MERA [3] provides unified testing conditions for models by compiling generation instructions for each task; however, the tasks themselves are oriented towards testing general comprehension. LIBRA [2] focuses on testing a model's ability to retain and retrieve information from a large context but is centered on short answers that do not require deep reasoning. Ru Arena General [16] focuses on pairwise model comparison rather than overall answer quality. Ping-Pong [4] evaluates the dialog abilities of models, which is important for interactive systems, but is not suitable for assessing the ability to conduct research and write coherent scientific-encyclopedic texts. At the same time, an entire class of tasks related to deep text analysis remains uncovered: creating detailed, structured, and factually accurate texts supported by a large number of sources.

The recent development of new agent capabilities, such as the emergence of the "Deep Research" function by OpenAI [9] or the development of the universal Storm algorithm [11], indicates a growing interest in conducting scientific research using large language models. This highlights the need to create new approaches for objectively evaluating the analytical capabilities of models. Existing benchmarks only partially address aspects critical for generating scientific-encyclopedic texts, such as the ability to extract relevant information from a set of documents, plan the structure of a future text, maintain text coherence and fluency, as well as ensure the factual consistency.

---

<sup>1</sup>Lomonosov Moscow State University, Moscow, Russia

<sup>2</sup>E-mail: dagrig14@yandex.ru

<sup>3</sup>E-mail: chdanorbis@yandex.ru

One of the closest studies in this area is the ResearchArena [6] benchmark, which formalizes the construction of an academic review; however, it is more aimed at testing the models' ability to select and organize relevant information and does not address their ability to generate coherent scientific-encyclopedic texts.

This paper proposes an approach for evaluating large language model abilities to analyze and produce scientific-encyclopedic texts. Our main contributions:

1. We collect "Ruwiki", a labeled dataset for generation of wikipedia-style articles in Russian language;
2. We propose new benchmark, RuWikiBench that measures large language model analytical capabilities in Russian language;
3. We evaluate the abilities of popular open-source large language models to generate Wikipedia-style articles.

The code and data of this work have been made publicly available<sup>4</sup>.

## 1. Dataset collection

To measure analytical capabilities of large language models we need a labeled corpus of texts that would contain various sources, the respective analysis of those sources and analysis-source mapping for factuality evaluation of claims. Wikipedia articles are the best candidates for such data because this genre simultaneously requires proving factual accuracy, acknowledging different perspectives, and expressing complex topics in common terms.

We utilize articles from Russian online encyclopedia "RuWiki" which distinguishing features are a large number of references to Russian-language sources, as well as stricter text filtering. These qualities ensure that the articles can be reliably verified and reproduced by human experts and therefore replicated by LLMs.

The data acquisition process included the following steps:

1. **Article Selection:** Articles on diverse topics containing a sufficient number of references to external sources were manually selected;
2. **Source Downloading:** For each article, the available sources it references were web-scraped;
3. **Splitting into Snippets:** To reproduce real Retrieval Augmented Generation (RAG) conditions, all texts were split into small fragments of approximately  $\approx 600$  words in length.

Fig. 1 shows a brief schematic of the source text extraction process. The sources were downloaded using the Python module `newspaper3k`<sup>5</sup>. A subset of "Ruwiki" articles  $B$  is taken as the initial corpus. To extract article's HTML code we use BeautifulSoup<sup>6</sup>. The obtained text is structured by splitting it into fragments corresponding to nested headings (H1, H2, H3, etc.), which preserves both the substantive part of the article and its hierarchical organization. Next, all external references cited in the "Notes" section are automatically extracted. Invalid links (e.g., 404 error) are excluded from further processing, and the text associated with them is removed, leaving only those sources that are actually accessible.

