# Supercomputing Frontiers and Innovations

## 2025, Vol. 12, No. 3

## Scope

- Future generation supercomputer architectures
- Exascale computing
- Parallel programming models, interfaces, languages, libraries, and tools
- Supercomputer applications and algorithms
- Novel approaches to computing targeted to solve intractable problems
- Convergence of high performance computing, machine learning and big data technologies
- Distributed operating systems and virtualization for highly scalable computing
- Management, administration, and monitoring of supercomputer systems
- Mass storage systems, protocols, and allocation
- Power consumption minimization for supercomputing systems
- Resilience, reliability, and fault tolerance for future generation highly parallel computing systems
- Scientific visualization in supercomputing environments
- Education in high performance computing and computational science

## Editorial Board

Guest Editor's Introduction for Special Issue on

# Supercomputing for Creating, Fine-tuning and Application of Large Language Models

Natalia Loukachevitch, DSc, Research Computing Center of Lomonosov Moscow State University, Moscow, Russia

I am pleased to introduce this special issue of Supercomputing Frontiers and Innovations devoted to Large Language Models, which recently have revolutionized a wide range of domains, from natural language understanding and generation to scientific discovery, coding, and education. Our issue is devoted to the application of LLMs to the analysis of Russian texts and demonstrates LLMs achievements and problems in various tasks.

The submissions can be approximately subdivided into three groups: LLM-based linguistic analysis, LLM-based emotional analysis and reasoning, and LLM-based applications.

In the first group of LLM-based linguistics analysis, Pavel Grashchenkov, Lada Pasko and Regina Nasyrova describe RuParam, a parametric dataset designed to evaluate if LLMs can distinguish between grammatical and ungrammatical sentences from the dataset. Dmitry Morozov, Anna Glazkova and Boris Iomdin evaluate state-of-the-art LLMs, including proprietary and open-weight models, using a prompt-based, few-shot learning approach in the task of Russian morpheme segmentation. Elena Shamaeva, Mikhail Tikhomirov and Natalia Loukachevitch study the possibilities of LLMs for syntactic parsing based on prompts with a single example (one-shot prompting).

In the second group of LLM-based emotion analysis and reasoning, Polina Iaroshenko and Natalia Loukachevitch present a comparative analysis of emotional evaluation of Russian nouns by large language models and native speakers, based on the ENRuN (Emotional Norms for Russian Nouns) database. Ivan Smirnov and colleagues study how to apply large language models for semantic role labeling of Russian emotive predicates through few-shot in-context learning combined with predicate-specific instructions. Elena Sidorova and colleagues address the problem of automatic extraction of argumentative structures in scientific communication texts in Russian. Anna Kuznetsova and colleagues evaluate modern open-source LLMs on the Russian intellectual game "What? Where? When?" They introduce a new dataset of 2600 questions (2018–2025), enriched with empirical human team success rates and annotated with structural and thematic clusters.

In the application group of papers, Denis Grigoriev, Daniil Khudiakov and Daniil Chernyshev present the Russian abstractive summarization dataset RuBookSum for long-form narrative summarization. Anastasia Kolmogorova, Elizaveta Kulikova and Vladislav Lobanov describe a system for the museum review analysis, which uses a two-step pipeline based on LLMs that initially extracts positive and negative keywords about each museum site and subsequently categorizes these keywords into predetermined categories. Nikolai Prokopyev, Marina Solnyshkina and Valery Solovyev address the problem of assessing terminological coherence by evaluating a corpus of textbooks against the Russian Federal State Educational Standard. The authors employ a hybrid methodology combining classical symbolic NLP methods for topic modeling (keyword extraction and term alignment) with qualitative analysis and use of modern large language models for items not found algorithmically.

# Contents

# RuParam: a Russian Parametric Dataset for LLM Evaluation

*Pavel V. Grashchenkov*[1] (iD), *Lada I. Pasko*[1] (iD), *Regina R. Nasyrova*[1] (iD)

We introduce RuParam, a parametric dataset designed to evaluate the acquisition of Russian by large language models (LLMs). This corpus mirrors the structure of the BLiMP family of datasets by containing minimal pairs of sentences. However, our goal was to expand its scope as much as possible by incorporating diverse phenomena from several domains of Russian grammar. A significant portion of the data originates from the Tests of Russian as a Foreign Language (TORFL); similar sources were not previously used for linguistic evaluation of LLMs. Additionally, this study details experimental findings involving six LLMs. These LLMs, sourced from multiple developers, vary in size and pretraining data, which affects their proficiency in Russian. We investigate how effectively these models handle universal, typological, and Russian-specific grammatical features. Our results indicate that while most of the models demonstrate relatively high performance, they struggle significantly with some of the Russian-specific categories.

*Keywords: Large Language Models, linguistic evaluation, minimal pairs, Russian, linguistic parameters, language acquisition.*

## Introduction

In theoretical linguistics, the ability to differentiate between grammatical and ungrammatical sentences (i.e. ones that conform to the rules of grammar of a certain language and ones that do not) has long been considered as a core part of human linguistic competence. A common way of illustrating grammaticality are minimal pairs – sentences that are identical in terms of their lexical content, but differ in whether they violate some grammatical constraint, cf. (1):

(1)     *The bird is singing* vs *\*The bird are singing.*

In recent years, such minimal pairs have moved beyond papers on theoretical linguistics into the field of large language model (LLM) evaluation, with BLiMP [28] being the pioneer in this area. This seems to be a logical step, since if LLMs are to accurately replicate human linguistic behavior, they should be tested on tasks that speakers of natural language are known to succeed at.

Since BLiMP, a lot of work has been done: similar corpora have been developed for various languages. We present **RuParam** – a Russian dataset of minimal pairs[2]. Although there have already been some successful attempts to create grammaticality datasets for Russian – RuCoLA [14] and RuBLiMP [21], we believe that our corpus makes a significant contribution to the field. RuParam addresses some of the shortcomings of other grammaticality corpora, such as semi-automatic data generation, only a small range of linguistic phenomena covered, scarce linguistic annotation, and natural variability in sentence acceptability.

**Our main contributions are as follows:**

- We introduce a new grammaticality dataset of 11,336 minimal pairs in Russian.
  - Our data are classified into 150 linguistic categories covering both universal and language-specific phenomena. The phenomena range from basic concepts of grammar, such as standard cases of subject-verb agreement and case assignment, to more nuanced cases, e.g. clitic placement, licensing of negative polarity items, allomorph distribution and so on.

---

[1]Lomonosov Moscow State University, Moscow, Russian Federation
[2]https://github.com/grapaul/RuParam

– One part of our dataset (8,039 pairs) originates from an independent source of grammaticality minimal pairs – multiple choice tasks from the Test of Russian as a Foreign Language (TORFL). To our knowledge, materials of language proficiency tests for non-native speakers had not previously been used for linguistic evaluation of LLMs.

– The other part of the dataset (3,297 pairs) was taken from Russian corpora (Russian National Corpus, RuConst) and manually modified by trained linguists.

- We evaluate six LLMs on our data using the method of metalinguistic prompting. Although none of the models reach 100% accuracy, some of them are very close to this threshold.

The article is structured as follows. Firstly, we provide some essential background on evaluating large language models using linguistic benchmarks. Secondly, we introduce the RuParam dataset and offer a comprehensive description thereof. Next, we present an experimental study, where RuParam was used to evaluate six different language models – those extensively trained on Russian data and others that were not. Lastly, we analyze and summarize our findings.

## 1. Related Work

BLiMP (Benchmark of Linguistic Minimal Pairs) [28] was the first wide-range dataset to use the grammaticality minimal pair format. BLiMP covers 12 linguistic phenomena of English, for which 67 minimal pair templates were created by linguists. The data were automatically generated using these templates, which enabled the massive size of the dataset (67K minimal pairs). However, this approach has certain limitations due to differences between generated and naturally occurring data. For example, automatically generated data fall short of corpus sentences in both length [5] and structural diversity [26], which makes the evaluation results not entirely representative. Recently, BLiMP-style datasets have been developed for a variety of languages: Chinese (CLiMP [30], SLING [19]), Dutch (BLiMP-NL [20]), Japanese (JBLiMP [16]), Russian (RuBLiMP [21]), Turkish (TurBLiMP [2]), and Urdu (UrBLiMP [1]). MultiBLiMP [12], in turn, covers 101 languages, focusing on just one linguistic phenomenon, namely verb-subject agreement. Some of these benchmarks use the original method of data generation (e.g. CLiPM), while others employ a more naturalistic approach, using examples from corpora as a starting point (e.g. MultiBLiMP, SLING, RuBLiMP, UrBLiMP).

Linguistic acceptability was used in LLM evaluation before the minimal pair approach. The predecessor of BLiMP, CoLA (Corpus of Linguistic Acceptability) [29], created for English, contains individual sentences tagged as either grammatical or ungrammatical; the sentences do not have a counterpart differing in grammaticality. All the sentences in CoLA, along with grammaticality judgements, come from works on theoretical linguistics such as articles and textbooks. As in the case of BLiMP, equivalents for CoLA have been created for many other languages, including Catalan (CatCoLA [3]), Chinese (CoLAC [10]), Danish (DaLAJ [27]), Japanese (JCoLA [17]), Hungarian (HuCoLA [13]), Italian (ItaCoLA [25]), Norwegian (NoCoLA [11]), Russian (RuCoLA [14]) and Spanish (EsCoLA [4]). Importantly for us, some of these datasets – DaLAJ and NoCoLA – use data from the field of second language (L2) acquisition. In both of them, the data come from L2 learner corpora. The ungrammatical sentences are those where L2 learners made mistakes, while grammatical ones are those corrected by native speakers.

In general, sources originating from the educational field have been extensively used for the task of LLM evaluation. MMLU [9] is one prominent example. It aims to evaluate the LLM's knowledge of factual information covering a wide range of subjects; the data stem from multiple-choice questions found in textbooks and examination materials from diverse fields. Among gram-

maticality datasets, RuCoLA includes ungrammatical sentences derived from tasks of Unified State Exam in Russian, aimed at high school graduates. However, the status of such sentences as ungrammatical is doubtful. They rather represent prescriptive norms that are not necessarily part of a native speaker's grammar – otherwise, there would be no point in using them as part of an exam for Russian schoolchildren.

Regarding the procedure of model evaluation for BLiMP and its equivalents, the preference of a model is standardly defined as the difference in the probability of sentences forming a minimal pair. This experimental design serves as a solution to the limitations of CoLA-style evaluation. Since CoLA views acceptability judgement as a binary classification task, it is necessary to train a supervised classifier prior to LLM evaluation. Other factors, besides grammaticality, such as sentence length and word frequency, are inevitably involved. On the contrary, BLiMP allows for separating the grammatical contrast from these additional factors. Furthermore, other approaches have been recently proposed for BLiMP-style data, such as metalinguistic prompting [18]. In this case, LLMs are treated as human subjects of linguistic experiments (see e.g. [6] for an overview of human acceptability judgement task): the prompt consists of an explicit verbal instruction to choose the more acceptable sentence out of the minimal pair. This method allows researchers to investigate whether LLMs have acquired human-like linguistic introspection.