Fig. 2 illustrates the schematic breakdown of an article<sup>7</sup> into key entities used in subsequent processing. During the data processing stage, text filtering is performed to ensure its correct

<sup>4</sup><https://github.com/Nejimaki-Tori/WikiBench>

<sup>5</sup><https://github.com/codelucas/newspaper>

<sup>6</sup><https://beautiful-soup-4.readthedocs.io/en/latest/>

<sup>7</sup><https://ru.ruwiki.ru/wiki/Python>

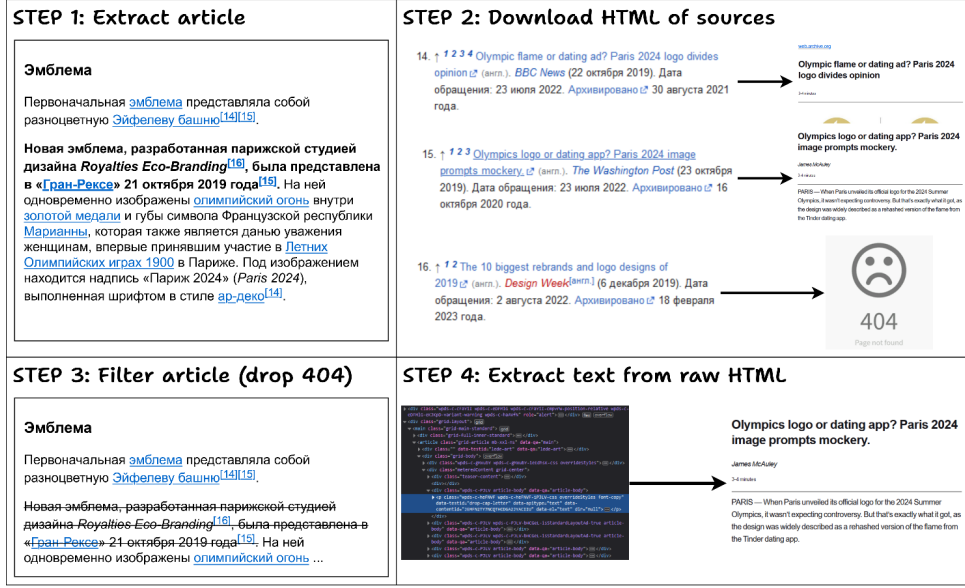


Figure 1. Source extract

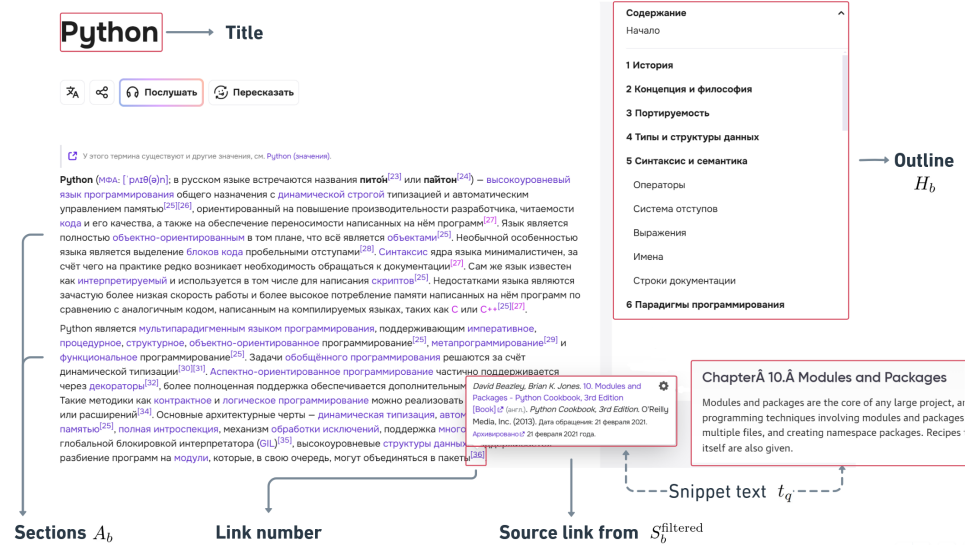


Figure 2. Main article entities

interpretation by the model. Each footnote (e.g., [1], [2]) is matched to a specific link corresponding to one of the available sources. This allows for precise identification of the link’s position in the article text and its use for filtering.

Table 1. Key characteristics of the collected dataset

Metric	RuWikiBench	ResearchArena
Number of articles	285	7,952
Number of sources	15,686	12,034,505
Total number of snippets	36,860	-
Average outline size (headings)	37	-
Average section size (words)	112	-

Based on the valid links  $S_b^{\text{filtered}}$ , filtered sets of paragraphs  $A_b^{\text{filtered}}$  and headings  $H_b^{\text{filtered}}$  are formed. That is, only content supported by the extracted sources is retained; everything else is removed. Only sources for which a text  $t_q$  of at least 1500 characters could be retrieved are kept, to filter out "noisy" responses from HTML pages such as errors (e.g., error 404) or blocking messages.  $A_b^{\text{filtered}}$  retains only those paragraphs that contain at least one link to a source for which text was successfully retrieved. Similarly,  $H_b^{\text{filtered}}$  is formed-only those headings under which at least one paragraph remains. The characteristics of the collected corpus are presented in Tab. 1.

## 2. Evaluation Methodology

To objectively assess the ability of language models to generate scientific and encyclopedic texts, it is necessary to replicate the real process of preparing encyclopedic content:

1. **Selection of relevant sources:** The model is given the article title and a set of snippets, among which it must identify and rank materials relevant to the topic by their significance;
2. **Article structure construction:** Based on the topic and selected sources, the model creates an outline with main sections in the Wikipedia-style;
3. **Section generation:** Article materials are distributed across sections, after which a summary of the relevant materials is generated for each section.

Each stage is evaluated independently of the previous ones, allowing for a quantitative measurement of the quality of each specific subtask.

### 2.1. Selecting relevant sources

One of the most effective search strategies [17] is the preliminary generation of an expected result (description) based on the original query (article title) to create an expanded search query. In our pipeline the description is generated in both Russian and English, to address available source text language diversity. The queries in both languages are then combined into a single textual query for a BM25-based search system.

Experiments were conducted with two approaches to expanded query generation:

1. **Ground-truth query based on the title and second-level headings:** A ground-truth query which is used for direct evaluation of the models' ranking abilities. We use LLaMa 3 70b model [15] to generate this query;
2. **Query generated from the title by the evaluated model:** Similar to real-world conditions, the LLM is fully responsible for the quality of the results and independently decides which search query to generate for BM25.

Examples of generated descriptions are shown in Tab. 2.

The documents selected by the BM25 query are sequentially passed to the large language model, which must classify each snippet as relevant (answer "yes") or non-relevant (answer "no"). To obtain numerical scores, we compare related wikipedia article titles of retrieved documents to title of article we seek to generate. The logarithmic probability of the tokens in the model's response is taken: if the answer is positive, the probability  $P(\text{yes})$  is used; if negative,  $1 - P(\text{no})$  is used. This approach allows ranking the retrieved documents by the model's confidence in their relevance: the higher the probability, the higher the model's confidence in the response, and the higher the document is ranked in the results.

**Table 2.** Comparison of english-translated descriptions of the article “C++” in two variants

Query Variant	Text
<b>Ground-truth</b>	The article "C++" is an overview of the C++ programming language, its history, structure, and features. It covers the main aspects of the language, including its standard library, differences from the C language, and further development. In addition, the article contains examples of C++ programs, comparisons with alternative programming languages, as well as critical analysis and discussion of the influence of C++ on the development of programming and existing alternatives. The article is intended for readers interested in the C++ language and its role in modern programming.
<b>Generation by title</b>	The article "C++" may be dedicated to the C++ programming language, one of the most popular and widely used programming languages in the world. The article may cover the basics of the language, its history, syntax and features, as well as its applications in various fields such as operating systems development, games, and web applications. In addition, the article may include information about C++ standards and libraries, as well as its comparison with other programming languages. The article can be useful both for beginner programmers and for experienced professionals who want to deepen their knowledge of C++. It may also include code examples and practical tips on using C++ in real projects.

## 2.2. Article structure planning

First, each text fragment (snippet) from the reference article source is converted by embedding model. Then, the snippets are clustered into potential section contents. To guarantee reproducibility KMeans algorithm is applied with the number of clusters equal to the number of second-level headings in the reference outline, and the centroids are initialized with the vector representations of these headings.