## 2. RuParam

### 2.1. Corpus Structure

RuParam includes 11,336 minimal pairs[3]. The dataset consists of two parts: the first part (8,039 pairs, 70.92%) is based on data from the Test of Russian as a Foreign Language[4] (TORFL); the second part (3,297 pairs, 29.08%) represents a parametric dataset created by linguists and based on naturally occurring sentences. The dataset covers the wide range of linguistic phenomena corresponding to 150 smaller categories. One of the goals of creating RuParam was to maximize the number of diagnostic grammatical features and to diversify the methods for obtaining contrast within each feature. This approach aims to enable not only an overall assessment of linguistic proficiency, but also a detailed analysis of the level of acquisition of specific grammar points.

As the data in the first part come from TORFL materials, it covers phenomena specific (although not necessarily unique) to Russian. The tasks in TORFL were independently created by experts in acquisition of Russian by learners with different native languages. Therefore, this part addresses crucial grammatical phenomena and is particularly relevant for testing multilingual LLMs. The purpose of the second part is twofold. First, it covers universal features (such as projectivity and island constraints) that are characteristic of natural language in general and are therefore absent from TORFL tasks. Second, it includes phenomena that are specific to Russian, but underrepresented in TORFL because of methodological reasons. In the next sections, we discuss the data generation procedure and the phenomena in more detail.

---

[3]The earlier version of RuParam, containing 4,382 minimal pairs, along with the results of LLM evaluation, was presented in [8].

[4] `https://testingcenter.spbu.ru/ru/materials.html`

`https://www.pushkin.institute/certificates/cct/tests-online/?ysclid=meqsg6x5b6171434445`

`https://test.tsu.ru/ru/trki`

## Part 1: TORFL

This part of the data originates from multiple-choice "Vocabulary/Grammar" tasks of TORFL. Each task consists of a sentence (or several sentences) with a gap. There are 3 or 4 options for filling in the gap, only one of which is correct. For example, the task in (2) tests the acquisition of adjective agreement. Only option B forms a grammatical sentence in Russian.

(2)   *Это очень ... здание.*
     this   very      building.N
     'It is a very ... building.'
     *A. высокая*
        tall.F
     *B. высокое*
        tall.N
     *C. высокий*
        tall.M

The minimal pairs were generated by filling in the gap with each option. The correct answer produces the grammatical member of a minimal pair, while the incorrect answers create the ungrammatical ones. The grammatical sentence was paired with all the ungrammatical options, generating two or three minimal pairs per task. Some of the original tasks from TORFL were excluded from the data because they required the use of prior context to choose the correct option (e.g. when a series of tasks represents a coherent text). TORFL covers all CEFR (Common European Framework of Reference) levels of language proficiency: A1, A2, B1, B2, C1, and C2. A1 corresponds to basic knowledge, and C2 is closest to native speaker proficiency. The complexity tags of the original tasks are included in the dataset. The quantity of data per level is shown in Tab. 1.

**Table 1.** Level distribution in TORFL subset of RuParam

| Level | A1 | A2 | B1 | B2 | C1 | C2 |
|---|---|---|---|---|---|---|
| Number of pairs | 1755 | 932 | 2413 | 1588 | 1147 | 204 |
| % | 21.84 | 11.59 | 30.00 | 19.75 | 14.27 | 2.54 |

While the TORFL tasks cover a wide variety of the Russian language phenomena, there is no linguistic annotation available to users. We identified the phenomena used in the tasks and provided annotation for each minimal pair. This was manually done by trained linguists, and each tag was verified by at least one other expert. The annotation included the following parameters, which are divided into 29 categories. We list the most important categories below:
- attributive and predicative agreement;
- lexical selection of all major parts of speech;
- government of different parts of speech;
- use of non-finite forms;
- use of aspect, tense, modality;
- use of coordinating and subordinating conjunctions;
- use of constructions with numerals;
- correctness of using particular parts of speech;
- use of copulas in nominal predications;

- grammaticality of voice forms;
- use of various types of pronouns;
- grammaticality of negative constructions.

Most of these parameters are found in tasks of different levels. Many categories, such as aspect, attributive agreement, and verbal selection, are present in tasks of all six levels.

For some pairs, more than one tag was appropriate. For example, in (3), the ungrammatical form *начинается* 'start.IPF.PRS' differs from the grammatical one *начался* 'start.PF.PST' in both aspect and tense. Such minimal pairs were replicated in the dataset, with only one grammatical category present in the annotation of each instance of the pair.

(3)   *gram*   *Спектакль* **начался** *давно,   вы   уже   опоздали.*
          performance start.PF.PST long ago you already are late
          'The performance started a long time ago, you are already too late.'

      *ungram* *Спектакль* **начинается** *давно,   вы   уже   опоздали.*
          performance start.IPF.PRS  long ago you already are late
          'The performance is starting a long time ago, you are already too late.'

## Part 2: Parametric Dataset

While the approach to the first part of the dataset was data-driven (we annotated pairs created independently for TORFL), the starting point for the second part was grammatical parameters. We used our experience in theoretical linguistics to identify categories that are important for Russian and human language in general, but which are insufficiently covered or not covered at all in TORFL tasks.

The total number of categories in the second part of RuParam is 121. The number of examples varies across different parameters, but each one is represented by at least 15 minimal pairs.
**Universal phenomena** included in the dataset are binding principles; island constraints (coordinate structure, complex NP, adjunct, subject, and others); projectivity.
Other categories are **specific to Russian**. Some of them, such as the directionality of branching, the use of null subjects, or *wh*-word placement (on the left periphery vs *in situ*), represent **parameters of typological variation**. The remaining phenomena covered in the dataset include the following: different types of agreement; clitic placement (P2-clitics, clitics forming conditional clauses); distribution of non-finite forms; licensing of different types of negative polarity items; case of nominal predication; control; voice; depictives; analytical tense forms; analytical comparative and superlative forms; morphophonological variation; double conjunctions; matching in free relatives; constructions with numerals; distribution of the short form of adjectives; genitive marking under negation, and others.

As in many other benchmarks with a similar purpose, grammatical sentences were derived from corpora. We used two corpora of Russian: RuConst [7] and Russian National Corpus, RNC [15]. These sentences were modified by experts in theoretical linguistics to ensure that an ungrammatical counterpart in a minimal pair violated some linguistic constraint. For example, the ungrammatical sentence in (4) is a case of non-projectivity. Although word order in Russian is relatively free, such sentences are ruled out. In (4), the adjective *главного* 'in.chief.GEN' was dislocated so that it is separated from the nominal head that it modifies (*редактора* 'editor.GEN') by another nominal phrase (*смены* 'change.GEN'). Each minimal pair was verified by at least one other expert to confirm that the expected grammatical contrast was present. The total number of errors did not exceed 1%.

(4)  *gram*   *O*   *причинах смены*   **главного**   *редактора не   сообщается.*
      about reasons    change.GEN in chief.GEN editor.GEN not is reported
      'The reasons for the change of the editor-in-chief are not reported.'

      *ungram*   *O*   *причинах* **главного**   *смены*   *редактора не   сообщается.*
      about reasons    in chief.GEN change.GEN editor.GEN not is reported
      Int. 'The reasons for the change of the editor-in-chief are not reported.'

# 3. Evaluation

We assessed the abilities of several LLMs to distinguish between grammatical and ungrammatical sentences in our dataset. The following subsections introduce the models that were tested and describe the experimental setup.

## 3.1. Models

The models for evaluation were chosen based on the following criteria: model size, and accessibility as well as quantity of Russian data used during the training procedure.

**The first group** included closed-source foundation models of large size that were extensively exposed to Russian during training. Consequently, we expected these models to provide the most reliable judgements:

- **GigaChat-2-Max**[5] – the largest and most efficient version of GigaChat models;
- **YandexGPT 5 Pro**[6] – likewise, the most advanced model of the YandexGPT family.

**The second group**, as opposed to the first one, consisted of open-source models with 7-8B parameters, based on the fruitful Qwen2.5 model [22]. Using these models, we aimed to investigate whether the ability to differentiate between sentences in a minimal pair is influenced by a smaller model size, multilingual pre-training, and different adaptation techniques:

- **Qwen2.5-7B-Instruct**[7] [22] – the instruction-tuned multilingual model pre-trained on a large-scale dataset and demonstrating high performance on various tasks from language understanding to coding;
- **T-lite-1.0**[8] – an adaptation of Qwen2.5 model for the Russian language, which was pre-trained in two stages using a combination of Russian and English texts, as well as instruction data, then fine-tuned to follow instructions and preferences;
- **RuadaptQwen2.5-7B**[9] [23, 24] – the modification of Qwen2.5-7B with the tokenizer better suited to the morphology of Russian. The model was also trained on Russian data and with Learned Embedding Propagation procedure.

Moreover, we added **Mistral-7B-Instruct-v0.3**[10] to the comparison, as it is also a multilingual 7B model with different pre-training data. However, it is more English-oriented than other models in terms of pretraining data, which could hinder its performance on our task.[11]

---

[5]https://developers.sber.ru/docs/ru/gigachat/models/updates
[6]https://yandex.cloud/ru/docs/foundation-models/concepts/yandexgpt/models
[7]https://huggingface.co/Qwen/Qwen2.5-7B-Instruct
[8]https://huggingface.co/t-tech/T-lite-it-1.0
[9]https://huggingface.co/RefalMachine/RuadaptQwen2.5-7B-Lite-Beta
[10]https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3
[11]Hereafter we will address the models by the first letter in the title: $G$(igaChat-2-Max), $Y$(andexGPT 5 Pro), $T$(-lite-1.0), $R$(uadaptQwen2.5-7B), $Q$(wen2.5-7B-Instruct), $M$(istral-7B-Instruct-v0.3).

Какое из двух предложений является правильным и грамматичным с точки зрения русского языка?

Предложение 1. <sent_one>

Предложение 2. <sent_two>

Ответь только одной цифрой 1 или 2, ничего не добавляя.

---

Which of the two sentences is correct and grammatical according to the Russian language?

Sentence 1. <sent_one>

Sentence 2. <sent_two>

Respond with only one number 1 or 2, without adding anything.

**Figure 1.** The prompt used for model evaluation

## 3.2. Setup

To evaluate a model's ability to make grammaticality judgements, it was prompted with the instruction in Fig. 1.

The model's response should have been either "1" or "2", denoting the number of the grammatical sentence. As our corpus was designed for diagnostic purposes, there was no training sample, so we only tested the models in zero-shot settings and did not study their abilities in the few-shot setup or after fine-tuning.

LLMs are prone to position bias [31], i.e. they tend to choose the first or the second option regardless of their contents. Therefore, each minimal pair was evaluated twice: with sentences given in their default order, and in the opposite one. The model's response was considered correct only if the model preferred the grammatical option in both iterations. The order of examples within the corpus is arbitrary, different grammatical categories are interleaved so that LLMs do not accumulate guesses about the type of ungrammaticality.

## 4. Results and Discussion

### 4.1. General Results

Table 2 summarizes the results of LLM evaluation. Some minimal pairs were rejected by the models due to ethical considerations, so no answer was given regarding the grammaticality of sentences. This is explained by the source of our data: some corpora examples, especially those coming from news, may contain discussions on sensitive topics, such as politics and health, cf. (5). However, the amount of filtered data was not significant in most cases. The second column of Tab. 2 presents the percentage of correct answers overall, while the third column shows the percentage after the filtered examples were excluded. The other two columns offer statistics for two parts of the dataset separately.