Next, five snippets closest to the cluster center are selected. This is done to reduce the influence of less relevant snippets on the final outline. The generation of mini-outlines for sections is carried out taking into account two key parameters: the context window size (to account for references and the overall semantics of the document) and two generation modes - directly from the texts and through preliminary generation of a brief cluster description. These generation modes enable different levels of abstraction: the direct mode preserves details with raw data, while the mode via preliminary cluster description improves consistency and reduces information duplication. At the final stage, all mini-outlines are combined into the final structured article outline.

## 2.3. Section generation

For each article section, we extract all snippets that were indicated as sources for the reference text of the section. These snippets are again converted into embeddings, and a pairwise similarity matrix is constructed as the product  $E \times E^\top$ . Elements with a similarity value above the threshold of 0.8 (empirically determined) are considered semantically close and are grouped together to avoid redundant repetitions during generation (e.g., when different sources paraphrase the same information). For each such semantic group, a hierarchical representation is built: the first five

texts are taken, and a brief description is generated based on them. This description is then supplemented using the next five texts, and so on, until a complete compressed representation of the group is obtained. Thus, only a set of brief descriptions remains - the most important information without excessive repetition. After this, the text of the section is generated based on the obtained group descriptions using a hierarchical summarization method [18].

### 3. Benchmark parameters

Below is a description of all data, models, hyperparameters, and procedures used to ensure reproducibility and analysis.

#### 3.1. Generation parameters

For all models, unless otherwise specified, the same generation parameters were used: temperature - 0.01, repetition penalty - 1.0, and top\_p - 0.9.

#### 3.2. Relevant source selection

Snippet indexing was performed using BM25<sup>8</sup> across the entire corpus of collected snippets without hyperparameter tuning (default values). For each relevant document, two non-relevant documents were selected (ratio 1:2) - this was done to improve evaluation robustness.

#### 3.3. Article structure planning

To build snippet embeddings we used `sergeyzh/BERTA`<sup>9</sup> model. Two context window options were considered: either a zero window (only the snippet itself) or one neighboring snippet to the left and right to expand the context. Heading similarity with reference headings was compared using cosine similarity: semantic correspondence was prioritized over exact wording or heading level. The comparison was performed against the cleaned article structure: all headings whose sections consisted entirely of text without downloadable sources were removed from the preprocessed text.

## 4. Metrics

Within the benchmark, two groups of metrics are considered: ranking metrics, which assess how well the model selects relevant sources, and text similarity metrics, which measure how closely the generated content matches the reference.

#### 4.1. Ranking Metrics

To evaluate the quality of the source list, we use NDCG@K [5] and R-Precision [1]:

$$\text{NDCG@K} = \frac{\text{DCG@K}}{\text{IDCG@K}}, \quad (1)$$

---

<sup>8</sup><https://github.com/xhluca/bm25s>

<sup>9</sup><https://huggingface.co/sergeyzh/BERTA>

$$\text{DCG@K} = \sum_{i=1}^K \frac{\text{rel}_i}{\log_2(i+1)}, \quad (2)$$

$$\text{IDCG@K} = \sum_{i=1}^K \frac{\text{rel}_i^{\text{IDEAL}}}{\log_2(i+1)}, \quad (3)$$

$$\text{R-Precision} = \frac{\sum_{i=1}^R \text{rel}_i}{R}, \quad (4)$$

where  $\text{rel}_i \in \{0, 1\}$  is the indicator of relevance for the document at position  $i$ ;  $\text{rel}_i^{\text{IDEAL}}$  is the same quantity in the ideal (fully sorted) ranking;  $R$  is the total number of relevant documents for the given query.

## 4.2. Text Similarity Metrics

The quality of generated sections and headings is evaluated with **BERTScore** [21]:

$$R_{\text{BERT}} = \frac{1}{|x|} \sum_{x_i \in x} \max_{\hat{x}_j \in \hat{x}} x_i^\top \hat{x}_j, \quad (5)$$

$$P_{\text{BERT}} = \frac{1}{|\hat{x}|} \sum_{\hat{x}_j \in \hat{x}} \max_{x_i \in x} x_i^\top \hat{x}_j, \quad (6)$$

$$F_{\text{BERT}} = \frac{2 P_{\text{BERT}} R_{\text{BERT}}}{P_{\text{BERT}} + R_{\text{BERT}}}, \quad (7)$$

where  $x$  is the reference text,  $\hat{x}$  is the generated text; each sentence is encoded using the model<sup>10</sup>, after which cosine similarity is computed.

ROUGE-L and BLEU were also considered for evaluating section generations.

**ROUGE-L** [7] is based on the length of the longest common subsequence (LCS) between the generated summary  $S$  and the reference  $R$ :

$$\text{Precision} = \frac{\text{LCS}(S, R)}{|S|}, \quad (8)$$

$$\text{Recall} = \frac{\text{LCS}(S, R)}{|R|}, \quad (9)$$

$$\text{ROUGE-L} = \frac{2 \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}, \quad (10)$$

**BLEU** [10] is an n-gram precision metric with a brevity penalty:

$$\text{BLEU}_N = \text{BP} \cdot \exp\left(\sum_{n=1}^N w_n \log p_n\right), \quad (11)$$

where  $p_n$  is the precision for  $n$ -grams,  $w_n$  are the weights, and BP is the brevity penalty.

<sup>10</sup><https://huggingface.co/sergeyzh/BERTA>

## 5. Experiments

### 5.1. Models

The experiments used the following large language models: RuadaptQwen2.5-7B-Lite-Beta [14], RuadaptQwen3-32B-Instruct-v2 [14], DeepSeek V3 [8], Qwen3-235B-A22B [20], tpro [13] and yagpt5lite [19]. In all tables, models are grouped by size, and the best results within each group are highlighted.

### 5.2. Results

Tab. 3 and Tab. 4 present the results of source selection evaluation. As the baseline we use BM25 retrieval without model-based reranking. In the first case (Tab. 3), where a pre-generated ground-truth search query was used for all models, the best results were achieved by DeepSeek V3, which indicates a strong ability to select relevant documents. In the second experiment (Tab. 4), where the query was based solely on the article title, tpro achieved the best results. The experiment showed that LLM query generation is not inferior in ranking quality to the reference setup. Presumably, this is because (as shown in Tab. 2) the queries turn out to be very similar to ground-truth: LLMs have knowledge of Wikipedia’s typical article structure from training and can therefore connect relevant concepts in the required format.

**Table 3.** Results of pure ranking ability evaluation

Model	NDCG	R-Pr
baseline (bm25)	88.81	62.51
<b>DeepSeek V3</b>	<b><u>95.42</u></b>	<b><u>83.86</u></b>
Qwen3-235B-A22B	94.49	82.42
RuadaptQwen3-32B-Instruct-v2	95.25	81.81
tpro	<u>95.42</u>	<u>83.53</u>
RuadaptQwen2.5-7B-Lite-Beta	88.26	62.26
yagpt5lite	<u>90.35</u>	<u>77.66</u>

**Table 4.** Results of evaluating BM25 query-generation ability

Model	BM25		Rerank	
	NDCG	R-Pr	NDCG	R-Pr
DeepSeek V3	88.39	60.65	<u>95.67</u>	<u>83.07</u>
Qwen3-235B-A22B	<u>89.17</u>	<u>62.98</u>	94.90	81.96
RuadaptQwen3-32B-Instruct-v2	85.39	52.80	95.82	81.62
<b>tpro</b>	<b><u>90.61</u></b>	<b><u>65.07</u></b>	<b><u>96.06</u></b>	<b><u>83.37</u></b>
RuadaptQwen2.5-7B-Lite-Beta	<u>88.81</u>	<u>62.51</u>	88.23	60.96
yagpt5lite	86.59	57.98	<u>90.27</u>	<u>77.65</u>

Overall, the models show fairly high performance at this stage, which may be due to the fact that an article title reflects its content well. In the best cases, up to 80% of the documents in the sample are relevant, which implies there is significant room for further improvement.