(5)    *Турецкие пограничники задержали судно, перевозившее тонну героина.*
       Turkish    border guards detained    ship    carrying      ton    of heroin
       'Turkish border guards detained a ship carrying a ton of heroin.'

We anticipated that the closed-source large-scale models with substantial exposure to Russian would exhibit superior performance. This expectation was confirmed by G, which had the

**Table 2.** Model evaluation results

| Model | Accuracy | Accuracy (filtered) | TORFL accuracy (filtered) | Parametric data accuracy (filtered) |
|---|---|---|---|---|
| G | 97.85 | 97.92 | 98.05 | 97.54 |
| Y | 95.98 | 97.61 | 97.47 | 97.94 |
| T | 87.79 | 91.06 | 90.81 | 91.66 |
| R | 89.13 | 89.13 | 89.29 | 88.75 |
| Q | 87.31 | 87.31 | 87.24 | 87.50 |
| M | 52.31 | 61.67 | 60.24 | 64.84 |

best results and was closely followed by Y. Although neither of the models reached 100% accuracy, they came close to this threshold. The open-source Qwen-based models (T, R, Q) achieved lower results and differed from each other by approximately two percentage points. M scored significantly lower than all other models, as expected given that it had less Russian pretraining data.

Regarding the difference between the two parts of RuParam, there seems to be none in terms of LLMs performance. The results for the two parts are approximately the same for all models. The only exception is M, which performs significantly better on the parametric part of the data (64.84% vs 60.24%). This may be due to the fact that the second part of the dataset includes some universal phenomena that do not necessarily require extensive pretraining on Russian data.

## 4.2. Results by TORFL Levels

The data in the first part of our dataset coming from TORFL was distributed among six CEFR complexity levels, ranging from A1 to C2. We expect that if the linguistic competence of LLMs is similar to the human one, the results of the models will correlate with the complexity levels: the more difficult the tasks, the lower the accuracy is. This prediction is partially borne out. The models' results are mostly in accordance with the levels, but some exceptions are present – this is shown in Tab. 3 (the results before filtration are given).

**Table 3.** Accuracy by TORFL levels

| | G | Y | T | R | Q | M |
|---|---|---|---|---|---|---|
| A1 | 98.69 | 98.69 | 91.00 | 91.80 | 90.55 | 52.05 |
| A2 | 98.61 | 97.53 | 89.06 | 90.88 | 88.63 | 48.28 |
| B1 | 98.14 | 97.93 | 87.53 | 90.05 | 88.11 | 47.91 |
| B2 | 98.11 | 97.36 | 88.41 | 88.22 | 85.77 | 54.97 |
| C1 | 97.82 | 95.82 | 85.35 | 87.71 | 84.39 | 45.16 |
| C2 | 89.71 | 89.22 | 69.12 | 68.14 | 69.61 | 39.22 |

## 4.3. Results by Linguistic Phenomena

The models demonstrate some common patterns in error frequency. The most complex phenomena come from the second part of the dataset, which consists of data created from corpus examples. Interestingly, most mistakes are made in Russian-specific categories by both groups

of models: industrial ($G$, $Y$) and open-source ($T$, $R$, $Q$, $M$). Table 4 shows the rating of most complex phenomena.

**Table 4.** The most complex linguistic phenomena.
'The category is marked "+" if it is among the top-10 error-prone categories for a model

| Category/Model | G | Y | T | R | Q | M | sum |
|---|---|---|---|---|---|---|---|
| GOV_LOC_1 | + | + | + | + | + | | 5 |
| SUPER_3 | + | + | + | | + | + | 5 |
| COND_1 | | | + | + | + | + | 4 |
| COND_5 | + | + | | | + | + | 4 |
| IMP_VAR | + | + | + | | | + | 4 |
| PREP_VAR | + | + | | + | + | | 4 |
| DISTR | + | + | + | | | | 3 |
| FUT_ASP_1 | | | + | + | | + | 3 |
| FUT_ASP_2 | | | | + | + | + | 3 |
| GOV_LOC_2 | + | | | | + | + | 3 |

Examples and descriptions of the most complex phenomena are given in Tabs. 5, 6, and 7; the data are presented as [*gram/ungram*].

**Table 5.** The description of the most complex linguistic phenomena along with examples illustrating them. Part 1

| | |
|---|---|
| GOV_LOC_1 | Some Russian nouns, such as *шкаф* 'closet', *угол* 'corner', *нос* 'nose' have a special form of the locative case which is required after certain prepositions (e.g. *в* 'in', *на* 'on'). For all other nouns, the prepositional case is used after these prepositions. The grammatical sentences in this category contain a noun in the locative case (e.g. *шкаф-у* 'closet-LOC'), while in their ungrammatical counterparts the regular prepositional case form is used instead (e.g. *шкаф-е* 'closet-PREP'). |
| | *26-летний мужчина прятался в* ***[шкафу/шкафе]****, где девочка* <br> 26-year-old man was hiding in [closet.LOC/closet.PREP] where girl <br><br> *хранит свою одежду.* <br> keeps REFL clothes <br> 'The 26-year-old man was hiding in the closet where the girl keeps her clothes.' |
| SUPER_3 | One of the ways to form the superlative degree form of an adjective in Russian is through the use of the circumfix *наи-...-ейш-*. In ungrammatical sentences, the second part of this circumfix (*-ейш-*) is omitted, while the first part (*наи-*) remains. |
| | *Дебаты – это* ***[наиважнейшая/наиважная]*** *часть* <br> debates COP [SUPER-important-SUPER-F.SG/SUPER-important-F.SG] part <br><br> *избирательного процесса.* <br> electoral process <br> 'Debates are the most important part of the electoral process.' |

As one can see, many complex phenomena are associated with allomorphy. To choose between allomorphs, one needs to know about the properties of individual lexemes (e.g. the presence of a special locative form) and about the context: both the lexical and the phonological features are important. We assume that it is the multifactorial nature of allomorphy that makes the

**Table 6.** The description of the most complex linguistic phenomena along with examples illustrating them. Part 2

| | |
|---|---|
| COND_1 | In counterfactual conditionals, both clauses must include the particle *бы* responsible for the conditional mood. The ungrammatical sentences are produced by eliminating this particle from subordinate clauses headed by *если* 'if'. |
| | *Все было бы нормально, если [**бы**/∅] разговоры подкреплялись* <br> everything was COND normal if [COND/∅] conversations were supported <br> *литературой.* <br> by literature <br> 'Everything would be fine if the conversations were supported by literature.' |
| COND_5 | In Russian, there is a construction with conditional semantics that includes a *wh*-word and the negative particle *ни*. The ungrammatical examples are created by omitting *ни*. |
| | *Как это [**ни**/∅] горько признать, мы еще до покупателя не доехали.* <br> how this [PART/∅] bitter admit we yet to buyer not reached <br> 'As much as I hate to admit it, we have not reached the buyer yet.' |
| IMP_VAR | The singular imperative form is formed by a suffix with two allomorphs: *-и* (*сохран-и* 'save-IMP') and *-∅* (*встань-∅* 'get.up-IMP'). The distribution is determined both phonologically and lexically: *-и* is generally used when a verb has stress on its inflection in the present tense, although there are many exceptions (*прыгн-и* 'jump-IMP', *вытян-и* 'draw-IMP'). The modification of sentences in this category consists of changing the required allomorph to the other one. |
| | *[**Сохрани**/**Сохрань**] свою жизнь ради людей, которые в тебя верят!* <br> [save-IMPER.1/save-IMPER.2] REFL life for people that in you believe <br> 'Save your life for the sake of the people who believe in you!' |
| PREP_VAR | Short prepositions ending with a consonant have an allomorph ending with *-о*, e.g. *с/со* 'with'. The choice of allomorph depends on the phonological conditions: if the word following the preposition begins with the same consonant the preposition ends with, *-о* is inserted (*с мнением* 'with an opinion' vs *со ссылкой* 'with reference'). Likewise, the preposition *о* 'about' has an allomorph that ends with a consonant *об* which is used before vowels (*о мнении* 'about the opinion' vs *об этом* 'about this'), and another lexically selective one *обо* (*о мнении* 'about the opinion' vs *обо мне* 'about me'). In the ungrammatical sentences, an incorrect allomorph of a preposition is used. |
| | *Об этом сообщает РИА Новости [**со**/**с**] ссылкой на режиссера.* <br> about it reports RIA Novosti [with(1)/with(2)] reference to director <br> 'This is reported by RIA Novosti with reference to the director.' |
| DISTR | This category deals with collective predicates (e.g. *пересекаться* 'cross'). Those require the use of a coordinated noun phrase or a plural noun as their subject. The ungrammatical sentences were altered to contain a semantically inappropriate subject. |
| | *[**Наши пути** не **пересекались**. / **Наш путь** не **пересекался**.]* <br> [our.PL path.PL not crossed.PL / our.SG path.SG not crossed.SG] <br> 'Our paths did not cross.' |

models struggle. For instance, models perform significantly better on agreement even though the difference between correct and incorrect forms is often only one character, just as in examples involving allomorphy. This is surprising given that most factors determining the choice of an

**Table 7.** The description of the most complex linguistic phenomena along with examples illustrating them. Part 3

| | |
|---|---|
| FUT_ASP_1 | Future tense in Russian is formed synthetically for perfective verbs (*совершит* 'make.PF.FUT') and analytically for imperfective verbs (*будет совершать* 'will make.IPF.INF'). In pairs of this category, the ungrammatical counterpart has a perfective infinitive in an analytical form instead of the imperfective one. |
| | *Планируется, что автомобиль будет [совершать/совершить]* planned that car will [make.IPF.INF/make.PF.INF] *вертикальный взлет и посадку.* vertical take off and landing 'It is planned that the car will take off and land vertically.' |
| FUT_ASP_2 | Just as in the case of FUT_ASP_1, the ungrammatical sentences in this category have an erroneous use of the analytical future form of a perfective verb. The difference is that the grammatical member of the minimal pair includes a correct synthetic form, but not an analytical one. |
| | *За моральный ущерб мужчина [получит/будет получить] 100 рублей.* for moral damage man [receive.PF.FUT/will receive.PF.INF] 100 rubles 'The man will receive 100 rubles for moral damage.' |
| GOV_LOC_2 | Similarly to GOV_LOC_1, this category deals with the locative/prepositional case distinction. The ungrammatical sentences demonstrate incorrect use of the locative form. |
| | *O красном [снеге/снегу] сообщали жители нескольких районов* about red [snow.PREP/snow.LOC] reported residents several districts *области.* region 'Residents of several districts of the region reported red snow.' |

allomorph are local (e.g. the first character of the next word), while by agreement the goal and the probe can be located at a considerable linear distance. However, in the case of agreement, it is mostly a single factor that matters: the morphological form of the goal (e.g. the verb form has to correspond to the morphological features of the noun in the nominative case).

## Conclusion

We introduced a new dataset for testing the competence of LLMs in Russian. Our corpus uses the minimal pair framework that is now common for the task of linguistic evaluation of LLMs. The data in this benchmark come from two sources: tasks of the Test of Russian as a Foreign Language and corpora of Russian. The first type of data is novel to the field, while the second one has been extensively used in similar datasets and is preferable to automatic data generation. Ungrammatical sentences in the first part of the dataset were independently created by the L2 acquisition experts, while those in the second part were manually generated specifically for the corpus by trained linguists. The dataset contains fine-grained linguistic annotation covering a wide range of categories. Some of them are universal, while others represent Russian-specific phenomena. All annotations were performed by experts in theoretical linguistics.

We evaluated six LLMs on our dataset. To do this, we used the metalinguistic prompting method, treating the models as if they were human subjects in linguistic studies. We found that large-scale models trained extensively on Russian demonstrate very high performance levels, close

to the 100% threshold. As it was expected, smaller open-source models achieved lower results. The evaluation results on the TORFL part of our dataset show that the degree of success of LLMs mostly correlates with CEFR complexity levels of the tasks. This finding demonstrates one similarity between linguistic competence of LLMs and humans. Although models generally achieve relatively good results, some categories proved to be problematic. These are Russian-specific phenomena from the part of the dataset generated using corpus data. Many of these phenomena involve morphophonological variation and the constraints on analytical forms. One of the most important findings is that models from different origins converge on exactly the same types of errors. This may not be a coincidence and could reveal insights into differences in linguistic competence of LLMs and human Russian speakers.

To conclude, RuParam is a novel, carefully designed source of data on the Russian language. We hope that it will be useful for further investigation into LLMs' grammar, as well as for model development.

## Limitations

- For evaluation we adopt Large Language Models, including foundational models, which constantly undergo modifications. Hence, later evaluations may differ from the results presented in the paper.
- While we present several findings that shed light on the common patterns in the linguistic competence of LLMs, further analysis is required to explain the reasons behind them. For instance, we assume that tokenization may affect the complexity of allomorphy. In addition, the analysis may be enriched by the data from other languages apart from Russian.

## Acknowledgements

## References

1. Adeeba, F., Dillon, B., Sajjad, H., Bhatt, R.: UrBLiMP: A Benchmark for Evaluating the Linguistic Competence of Large Language Models in Urdu (2025), `https://arxiv.org/abs/2508.01006`

2. Başar, E., Padovani, F., Jumelet, J., Bisazza, A.: TurBLiMP: A Turkish Benchmark of Linguistic Minimal Pairs. In: Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing. p. 16506–16521. Association for Computational Linguistics, Suzhou, China (2025). `https://doi.org/10.34810/data1393`

3. Bel, N., Punsola, M., Ruiz-Fernández, V.: CatCoLA, Catalan Corpus of Linguistic Acceptability. Procesamiento del Lenguaje Natural 73, 177–190 (2024). `https://doi.org/10.34810/data1393`

4. Bel, N., Punsola, M., Ruiz-Fernández, V.: EsCoLA: Spanish corpus of Linguistic Acceptability. In: Calzolari, N., Kan, M.Y., Hoste, V., *et al.* (eds.) Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024). pp. 6268–6277. ELRA and ICCL, Torino, Italia (May 2024). `https://doi.org/10.34810/data1138`

5. Daultani, V., Martínez, H.J.V., Okazaki, N.: Acceptability Evaluation of Naturally Written Sentences. Journal of Information Processing 32, 652–666 (2024). `https://doi.org/10.2197/ipsjjip.32.652`

6. Featherston, S.: Response Methods in Acceptability Experiments, p. 39–61. Cambridge Handbooks in Language and Linguistics, Cambridge University Press (2021). `https://doi.org/10.1017/9781108569620`

7. Grashchenkov, P.: RuConst: A Treebank for Russian. Lomonosov Philology Journal. Series 9. Philology 3, 94–112 (2024). `https://doi.org/10.55959/MSU0130-0075-9-2024-47-03-7`, (in Russian)

8. Grashchenkov, P., Pasko, L., Studenikina, K., Tikhomirov, M.: Russian parametric corpus RuParam. Scientific and Technical Journal of Information Technologies, Mechanics and Optics 24(6), 991–998 (2024). `https://doi.org/10.17586/2226-1494-2024-24-6-991-998`, (in Russian)

9. Hendrycks, D., Burns, C., Basart, S., *et al.*: Measuring Massive Multitask Language Understanding. Proceedings of the International Conference on Learning Representations (ICLR) pp. 11260–11285 (2021). `https://doi.org/10.18653/v1/2024.findings-acl.671`

10. Hu, H., Zhang, Z., Huang, W., *et al.*: Revisiting acceptability judgements (05 2023). `https://doi.org/10.48550/arXiv.2305.14091`

11. Jentoft, M., Samuel, D.: NoCoLA: The Norwegian Corpus of Linguistic Acceptability. In: Alumäe, T., Fishel, M. (eds.) Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa). pp. 610–617. University of Tartu Library, Tórshavn, Faroe Islands (May 2023), `https://aclanthology.org/2023.nodalida-1.60/`

12. Jumelet, J., Weissweiler, L., Bisazza, A.: MultiBLiMP 1.0: A Massively Multilingual Benchmark of Linguistic Minimal Pairs (04 2025). `https://doi.org/10.48550/arXiv.2504.02768`

13. Ligeti-Nagy, N., Ferenczi, G., Héja, E., *et al.*: HuLU: Hungarian Language Understanding Benchmark Kit. In: Calzolari, N., Kan, M.Y., Hoste, V., *et al.* (eds.) Proceedings of the

2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024). pp. 8360–8371. ELRA and ICCL, Torino, Italia (May 2024), `https://aclanthology.org/2024.lrec-main.733/`

14. Mikhailov, V., Shamardina, T., Ryabinin, M., *et al.*: RuCoLA: Russian Corpus of Linguistic Acceptability. In: Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. p. 5207–5227. Association for Computational Linguistics (2022). `https://doi.org/10.18653/v1/2022.emnlp-main.348`

15. Savchuk, S.O., Arkhangelskiy, T., Bonch-Osmolovskaya, A.A., *et al.*: Russian National Corpus 2.0: New opportunities and development prospects. Voprosy Jazykoznanija 2, 7–34 (2024). `https://doi.org/10.31857/0373-658X.2024.2.7-34`, (in Russian)

16. Someya, T., Oseki, Y.: JBLiMP: Japanese Benchmark of Linguistic Minimal Pairs. In: Vlachos, A., Augenstein, I. (eds.) Findings of the Association for Computational Linguistics: EACL 2023. pp. 1581–1594. Association for Computational Linguistics, Dubrovnik, Croatia (May 2023). `https://doi.org/10.18653/v1/2023.findings-eacl.117`

17. Someya, T., Sugimoto, Y., Oseki, Y.: JCoLA: Japanese Corpus of Linguistic Acceptability. In: Calzolari, N., Kan, M.Y., Hoste, V., Lenci, A., Sakti, S., Xue, N. (eds.) Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024). pp. 9477–9488. ELRA and ICCL, Torino, Italia (May 2024), `https://aclanthology.org/2024.lrec-main.828/`

18. Song, S., Hu, J., Mahowald, K.: Language Models Fail to Introspect About Their Knowledge of Language (2025), `https://arxiv.org/abs/2503.07513`

19. Song, Y., Krishna, K., Bhatt, R., Iyyer, M.: SLING: Sino Linguistic Evaluation of Large Language Models. In: Goldberg, Y., Kozareva, Z., Zhang, Y. (eds.) Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. pp. 4606–4634. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates (Dec 2022). `https://doi.org/10.18653/v1/2022.emnlp-main.305`

20. Suijkerbuijk, M., Prins, Z., Kloots, M.d.H., *et al.*: BLiMP-NL: A Corpus of Dutch Minimal Pairs and Acceptability Judgments for Language Model Evaluation. Computational Linguistics. P. 1–35 (05 2025). `https://doi.org/10.1162/coli_a_00559`

21. Taktasheva, E., Bazhukov, M., Koncha, K., *et al.*: RuBLiMP: Russian Benchmark of Linguistic Minimal Pairs. In: Al-Onaizan, Y., Bansal, M., Chen, Y.N. (eds.) Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing. pp. 9268–9299. Association for Computational Linguistics, Miami, Florida, USA (Nov 2024). `https://doi.org/10.18653/v1/2024.emnlp-main.522`

22. Team, Q.: Qwen2.5: A party of foundation models (September 2024), `https://qwenlm.github.io/blog/qwen2.5/`

23. Tikhomirov, M., Chernyshev, D.: Impact of Tokenization on LLaMa Russian Adaptation. In: 2023 Ivannikov Ispras Open Conference (ISPRAS). pp. 163–168 (2023). `https://doi.org/10.1109/ISPRAS60948.2023.10508177`

24. Tikhomirov, M., Chernyshev, D.: Facilitating large language model Russian adaptation with Learned Embedding Propagation. Journal of Language and Education 10(4), 130–145 (2024). `https://doi.org/10.17323/jle.2024.22224`

25. Trotta, D., Guarasci, R., Leonardelli, E., Tonelli, S.: Monolingual and Cross-Lingual Acceptability Judgments with the Italian CoLA corpus. In: Findings of the Association for Computational Linguistics: EMNLP 2021. p. 2929–2940. Association for Computational Linguistics (2021). `https://doi.org/10.18653/v1/2021.findings-emnlp.250`

26. Vázquez Martínez, H.J., Heuser, A., Yang, C., Kodner, J.: Evaluating Neural Language Models as Cognitive Models of Language Acquisition. In: Hupkes, D., Dankers, V., Batsuren, K., *et al.* (eds.) Proceedings of the 1st GenBench Workshop on (Benchmarking) Generalisation in NLP. pp. 48–64. Association for Computational Linguistics, Singapore (Dec 2023). `https://doi.org/10.18653/v1/2023.genbench-1.4`

27. Volodina, E., Mohammed, Y.A., Klezl, J.: DaLAJ - a dataset for linguistic acceptability judgments for Swedish. In: Proceedings of the 10th Workshop on NLP for Computer Assisted Language Learning. p. 28–37. LiU Electronic Press (2021), `https://aclanthology.org/2021.nlp4call-1.3/`

28. Warstadt, A., Parrish, A., Liu, H., *et al.*: BLiMP: The Benchmark of Linguistic Minimal Pairs for English. Transactions of the Association for Computational Linguistics 8, 377–392 (07 2020). `https://doi.org/10.1162/tacl_a_00321`

29. Warstadt, A., Singh, A., Bowman, S.R.: Neural Network Acceptability Judgments. Transactions of the Association for Computational Linguistics 7, 625–641 (09 2019). `https://doi.org/10.1162/tacl_a_00290`

30. Xiang, B., Yang, C., Li, Y., *et al.*: CLiMP: A Benchmark for Chinese Language Model Evaluation. In: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume. pp. 2784–2790. Association for Computational Linguistics, Online (Apr 2021). `https://doi.org/10.18653/v1/2021.eacl-main.242`

31. Zheng, L., Chiang, W.L., Sheng, Y., *et al.*: Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (eds.) Advances in Neural Information Processing Systems. vol. 36, pp. 46595–46623. Curran Associates, Inc. (2023). `https://doi.org/10.5555/3666122.3668142`

# Large Language Models versus Native Speakers in Emotional Assessment of Russian Words

*Polina V. Iaroshenko*[1] (iD), *Natalia V. Louckachevitch*[1] (iD)

The paper presents a comparative analysis of emotional evaluation of Russian nouns by large language models and native speakers. Based on the ENRuN (Emotional Norms for Russian Nouns) database, which contains ratings of 1,800 nouns across five basic emotions (happiness, sadness, anger, fear, and disgust), the research compares human assessments with evaluations provided by seven large language models (Llama-3-70B, Qwen 2.5-32B, YandexGPT 5 Lite, RuadaptQwen2.5-7B, RuadaptQwen2.5-32B-Pro-Beta, T-pro, T-lite). Although some models demonstrated relatively high correlation with human assessments, persistent systematic deviations were observed across all tested models. The analysis reveals significant differences in emotional perception during word evaluation: the models demonstrate a tendency to hyperbolise negative emotions and show variability in assessing positive emotions, particularly when analysing words related to sensitive topics (violence, religion, obscene vocabulary). The findings indicate that the closest alignment with human evaluations is achieved when there is a balance between the model's size and the quality of its language adaptation.