**Table 5.** Results of outline generation

Model	Mean BERTScore F1		
	Direct		Description
	no neighbors	one neighbor	
<b>DeepSeek V3</b>	<b><u>63.51</u></b>	<b><u>62.93</u></b>	<b><u>65.50</u></b>
Qwen3-235B-A22B	60.86	59.06	62.66
RuadaptQwen3-32B-Instruct-v2	60.12	<u>60.04</u>	<u>62.91</u>
tpro	<u>60.32</u>	59.09	60.75
RuadaptQwen2.5-7B-Lite-Beta	<u>60.03</u>	58.21	<u>61.58</u>
yagpt5lite	59.72	<u>60.07</u>	60.25

Tab. 5 presents the results of article structure planning. Direct - generation from the cluster snippets, with neighbor context as indicated; Description - generation via a preliminary description of all cluster elements. The results show that with preliminary description generation, all models consistently improve in quality. RuadaptQwen3 shows the largest gain, rising to second place and effectively matching the results of the larger model, Qwen3-235B-A22B. DeepSeek V3 remains the leader, showing a substantial margin over the others. At the bottom in quality are RuadaptQwen2.5-7B-Lite-Beta and yagpt5lite. At the same time, yagpt5lite, with only 8 billion parameters, delivers results comparable to a 32-billion-parameter model. Tab. 6 shows a comparison of a small excerpt of the reference and generated outlines. A common issue across all models was excessive heading hierarchy depth. On “Ruwiki”, headings were rarely deeper than level three; however, the models often created fourth- and fifth-level headings, implying that all information belongs in one large section, even though it may differ somewhat in meaning and, in the original outline, would correspond to unrelated headings.

**Table 6.** Comparison of two englis-translated article outlines

Generated	Reference
# Introduction to Python	# Python
## Language Overview	## History
### History and Key Aspects	## Concept and Philosophy
#### Main Features and Implementations	## Portability
# Python Basics	## Data Types and Structures
## Syntax and Semantics	## Syntax and Semantics
### Data Types and Structures	### Indentation System
#### Numbers, Lists, Dictionaries	### Expressions
and Object-Oriented Programming	### Names
# Advanced Python Topics	### Documentation Strings
## Flow Control and Multithreading	## Programming Paradigms
...	...

Tab. 7 reports the evaluation of section-generation quality. The difference between model performance may be considered marginal, but this is due to the sensitivity of the metric used. Sections for which the algorithm did not select a single relevant snippet were excluded from the final metrics. The best overall results were demonstrated by Qwen3-235B-A22B; however, in terms

of ROUGE-L and BLEU, RuadaptQwen3-32B-Instruct-v2 shows better structural consistency and greater phrase overlap with the reference. The yagpt5lite model performs above average, especially on BLEU, at a much smaller size, whereas tpro shows the lowest values across all metrics.

**Table 7.** Results of section generation

Model	Mean F1	Mean ROUGE-L	Mean BLEU
DeepSeek V3	53.48	14.34	2.81
Qwen3-235B-A22B	<b>53.74</b>	<u>14.63</u>	<u>3.07</u>
<b>RuadaptQwen3-32B-Instruct-v2</b>	53.21	<b>15.46</b>	<b>3.40</b>
tpro	53.15	13.58	2.27
RuadaptQwen2.5-7B-Lite-Beta	52.99	12.29	2.11
yagpt5lite	<u>53.43</u>	<u>14.85</u>	<u>3.16</u>

**Table 8.** Comparison of english-translated texts of two sections

Model	Text
DeepSeek V3	<b>COVID-19 is an infectious disease</b> , ... which led to a <b>global pandemic</b> that began in 2020. Initially presenting with respiratory symptoms such as cough, fever, and shortness of breath, the disease can cause severe complications, ... The virus is highly transmissible, is presumed to have a zoonotic origin, and spread rapidly worldwide. To control the pandemic, the WHO recommends vaccination, mask-wearing, social distancing, and hand hygiene, with vaccine effectiveness against the original strain reaching 85% or higher. Although COVID-19 in children more often has a mild course, severe cases are possible, including multisystem inflammatory syndrome.
yagpt5lite	<b>COVID-19 is a pandemic</b> caused by the novel coronavirus SARS-CoV-2. As of January 14, 2022, the WHO had confirmed about <b>318,648,834</b> cases of COVID-19 worldwide, including <b>5,518,343</b> deaths. The first COVID-19 vaccine was introduced in December 2020. On December 2, 2020, the United Kingdom became the first country to approve the Pfizer-BioNTech (BNT162) vaccine, which the WHO authorized for emergency use. SARS-CoV-2 is considered more contagious than SARS-CoV and quickly spread around the world after several infection cases in Wuhan, China. The pathogenesis of SARS-CoV-2 is associated with inflammatory responses that adversely affect the lungs and cause symptoms such as cough, fever, general malaise, shortness of breath, and respiratory failure.

For a qualitative evaluation of section-generation quality, we consider the introductory parts of the article “COVID19” produced by DeepSeek V3 and yagpt5lite, respectively, shown in Tab. 8. Despite some semantic inaccuracies (for example, the statement “COVID-19 is a pandemic,” whereas in reality it is a disease), yagpt5lite demonstrates a quite solid result. Its text falls short of DeepSeek V3’s version in terms of coverage and systematic exposition, but contains more numerical data and concrete facts. At the same time, the material generated by DeepSeek V3

tends to resemble an encyclopedic article, whereas the yagpt5lite version is closer in style to a technical report on the disease.

## Conclusion

This paper proposes RuWikiBench benchmark for evaluating the analytical capabilities of large language models in generating scientific and encyclopedic texts in Russian. The core of the proposed evaluation system is a three-stage process, consisting of three independent systems that naturally arise when creating articles on a specific topic. Relying on a filtered "Ruwiki" corpus and a clearly defined evaluation methodology, the proposed benchmark establishes a foundation for further research in the application of language models to the task of generating scientific and encyclopedic text in Russian language.

Experiments showed that with a fixed search query, DeepSeek V3 demonstrates the best source selection quality, significantly outperforming BM25 without reranking. At the structure planning stage, it was found that adding a preliminary cluster description consistently improves the quality of outlines for all models, including DeepSeek V3, which demonstrated the best understanding of the process. All models showed comparable section generation quality; however, RuadaptQwen3-32B-Instruct-v2 leads in ROUGE-L and BLEU metrics, indicating a text structure more consistent with the reference. The work shows that models possess significant potential, but their reliable application requires further development of methods for analyzing and structuring review materials.

## Acknowledgements

The study was supported by grant No. 25-11-00191 from the Russian Science Foundation. The work was carried out using the supercomputer "MSU-270" of the Lomonosov Moscow State University.

*This paper is distributed under the terms of the Creative Commons Attribution-Non Commercial 3.0 License which permits non-commercial use, reproduction and distribution of the work without further permission provided the original work is properly cited.*

## References

1. Buckley, C., Voorhees, E.: Evaluating evaluation measure stability. SIGIR Forum (ACM Special Interest Group on Information Retrieval) (10 2000). <https://doi.org/10.1145/345508.345543>
2. Churin, et al.: "long input benchmark for russian analysis.". CoRR (2024)
3. Fenogenova, A., Chervyakov, A., Martynov, N., Kozlova, A., Tikhonova, M., Akhmetgareeva, A., Emelyanov, A., Shevelev, D., Lebedev, P., Sinev, L., Isaeva, U., Kolomeytseva, K., Moskovskiy, D., Goncharova, E., Savushkin, N., Mikhailova, P., Minaeva, A., Dimitrov, D., Panchenko, A., Markov, S.: MERA: A comprehensive LLM evaluation in Russian. In: Ku, L.W., Martins, A., Srikumar, V. (eds.) Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 9920–9948. Association for Computational Linguistics, Bangkok, Thailand (Aug 2024).