*Keywords: Large Language Models, human-likeness, emotional intelligence, Russian language.*

## Introduction

Currently, large language models (LLMs) are one of the key tools for addressing numerous NLP tasks. One of the relevant research directions in the field is the emotion analysis [2]. Within this framework, questions regarding the emotional intelligence of LLMs are gaining increasing prominence in the research community [7]. Studies in this area combine both applied and fundamental aspects. From a practical perspective, the emotional alignment of LLMs is essential, for instance, to enhance the quality of communication between humans and AI assistants, which are actively employed in various domains (medicine, education, entertainment, etc.), as well as for using language models in annotating emotion-related data (see, for example, the review by [13]). From a theoretical standpoint, research interest lies in analysing the emotional behaviour of LLMs in various situations, understanding how language models process emotions, and comparing human emotional reactions with those of LLMs [8].

Thus, the need to study the emotional behaviour of LLMs and compare it with human responses is increasing. This includes the pertinent question of how models' emotional behaviour varies depending on specific languages and value systems characteristic of their native speakers. The reproduction of certain biases by LLMs has been repeatedly noted in numerous studies (such as geopolitical [15] or gender stereotypes, particularly concerning the emotional behaviour of men and women [3]).

The aim of the study is to compare emotional assessments of Russian words by native speakers against assessments of the same words by LLMs. With this aim, we utilise the ENRuN (Emotional Norms for Russian Nouns) database [16] – a dataset comprising emotional ratings for 1,800 Russian nouns provided by native Russian speakers.

The article is organized as follows. Section 1 presents a brief overview of the related work. Section 2 describes the data utilised in the study. In Section 3, we outline the methodological framework adopted for the study. Section 4 details the process of prompt engineering and hyper-

---

[1] Research Computing Center, Lomonosov Moscow State University, Moscow, Russian Federation

parameter tuning. Section 5 presents the main findings. Conclusion summarizes the study and points directions for further work.

# 1. Related Work

Earlier studies focusing on emotions expressed in text primarily addressed issues related to emotion recognition and classification [1, 4, 5, 9].

With the advancement of LLMs, the challenge has evolved beyond mere emotion recognition to include the production of appropriate emotional responses by models. Consequently, numerous studies have been conducted to examine the emotional intelligence of LLMs [18]. New types of benchmarks aimed at assessing LLMs' emotional intelligence are emerging. The work of [14] introduces EmoBench (available in English and Chinese), which encompasses tasks in two areas: Emotional Understanding and Emotional Application. Both areas provide LLMs with brief descriptions of life scenarios; however, the first area requires identifying emotions and their causes (emphasising complex, emotionally ambiguous situations), while the second one demands selecting situation-appropriate responses. An example scenario from the Emotional Understanding section reads: "After a long day of terrible events, Sam started laughing hysterically when his car broke down". The authors' experiment revealed a significant gap between the performance of LLMs tested on EmoBench and the responses of human participants. The work of [6] presents EmotionQueen, a benchmark for assessing LLM empathy. In this benchmark, LLMs are tasked with responding to statements containing information about various life events. The authors note that while earlier research focused primarily on recognising explicit emotions, the field of measuring LLMs' emotional intelligence now concentrates on more profound and complex analysis, including the understanding of implicit emotions not directly expressed in user statements, or mixed emotions in situations involving multiple events with different emotional connotations. The results of this EmotionQueen experiment, however, demonstrated that some LLMs, particularly LLaMA2-70B and Claude2, can surpass human levels of empathy.

LLMs' emotional intelligence is frequently evaluated using psychometric tests designed for humans. For example, in Dalal et al. 2025, LLMs' emotional intelligence was assessed using the Situational Test of Emotional Understanding (STEU) [12], where respondents are presented with situation descriptions and must explain what feelings a person should experience under these circumstances. The study [7] revealed that LLMs deviated from reference answers (human responses) in 33% of cases. It was also noted that in several instances, the models offered reasonable alternative emotional assessments for various situations. In [8], a dataset describing various situations was created to evaluate LLMs' empathetic capabilities, with models being required to assess these situations in terms of the emotions they evoke. Human responses served as the gold standard. The researchers observe that while the models' reactions to the proposed situations can generally be characterised as appropriate, none of the tested LLMs demonstrated results sufficiently close to human references.

In summary, the issue of emotional alignment of LLMs is becoming increasingly significant. The described studies mainly examine the adequacy of LLM responses to various life circumstances. The present study also aims to compare LLM responses with those of human participants; however, instead of situation descriptions or utterances, Russian nouns will serve as stimuli for the models.

## 2. Data

The current version[2] of the ENRuN (Emotional Norms for Russian Nouns) database [16] contains emotional ratings of 1,800 Russian nouns. Each word was rated within the dimensional (valence and arousal) and categorical approach (happiness, sadness, anger, fear, and disgust). For each word, the mean values, standard deviations, minimum and maximum scores for each parameter, and the number of people[3] who rated the word are presented.

The ENRuN word list was compiled based on the frequency dictionary of the Russian language [10], with lexical items selected according to several formal criteria (such as word length, frequency, etc.). The lexical composition of the list is notably diverse: it includes both neutral words ("magazine", "calculator", "marble") and sensitive terms related to health, religion, or moral values ("alcoholism", "atheism", "looting").

For this study, we used averaged respondent ratings obtained through a categorical approach survey, measuring the degree of association between words and specific emotions (happiness, sadness, anger, fear, and disgust) on a five-point scale (see Tab. 1 for several examples).

**Table 1.** Example of the ENRuN Data Presentation

| Word | Happy | Sad | Anger | Fear | Disgust |
|------------|-------|-------|-------|-------|---------|
| Professor | 1.615 | 0.462 | 0.231 | 0.885 | 0.385 |
| Friendship | 4.524 | 1.524 | 0.476 | 0.476 | 0.143 |
| Garbage | 0.043 | 0.391 | 0.913 | 0.783 | 4.174 |

Thus, we observe that the word "professor" is relatively neutral and does not trigger strong emotional associations from respondents, while the word "friendship" is rated as joyful, and the word "garbage" ("pomoika") shows a high rating for the emotion of disgust.

An earlier, publicly available version of the database [11], containing ratings for 378 words, presents the instruction given to respondents during the categorical approach experiment:

*"Please rate using the scale from 0 to 5 to which extent, in your opinion, each word is related to emotions of happiness, fear, disgust, anger, and sadness. You will have to fill out the tables below. Words are in the rows and emotions are in the columns. If you think that the given word is not related at all to the given emotion, write "0". If you think that the given word is very much related to the given emotion, write "5". You can also use all the intermediate values of this scale. You have to give five ratings for each scale indicating as to how strongly the given word is related to happiness (1st row), fear (2nd row), disgust (3rd row), anger (4th row), and sadness (5th row). If necessary, you can give high ratings in several columns for the same word".*

This instruction, originally formulated for human respondents, will serve as the basis for developing prompts for LLMs.

## 3. Methodology

The core idea of the experiment is to task large language models with assessing words from the ENRuN database in terms of their associations with emotions (happiness, sadness, anger,

---

[2]The current version of the database can be provided to researchers upon request.

[3]The database development is an ongoing process. The current analysis incorporates responses from a sample of 692 participants at the time of manuscript submission.

fear, and disgust) on a five-point scale. This approach yields results that can be compared with human assessments available in the ENRuN database.

To compare assessments between human respondents and LLMs, the following metrics were employed: Pearson correlation coefficient, Spearman correlation coefficient, and standardised difference. Human assessment serves as the reference standard in this case. The Pearson correlation coefficient helps determine how accurately LLMs reproduce general trends in emotional word assessments. The Spearman correlation coefficient is included in the analysis as it is less sensitive to outliers and non-linear associations, which is crucial when working with emotional assessments, where the association between human and model ratings may be non-linear. The standardised difference (Std Diff) was selected to quantify absolute differences between LLM and human responses, enabling the identification of consistent discrepancies in LLM assessments. This comprehensive approach provides a more complete picture of how successfully LLMs can reproduce human assessments of words' emotional content.

Various categories of LLMs were selected for the study: models from Russian developers (YandexGPT 5 Lite[4]), including adapted models (T-lite-it-1.0[5], T-pro-it-1.0[6] – from the Qwen 2.5 family; RuadaptQwen2.5-7B[7] – adaptation to Russian of T-lite-it-1.0, RuadaptQwen2.5-32B-Pro-Beta[8] – adaptation to Russian of T-pro-it-1.0), as well as models of foreign origin (Qwen 2.5[9], Llama-3[10]).

The models selected also differ in their parameter count and include the following: small models (7-8B) – YandexGPT 5 Lite, T-lite, RuadaptQwen2.5-7B; medium-sized models (32B) – T-pro, Qwen 2.5, RuadaptQwen2.5-32B-Pro-Beta; and a large model (70B) – Llama-3.

This selection of models enables evaluation of whether the alignment between model and human responses correlates with language adaptation or parameter count.

In the preparatory phase of the experiment, only the base model was used. RuadaptQwen2.5-32B-Pro-Beta, specifically adapted for Russian [17], was selected as the base model for the study. The base model was used to identify the most effective prompt and optimal hyperparameters for collecting emotional word assessments. Throughout the experiment, emotional assessments were collected from all models using the selected prompt and hyperparameters. The obtained model responses were compared both with each other and with human assessments of the nouns.

## 4. Experimental Settings

**Prompt Selection.** The prompt for this task was developed based on the instruction given to respondents who participated in word assessment for the ENRuN database; the complete instruction text is provided in Section 2.

Three prompt variants were tested: "min", "base", and "detailed". The "min" variant was the shortest, containing only the task description without additional information. The "base" variant included both the task and a brief description of the role the model should assume when answering questions. The "detailed" variant contained the most comprehensive role description, emphasising internal motivation and the significance of survey participation.

---

[4]https://huggingface.co/yandex/YandexGPT-5-Lite-8B-instruct
[5]https://huggingface.co/AnatoliiPotapov/T-lite-instruct-0.1
[6]https://huggingface.co/t-tech/T-pro-it-1.0
[7]https://huggingface.co/RefalMachine/RuadaptQwen2.5-7B-Lite-Beta
[8]https://huggingface.co/RefalMachine/RuadaptQwen2.5-32B-Pro-Beta
[9]https://huggingface.co/Qwen/Qwen2.5-32B-Instruct
[10]https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct

Below is the complete text of the "detailed" prompt.

*Assume the ROLE and complete the TASK.*

*ROLE:*

*You are an ordinary person who speaks Russian and lives in Russia. You have been invited to participate in an experiment by scientists from the Laboratory of Cognitive Research. The experiment is conducted to study how Russian native speakers evaluate various words in terms of their emotional content. You are very interested in participating in the research. You answer questions attentively, focusing intently and sincerely. You understand that your responses are crucial for the experiment.*

*TASK:*

*Please rate on a scale from 0 to 5 how much, in your opinion, the word for evaluation is associated with emotions such as happiness, fear, disgust, anger, and sadness. If you think the word is not at all associated with a given emotion, assign "0"; if you believe the word is very strongly associated with the emotion, assign "5". You may also use all intermediate values on this scale. You can use any decimal values between 0 and 5 (for example, 2.5, 3.7, 4.8, etc.). Thus, each word requires 5 ratings: how much it is associated with happiness, sadness, anger, fear, and disgust. If necessary, you may assign high ratings in several emotion categories for the same word or assign zero values across all categories if the word evokes no emotions for you.*

*The answer should include only five numerical ratings for emotions separated by spaces in the following order: first rating for HAPPINESS, second for FEAR, third for DISGUST, fourth for ANGER, fifth for SADNESS. The answer should NOT include additional comments.*

The "min" prompt does not include the ROLE section, while the "base" prompt includes a condensed version of the role description: "You are a person, a Russian native speaker participating in a psychological experiment". The TASK section remains identical across all prompts. Model responses for each of the three prompt variants were compared with human assessments from the ENRuN database using the following metrics: Pearson correlation coefficient, Spearman correlation coefficient, and Std Diff. In selecting the optimal prompt, we aimed for minimal Std Diff values alongside high Pearson and Spearman correlations. The testing revealed that the "detailed" variant proved most effective, as it achieved the highest Pearson coefficient and lowest Std Diff, while only marginally falling behind the "base" prompt in Spearman coefficient, with the difference being insignificant (see Tab. 2).

**Table 2.** Evaluation of Different Prompt Variants

| Prompt | Pearson | Spearman | Std Diff |
|---|---|---|---|
| min | 0.56 | 0.41 | 1.11 |
| base | 0.64 | **0.57** | 1.07 |
| detailed | **0.67** | 0.55 | **1.01** |

**Hyperparameters**. When selecting hyperparameters, particular attention was paid to the temperature parameter. At higher temperature values (0.4–0.5), significant deviations from the given instructions were observed. This was especially characteristic when evaluating words semantically related to typical language model tasks ("enumeration", "decoding", "scale", "format", "input", etc.). In such cases, instead of evaluating the proposed word according to specified emotion categories within the established scale, the model showed a tendency to generate and

evaluate its own words. At low temperature values (0.1–0.3), the opposite problem was observed: the model demonstrated a tendency towards insufficient variability, assigning uniform ratings across all emotion categories (for example, predominantly using a rating of "3" while ignoring other scale values).

To achieve an optimal balance between excessive variability and undesirable uniformity in responses, a temperature value of 0.2 was selected in combination with a repetition penalty of 1.2, which helped minimise the aforementioned issues.

**Response Variability**. To test the model's response variability using the "detailed" prompt, three series of responses were collected. Analysis of consistency between series showed a moderately high degree of response reproducibility: the average Pearson correlation coefficient between the three iterations was 0.71, and the Spearman correlation coefficient was 0.66. The standard difference between response distributions (Std Diff = 1.03) indicates a certain variability in model responses. Words with the highest degree of variation in ratings between iterations were identified. To investigate the most notable rating discrepancies, an additional prompt was developed. This prompt retained the role description from the "detailed" version but included a new task: to explain the reasoning behind the assigned ratings. The explanations provided by the model could be characterised as generally appropriate in most cases. Consider the example of the word "nostalgia", which received the highest score in the "Sadness" category in the first iteration and in the "Anger" category in the second iteration (see Tab. 3).

**Table 3.** A word example of inter-iteration rating variability in model assessments

| Iteration | Sadness | Anger | Model's explanation |
|---|---|---|---|
| 1 | 4.5 | 0 | Sadness (4.5/5): the primary emotional colouring of nostalgia is indeed sadness. It is a feeling of loss of something precious and unique that can never be recovered. Hence, the oppressive note of sadness is quite strongly present. |
| 2 | 0 | 4.5 | Anger (4.5/5): the high rating for anger is attributed to nostalgia's capacity to evoke irritation or even rage regarding how much everything has changed, what opportunities were missed, or losses that had to be endured. This dissatisfaction with the current state of affairs, viewed through the lens of an idealised perception of the lost period, triggers strong emotions of discontent with oneself or circumstances. |

In this case, some variability in model ratings may be a positive factor, as it partially mimics differences in word assessments by actual respondents. To test this hypothesis, human assessments were compared with each of the three response series separately, and then with the averaged indicators across the three series (see Tab. 4).

Thus, averaging the results of three iterations demonstrates more balanced and stable results. Although the first iteration shows higher correlations individually, averaging reduces the standard deviation. For further comparison of model responses with those of actual respondents, the average ratings across three iterations were used.

**Table 4.** Model-human rating differences by iteration and mean values

| Iteration | Pearson | Spearman | Std Diff |
|---|---|---|---|
| 1 | 0.69 | 0.59 | 1.02 |
| 2 | 0.55 | 0.46 | 1.16 |
| 3 | 0.56 | 0.46 | 1.15 |
| Mean (3 iter.) | 0.67 | 0.57 | 0.93 |

# 5. Experiment: Comparison of Responses from Different LLMs

Using the "detailed" prompt, three series of responses were collected from each model, followed by obtaining averaged responses. It should be noted that Llama-3 refused to evaluate 4 words from the proposed list: three were instances of obscene language, and one was a colloquial term for an infectious disease (gonorrhoea). The reasons given for refusal were inappropriate vocabulary in the case of obscene words, and inability to provide recommendations regarding illegal content in the case of the disease term. Notably, while the ENRuN database word list contained other obscene words and disease terms (such as syphilis), Llama-3 did provide ratings for these.

For comparing Llama-3's responses with human assessments and other models, zero values were assigned across all emotions for the words it refused to evaluate.

**Comparison of LLMs with Human Assessment**. The averaged responses were compared with human assessments of nouns from the ENRuN database. The comparison results are presented in Tab. 5).

**Table 5.** Evaluation of LLM Responses Compared to Human Ratings

| Model | Pearson | Spearman | Std Diff |
|---|---|---|---|
| RuadaptQwen2.5-32B-Pro-Beta | **0.67** | 0.57 | **0.93** |
| YandexGPT 5 Lite | 0.62 | **0.58** | 1.05 |
| T-pro | 0.55 | 0.47 | 1.08 |
| Qwen 2.5-32B | 0.57 | 0.48 | 1.15 |
| Llama-3 | 0.61 | 0.55 | 1.16 |
| RuadaptQwen2.5-7B | 0.41 | 0.33 | 1.18 |
| T-lite | 0.35 | 0.29 | 1.22 |

RuadaptQwen2.5-32B-Pro-Beta demonstrates the highest correlation with human responses, showing the highest Pearson coefficient (0.66) and one of the highest Spearman coefficients (0.56). This is further confirmed by the lowest standard deviation (0.93) among all models. YandexGPT 5 Lite shows the second-best result with Pearson coefficient of 0.62 and Spearman coefficient of 0.58, indicating good alignment with human responses. T-lite shows the lowest correlation values (Pearson: 0.35, Spearman: 0.29) and the highest standard deviation (1.22), indicating substantial divergence from human assessments.

When comparing human assessments with LLM responses, the following trends were identified across emotion categories:

- positive emotions are represented by a single class – "happiness". This emotion shows the greatest variability between models. Most models tend to overestimate ratings compared to humans: YandexGPT 5 Lite (ratings higher than human assessments in 83.67% of cases),

RuadaptQwen2.5-7B (85.83%), T-lite (76.33%), Llama-3 (73.39%). However, other models demonstrate the opposite tendency, underestimating ratings compared to humans: T-pro (73.83%), Qwen2.5-32B (64.89%), RuadaptQwen2.5-32B-Pro-Beta (59.67%).;

- in assessing negative emotions, the greatest consistency is observed in the "fear" category, although most models still tend to underestimate it in more than half of cases: T-lite (61.33%), Llama-3 (68.28%), Qwen2.5-32B (66.44%);

- for "disgust" and "anger" categories, there is a tendency to overestimate ratings compared to humans. For "disgust": the most pronounced overestimation is shown by RuadaptQwen2.5-7B (87.94%) and T-lite (79.44%). For "anger", the most pronounced overestimation is recorded in RuadaptQwen2.5-7B (87.78%), T-lite (84.39%), and T-pro (81.56%).

Thus, RuadaptQwen2.5-7B, YandexGPT 5 Lite, and T-lite tend to systematically overestimate across most emotions, while Llama-3 shows the opposite tendency, more frequently underestimating. RuadaptQwen2.5-32B-Pro-Beta demonstrates the most balanced assessments. Qwen2.5-32B and T-pro show mixed patterns with a predominance of underestimation for positive emotions.

Based on the comparison results between model and human respondent answers, 10 words were identified for each model, where assessments showed maximum absolute difference compared to human ratings. In total, 42 unique words appeared in the top-10 lists across different LLMs. Analysis of this vocabulary revealed the prevalence of certain semantic categories: religion – "funeral service" ("otpevanie"), "Satan"; manifestations of violence – "torture", "suffering", "slaughter" ("boinya"), "slap" ("poshchechina"); taboo subjects (obscene language, vocabulary related to narcotic or poisonous substances, immoral activities). The identified vocabulary has predominantly negative connotations. Among the most frequent words, "Satan" appears in 6 out of 7 models, and words such as "villainy", "downfall", and "torture" each appear in 4 models' lists. In the case of "Satan", the general tendency to overestimate negative emotional classes is confirmed. Consistent overestimation by all models is recorded for "sadness" and "disgust" classes, meaning models tend to consider the word "Satan" much sadder and more disgusting than humans do.



**Figure 1.** Heatmap illustrating differences between LLM responses according to the Std Diff metric

**Comparison between LLMs**. The averaged model responses were compared with each other using the same metrics. The strongest correlations were observed between the following model pairs: Qwen2.5.32b and T-pro (Pearson=0.89, Spearman=0.83); RuadaptQwen2.5-32B-Pro-Beta and T-pro (Pearson=0.86, Spearman=0.78); RuadaptQwen2.5-32B-Pro-Beta and Qwen2.5.32b (Pearson=0.85, Spearman=0.78). This indicates substantial similarity in their emotional assessment of Russian nouns, likely due to T-pro being based on the Qwen2.5 model family, and RuadaptQwen2.5-32B-Pro-Beta being the adaptation to Russian of T-pro.

Notably, the lowest correlations were found between: T-lite and Llama 3 (Pearson=0.46, Spearman=0.38); RuadaptQwen2.5-32B-Pro-Beta and T-lite (Pearson=0.47, Spearman=0.40); Qwen2.5-32B and T-lite (Pearson=0.47, Spearman=0.40). This result may indirectly demonstrate the influence of model parameter count on the similarity of its emotional assessments to human ones, as Llama 3 contains the highest number of parameters (70B) among the models in the experiment. RuadaptQwen2.5-32B-Pro-Beta and Qwen2.5-32B are also characterised by a relatively high parameter count. Standard deviation ranges from 0.09 (T-pro vs YandexGPT 5 Lite) to 0.44 (RuadaptQwen2.5-7B vs Llama 3), indicating significant variability in the scale of differences between models. The difference in model responses according to the standard deviation metric is presented in Fig. 1.