<https://doi.org/10.18653/v1/2024.acl-long.534>, <https://aclanthology.org/2024.acl-long.534/>

4. Gusev, I.: Pingpong: A benchmark for role-playing language models with user emulation and multi-model evaluation (2025), <https://arxiv.org/abs/2409.06820>
5. Järvelin, K., Kekäläinen, J.: Cumulated gain-based evaluation of ir techniques. *ACM Trans. Inf. Syst.* **20**, 422–446 (2002), <https://api.semanticscholar.org/CorpusID:1981391>
6. Kang, H., Xiong, C.: Researcharena: Benchmarking large language models’ ability to collect and organize information as research agents (2025), <https://arxiv.org/abs/2406.10291>
7. Lin, C.Y.: ROUGE: A package for automatic evaluation of summaries. In: *Text Summarization Branches Out*. pp. 74–81. Association for Computational Linguistics, Barcelona, Spain (Jul 2004), <https://aclanthology.org/W04-1013/>
8. Liu, A., et al.: Deepseek-v3 technical report. *CoRR* (2024)
9. OpenAI: Introducing deep research. <https://openai.com/index/introducing-deep-research/> (2024), accessed: 2025-07-30
10. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Isabelle, P., Charniak, E., Lin, D. (eds.) *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. pp. 311–318. Association for Computational Linguistics, Philadelphia, Pennsylvania, USA (Jul 2002). <https://doi.org/10.3115/1073083.1073135>, <https://aclanthology.org/P02-1040/>
11. Shao, Y., Jiang, Y., Kanell, T.A., Xu, P., Khattab, O., Lam, M.S.: Assisting in writing wikipedia-like articles from scratch with large language models (2024), <https://arxiv.org/abs/2402.14207>
12. Shavrina, T., Fenogenova, A., Anton, E., Shevelev, D., Artemova, E., Malykh, V., Mikhailov, V., Tikhonova, M., Chertok, A., Evlampiev, A.: RussianSuperGLUE: A Russian language understanding evaluation benchmark. In: Webber, B., Cohn, T., He, Y., Liu, Y. (eds.) *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. pp. 4717–4726. Association for Computational Linguistics, Online (Nov 2020). <https://doi.org/10.18653/v1/2020.emnlp-main.381>, <https://aclanthology.org/2020.emnlp-main.381/>
13. T-Bank: T-bank has opened access to its own russian-language language model in the 7–8 billion parameter weight category. <https://www.tbank.ru/about/news/20072024-t-bank-opened-access-its-own-russian-language-language-model-weight-category-of-7-8-billion-parameters/> (2024), accessed: 2025-08-21
14. Tikhomirov, M., Chernyshov, D.: Facilitating large language model russian adaptation with learned embedding propagation. *Journal of Language and Education* **10**(4), 130–145 (Dec 2024). <https://doi.org/10.17323/jle.2024.22224>, <https://jle.hse.ru/article/view/22224>
15. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., Lample, G.: Llama: Open and efficient foundation language models (2023), <https://arxiv.org/abs/2302.13971>

16. VikhrModels: Rullm arena: Russian llm evaluation benchmark. [https://github.com/VikhrModels/ru\\_llm\\_arena](https://github.com/VikhrModels/ru_llm_arena) (2024), accessed: 2025-07-30
17. Wang, X., Wang, Z., Gao, X., Zhang, F., Wu, Y., Xu, Z., Shi, T., Wang, Z., Li, S., Qian, Q., Yin, R., Lv, C., Zheng, X., Huang, X.: Searching for best practices in retrieval-augmented generation. In: Al-Onaizan, Y., Bansal, M., Chen, Y.N. (eds.) Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing. pp. 17716–17736. Association for Computational Linguistics, Miami, Florida, USA (Nov 2024). <https://doi.org/10.18653/v1/2024.emnlp-main.981>, <https://aclanthology.org/2024.emnlp-main.981/>
18. Wu, J., Ouyang, L., Ziegler, D.M., Stiennon, N., Lowe, R., Leike, J., Christiano, P.: Recursively summarizing books with human feedback (2021), <https://arxiv.org/abs/2109.10862>
19. Yandex: Yandexgpt 5 with reasoning mode. <https://ya.ru/ai/gpt> (2025), accessed: 2025-07-30
20. Yang, A., et al.: Qwen3 technical report (2025), <https://arxiv.org/abs/2505.09388>
21. Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y.: BERTScore: Evaluating text generation with BERT. In: International Conference on Learning Representations (ICLR) (2020), <https://openreview.net/forum?id=SkeHuCVFDr>