# Conclusion

This study of emotional assessment of Russian nouns by language models has revealed substantial differences between machine and human perception of words' emotional content. Although some models demonstrated relatively high correlation with human assessments, persistent systematic deviations were observed across all tested models.

The research findings highlight two significant patterns in LLMs' emotional processing. First, there is a consistent tendency to hyperbolise negative emotions ("disgust", "anger", "sadness"). Second, models display considerable variability in assessing positive emotions ("happiness"), suggesting fundamental disparities in their emotional perception mechanisms. These differences are particularly evident in the assessment of words related to negative events, taboo subjects, and religious vocabulary, where the greatest discrepancies with human assessments were observed.

The obtained results indicate that model size, while being a significant factor, is not the sole determinant for achieving high correlation with human responses. This is illustrated by Llama-3 (70B) which, despite being the largest among the studied models, showed average results. Optimal performance is achieved through a balanced approach that considers both model size and the quality of language adaptation. The highest correlation indicators were achieved by RuadaptQwen2.5-32B-Pro-Beta, which is characterised by both a relatively large number of parameters (32B) and targeted adaptation to Russian. However, the smaller Ruadapt family model, Ruadapt Qwen2.5-7B, despite its language adaptation, showed one of the lowest results.

In the course of further work on this research, we intend to expand the number of LLMs under examination (for instance, by incorporating larger-scale models). Additionally, we plan to conduct a more detailed investigation of word groups for which LLM evaluations diverge most significantly from human assessments, as this area holds considerable research potential.

## Limitations

In the present study, the testing of various hyperparameters and prompt engineering was conducted using RuadaptQwen2.5-32B-Pro-Beta as a base model. This approach may confer a certain advantage upon the aforementioned model in comparison with the others.

## Acknowledgements

## References

1. Acheampong, F.A., Wenyu, C., Nunoo-Mensah, H.: Text-based emotion detection: Advances, challenges, and opportunities. Engineering Reports 2, e12189 (2020). `https://doi.org/10.1002/eng2.12189`

2. Plaza-del Arco, F.M., Cercas Curry, A.A., Cercas Curry, A., Hovy, D.: Emotion analysis in NLP: Trends, gaps and roadmap for future directions. In: Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation. pp. 5696–5710 (2024)

3. Plaza-del Arco, F.M., Curry, A.C., Curry, A., *et al.*: Angry men, sad women: Large language models reflect gendered stereotypes in emotion attribution. In: Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics. vol. 1, pp. 7682–7696. Association for Computational Linguistics, Bangkok (2024). `https://doi.org/10.18653/v1/2024.acl-long.415`

4. Bostan, L.A.M., Klinger, R.: An analysis of annotated corpora for emotion classification in text. Tech. rep., Otto-Friedrich-Universität, Bamberg (2024)

5. Cavicchio, F.: Emotion Detection in Natural Language Processing. Springer, Cham (2025). `https://doi.org/10.1007/978-3-031-72047-5`

6. Chen, Y., Yan, S., Liu, S., *et al.*: EmotionQueen: A benchmark for evaluating empathy of large language models. In: Findings of the Association for Computational Linguistics: ACL 2024. pp. 2149–2176. Association for Computational Linguistics, Bangkok, Thailand (2024). `https://doi.org/10.18653/v1/2024.findings-acl.128`

7. Dalal, D., Negi, G., Picca, D.: LLMs and emotional intelligence: Evaluating emotional understanding through psychometric tools. In: Proceedings of the 33rd ACM Conference on User Modeling, Adaptation and Personalization. pp. 323–328. UMAP '25 (2025). `https://doi.org/10.1145/3699682.3728315`

8. Huang, J., Lam, M.H., Li, E.J., *et al.*: Apathetic or empathetic? Evaluating LLMs' emotional alignments with humans. Advances in Neural Information Processing Systems 37, 97053–97087 (2024)

9. Kazyulina, M., Babii, A., Malafeev, A.: Emotion classification in Russian: Feature engineering and analysis. In: Analysis of Images, Social Networks and Texts, AIST 2020. Lecture Notes in Computer Science, vol. 12602, pp. 135–148 (2021). `https://doi.org/10.1007/978-3-030-72610-2_10`

10. Lyashevskaya, O.N., Sharov, S.A.: Frequency Dictionary of Modern Russian Language (based on the materials of the Russian National Corpus). Azbukovnik, Moscow (2009), (in Russian)

11. Lyusin, D., Sysoeva, T.A.: ENRuN database: Emotional ratings of Russian nouns. Experimental Psychology 18(2), 206–219 (2025). `https://doi.org/10.17759/exppsy.2025180212`, (in Russian)

12. MacCann, C., Roberts, R.D.: New paradigms for assessing emotional intelligence: theory and data. Emotion 8(4), 540–551 (2008). `https://doi.org/10.1037/a0012746`

13. Raj, P.: A literature review on emotional intelligence of large language models (LLMs). International Journal of Advanced Research in Computer Science 15(4) (2024). `https://doi.org/10.26483/ijarcs.v15i4.7111`

14. Sabour, S., Liu, S., Zhang, Z., *et al.*: EmoBench: Evaluating the emotional intelligence of large language models. In: Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics. vol. 1, pp. 5986–6004. Association for Computational Linguistics, Bangkok (2024). `https://doi.org/10.18653/v1/2024.acl-long.326`

15. Salnikov, M., Korzh, D., Lazichny, I., *et al.*: Geopolitical biases in LLMs: what are the "good" and the "bad" countries according to contemporary language models. `https://arxiv.org/abs/2506.06751` (2024)

16. Sysoeva, T.A., Lyusin, D.V.: Development of an extended database with emotional ratings of nouns ENRuN-2: successes, problems and prospects. In: Vladimirov, I., Korovkin, S. (eds.) Psychology of Cognition: Proceedings of the All-Russian Scientific Conference. pp. 316–320. YARSU, Yaroslavl (2024), (in Russian)

17. Tikhomirov, M., Chernyshov, D.: Facilitating large language model Russian adaptation with learned embedding propagation. Journal of Language and Education 10(4), 130–145 (2024). `https://doi.org/10.17323/jle.2024.22224`

18. Wang, X., Li, X., Yin, Z., *et al.*: Emotional intelligence of large language models. Journal of Pacific Rim Psychology 17, 18344909231213958 (2023). `https://doi.org/10.1177/18344909231213`

# LLM for Semantic Role Labeling of Emotion Predicates in Russian

*Ivan V. Smirnov*[1] (iD), *Daniil S. Larionov*[1] (iD), *Elena N. Nikitina*[1] (iD), *Grigory A. Kazachonok*[2] (iD)

Semantic role labeling (SRL) for morphologically rich languages, such as Russian, faces significant challenges due to complex case marking systems, free word order, and limited annotated resources. These challenges are particularly acute for emotion predicates, which require specialized linguistic expertise to capture distinctions between roles denoting those who feel, causes and objects of feelings. We propose a novel approach that leverages large language models to address SRL for Russian emotion predicates through few-shot in-context learning combined with predicate-specific instructions. Our method was evaluated on a manually annotated dataset of 169 sentences containing six emotion predicate groups extracted from Russian social media texts. We compared three state-of-the-art LLMs (Claude 3.7 Sonnet, GPT-5 Mini, and DeepSeek V3) against a RuELECTRA-based trained sequence labelling baseline using both exact and partial matching criteria. Claude 3.7 achieved the highest performance with 74.85% F1 score on partial matching, substantially outperforming the baseline (22.67%). For general predicates on FrameBank, our adapted method with GPT-5 Mini reached 85.0% F1 compared to the previous state-of-the-art of 80.1%. The LLM-based approach successfully handles complex linguistic phenomena, including syntactic zeros and multi-word arguments, while requiring minimal manually annotated training data. We demonstrate that LLM-based methods can significantly advance SRL for Russian by reducing dependency on large-scale annotated corpora while achieving competitive performance.

*Keywords: semantic role labeling, llm, russian language, deep learning, neural networks.*

## Introduction

Semantic Role Labeling (SRL) is a fundamental task in natural language processing that aims to identify the semantic relationships between predicates and their arguments in sentences [11]. Traditional approaches to SRL have largely relied on supervised learning methods trained on carefully annotated corpora such as FrameBank [14]. However, for morphologically rich languages like Russian, the task presents significant challenges due to complex case marking systems, relatively free word order, and limited availability of annotated resources. The annotation scarcity problem is a significant bottleneck, particularly for complex semantic classes such as emotion predicates. Emotion predicates, which include verbs of fear and emotional attitude, and psychological states (such as "пугать-пугаться" (to frighten – to be frightened), "бояться" (to fear), "нравиться" (to like), "любить" (to love), etc.), present unique challenges due to their complex argument structures [17] and the subjective nature of emotional experiences, which specifically appears in syntactic zeros of experiencer (the argument denoting those who undergo emotions): *Пугает неопределенность* (Uncertainty frightens ∅); *вид устрашает...как в Чернобыле..* (The view frightens ∅ like in Chernobyl); *Рост тарифов ЖКХ не страшил бы так, если бы в городе создавались новые высокотехнологичные рабочие места* (The grow of utilities rates would not have frightened ∅ so much if the new high tech jobs had been created)[3] [16].

The annotation of emotion predicates requires specialized linguistic expertise to capture the subtle distinctions between different types of emotional roles, such as experiencers, stimuli, and

---

[1]Federal Research Center "Computer Science and Control" of Russian Academy of Sciences, Moscow, Russian Federation
[2]Moscow Institute of Physics and Technology, Dolgoprudny, Russia
[3]Here and everywhere else examples are provided with original authors grammar.

targets of emotions [19], as well as their superficial expressions. See, for example, split of Causator role into two arguments: **Мэр** *удивил горожан* **гранитными бордюрами**, clause and infinitive arguments: *Владимир, а вам нравится,* **когда вас с кем сравнивают?**; *боится* **лишние секунды потерять**. Traditional approaches rely heavily on manually curated training data, which is both expensive to produce and is limited in coverage. This creates a particular challenge for languages like Russian, where comprehensive emotional semantic resources are scarce compared to English.

Recent advances in Large Language Models (LLMs) have demonstrated remarkable capabilities across various natural language processing tasks through in-context learning and prompting-based approaches [12]. These models show particular promise for tasks that benefit from semantic understanding and structured reasoning, making them natural candidates for semantic role labeling applications.

We propose a novel approach that leverages the linguistic capabilities of LLMs to address the challenges of semantic role labeling in Russian, particularly for the domain of emotion predicates. Our method combines two key techniques: few-shot learning and instruction retrieval. Few-shot learning enables the model to generalize from a small number of examples, while instruction retrieval allows the system to access relevant linguistic knowledge and annotation guidelines dynamically.

Our contributions advance the state of semantic role labeling for Russian by demonstrating that LLM-based approaches can achieve competitive performance while reducing the dependency on large-scale manually annotated corpora. The proposed framework provides flexibility for different application scenarios: the expert-curated instructions and examples for high-precision and low-data cases in specialized domains, and the generic retrieval approach for broader coverage where more annotated data is available. We release the source code of our approach at `https://github.com/ru-nlp/llm-for-srl`.

The article is organized as follows. Section 1 reviews related work on LLM-based approaches to semantic role labeling, the FrameBank resource for Russian, and previous research on emotion predicates. Section 2 describes our methodology, including the construction of the Russian emotion predicates dataset, the configuration of three evaluated LLMs (Claude 3.7 Sonnet, GPT-5 Mini, and DeepSeek V3) and the RuELECTRA-SRL baseline, our few-shot prompting strategy with predicate-specific instructions, and the evaluation protocol with exact and partial matching criteria. We also describe the adaptation of our approach to general predicates on FrameBank. Section 3 presents experimental results for both emotion predicates and general semantic role labeling, with detailed analysis of model performance across different semantic roles and matching criteria. Section 4 discusses the computational and linguistic implications of our findings, examining how LLMs handle complex phenomena such as syntactic zeros, anaphoric references, multi-word arguments, and irregular constructions in social media texts. The Conclusion summarizes our results and outlines directions for future research.

## 1. Related Work

### 1.1. LLM for SRL

[8] contains a comprehensive survey on various methods and applications of semantic role labeling. Large language models have become the dominant paradigm for solving NLP tasks, and naturally, almost all of the recent approaches to SRL employ LLMs in some way or form.

Current approaches integrate LLMs into the SRL pipeline in several innovative ways. A common technique involves using the embeddings generated by an LLM as rich feature inputs for a downstream classifier, often enhanced with syntactic information such as dependency parses. This hybrid approach, exemplified by [11], combines deep semantic understanding with structural grammatical cues.

The current state-of-the-art results for both English and Chinese, as demonstrated by [12], are achieved through a more integrated method. Their model employs prompts to a fine-tuned LLM, equipped with a self-correction mechanism and a searchable database to improve accuracy and consistency.

The prompting paradigm itself is a substantial area of study. The research by [9] explores a few-shot prompt-based approach, analyzing the inherent capability of LLMs to understand semantic structure without extensive task-specific training. Similarly, [22] investigates a zero-shot technique for SRL and sentiment analysis, further probing the model's cross-lingual abilities by testing if semantic roles are preserved when translating sentences from English to Arabic.

In [9], a few-shot prompt-based approach is introduced, with a discussion of LLMs' capabilities for understanding semantics. [22] features a zero-shot prompting technique for semantic role labeling and sentiment analysis in English. The authors also investigate whether LLMs can correctly translate English sentences into Arabic, preserving the semantic role labels.

The application of these techniques also extends to specialized domains. [4] evaluate prompt-based SRL within the legal domain, comparing a general-purpose LLM to one fine-tuned specifically on legal corpora. Their findings indicate that domain-specific adaptation yields significant gains in performance and efficiency for processing complex legal texts. Finally, some research moves beyond pure prompting architectures. For instance, [27] introduces a novel framework that combines prompts to LLMs with Graph Neural Networks, aiming to capture both semantic and relational dependencies between arguments.

## 1.2. FrameBank

FrameBank [14] is a semantically annotated database of Russian sentences primarily based on the Berkeley FrameNet project [3]. It serves as a valuable resource for automating SRL and has been widely used in Russian NLP research.

FrameBank has been instrumental in training various SRL systems for Russian. [10] utilized FrameBank to train a semantic role classifier based on a Support Vector Machine (SVM) model, demonstrating the FrameBank's utility for traditional machine learning approaches. In [23], the authors experimented with training a neural network on this dataset. In a more recent work, [11] employed FrameBank to train a neural network encoder specifically designed to identify arguments and assign them their correct semantic roles.

## 1.3. Emotion Predicates

In [18], different approaches to emotion identification are compared. In short, there are three main approaches: the first categorizes emotions across entire texts, the second labels separate emotionally charged words, and the third classifies emotions within clauses. [25], similar to the third approach and our own, studies emotions within semantic frames: a frame represents an event that triggers an emotional response.

Authors in [7] present SRL4E, a unified evaluation framework that consolidates six heterogeneous emotion datasets under a standard annotation scheme based on Plutchik's emotions, enabling consistent training and evaluation of systems that identify not only emotions, but also their semantic constituents (experiencer, target, and stimulus) within text.

## 2. Methodology

### 2.1. Russian Language Dataset of Emotion Predicates

We constructed a specialized evaluation dataset[4] for Russian semantic role labeling, focusing on psychological predicates. The dataset comprises 169 manually annotated sentences extracted from Russian social media and informal text sources. Our annotation targets six predicate groups representing emotional and psychological states. They are verbs of fear: *пугать* (frighten), *ужасать* (horrify), *бояться* (fear), *опасаться* (be apprehensive), *страшить* (intimidate), and verbs of emotional attitude: *нравиться/любить* (like/love).

Each sentence was annotated by an expert linguist following a predefined semantic role taxonomy comprising five primary roles:

- **Experiencer**: The entity experiencing the psychological state (**Его** *страшит неопределенность;* **Он** *страшится будущего;* **Он** *любит девушку;* **Ему** *нравится девушка*);
- **Causator**: The entity or event that triggers the psychological response (*Его страшит* **неопределенность***; Он страшится* **будущего**);
- **Instrument**: The means or medium through which the Causator induces the response (*Будущее пугает* **неопределенностью**);
- **Deliberative**: The entity about whose welfare the Experiencer is concerned (typically marked by the Russian preposition *за*) (*Он боится* **за сына**);
- **Object**: The entity or event towards which the Experiencer feels the attitude (*Он любит* **девушку***; Ему нравится* **девушка**).

### 2.2. Model Configuration

We have evaluated several available LLMs, including both closed and open-weight ones:

- Anthropic Claude Sonnet 3.7 [1] – a proprietary State-of-the-Art (SOTA) LLM, developed by Anthropic. Excels at instruction following and is a particularly powerful tool for non-English language processing. The model is a hybrid reasoner; however, we have specifically disabled reasoning in our experiments;
- OpenAI GPT-5 Mini[5] – a mini variant of SOTA LLM from OpenAI. The model is reasoning-based, which means that the reasoning component cannot be disabled in it. Thus, we set it to minimal reasoning effort;
- DeepSeek V3 [13] – an open-weight non-reasoning LLM. The model contains 671 billion parameters and requires substantial hardware resources for running: up to 16 H100/A100 GPUs. Despite the size, this model is particularly interesting due to its open availability, which allows practitioners to utilize their existing HPC resources.

---

[4]https://huggingface.co/datasets/dl-ru/srl-emotion-predicates
[5]https://openai.com/index/gpt-5-system-card/

We established a baseline using RuELECTRA-SRL [2], a transformer-based model specifically fine-tuned for Russian semantic role labeling through token classification using a dataset from the previous work[6]. This approach mimics the functionality of NER models with BIO-style annotation to capture multi-word arguments.

## 2.3. Prompting Strategy

Our LLM-based approach is focused on a few-shot in-context learning approach with predicate-specific demonstrations. The prompting template for LLM consists of four main components:

1. **System Prompt**: Role specification as a native Russian linguist with explicit instructions for null-role handling ("No-Roles#No-Roles");
2. **Rule Specification**: JSON-formatted semantic role definitions tailored to each predicate group;
3. **Few-shot Examples**: All available training instances from the target predicate group, formatted as input-output pairs;
4. **Target Query**: The sentence requiring semantic role analysis.

The output format specification requires the model to generate role annotations as "- argument#role" pairs, facilitating parsing and evaluation. See example of the prompt in Fig. 1.

## 2.4. Evaluation Protocol

We evaluated model performance using both exact and partial matching criteria. Partial matching is necessary to capture additional aspects of model performance in multi-word arguments, particularly in clauses or phrases. It is not always possible to perfectly align LLMs with expert annotators on what should be deemed an argument in a multi-word case. The evaluation protocol is implemented as follows:

**Exact Matching**: A predicted argument-role pair $(a_p, r_p)$ is considered correct if and only if there exists a ground truth pair $(a_t, r_t)$ such that $\text{normalize}(a_p) = \text{normalize}(a_t)$ and $r_p = r_t$, where $\text{normalize}(\cdot)$ applies lowercase conversion and whitespace normalization.

**Partial Matching**: A prediction is considered correct if $r_p = r_t$ and either $\text{normalize}(a_p) \subseteq \text{normalize}(a_t)$ or $\text{normalize}(a_t) \subseteq \text{normalize}(a_p)$ (substring containment).

For each matching criterion, we computed:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \tag{1}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \tag{2}$$

$$\text{F}_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}. \tag{3}$$

We report both macro-averaged metrics across all instances and per-role micro-averaged metrics to assess performance across semantic roles.

---

[6]`https://huggingface.co/datasets/dl-ru/srl-data-ru`

```
[System]
You are native russian linguist specializing in semantic
role labelling. If there are no roles:  reply with
- No-Roles#No-Roles.

You must follow these rules in your work:
{
"semantic_roles":  {
"experiencer":  "Entity experiencing fear or dread.",
"causator":  "Entity or event that instills fear or dread.",
"instrument":  "Means or medium used by the causator
to induce fear (when present)."
}
}
```

[User] Полиция пытается устрашить преступников новыми мерами

[Assistant]
- Полиция#Causator
- преступников#Experiencer
- мерами#Instrument

[User] Не нужно устрашать малышей рассказами о монстрах

[Assistant]
- малышей#Experiencer
- рассказами#Instrument

*... [16 more examples] ...*

[User] Сейчас «выздоравливать» начнут сотнями статистику
устрашающую надо улучшать срочно.

**Figure 1.** Few-shot prompt for semantic role labeling of Russian Emotion predicates. The system message defines three semantic roles (Experiencer, Causator, Instrument), followed by 18 labeled examples with user inputs and assistant responses demonstrating various linguistic patterns and colloquial usage

## 2.5.  General Role Labeling

We have additionally tried the few-shot approach described above for labeling any predicates, not only the emotion ones. To do this, we have used the semantically annotated FrameBank dataset [14]. First, the dataset is filtered so that every sentence contains exactly one annotated predicate and every predicate appears at least 10 times in the dataset. Then, a portion of the sentences is removed to be later used as examples. The rule is that for every predicate, there must be at least 5 examples, and for every semantic role that its arguments can take on, there must be at least one example.

We have employed Gemini 2.5 Flash, GPT-5 Mini, and DeepSeek-V3 for this problem. The prompt template contained additional instructions concerning the model's behavior: in what

form it should answer, how it should identify the arguments if multiple words fit, etc. It also contained the target sentence, the predicate, and examples of semantic roles for that predicate.

Since RuELECTRA-SRL only deals with emotion predicates, a different baseline had to be chosen. We chose the approach from [11]. This approach is based on a pre-trained language model, fine-tuned on FrameBank, and it has the high SRL score on FrameBank corpora. For evaluation, we gave the same 10000 randomly selected sentences to every model.

# 3. Results

## 3.1. Semantic Role Labeling for Emotion Predicates

**Table 1.** Semantic Role Labeling Evaluation Results. For per-role results we present F1 score. For overall results we present macro-averaged F1 score, precision and recall. **Bold numbers** indicates best scores in each category across models with exact matching. <u>Underlined numbers</u> indicate best score with partial matching

| Role/Metric | Claude 3.7 | | GPT-5 Mini | | DeepSeek-V3 | | Baseline | |
|---|---|---|---|---|---|---|---|---|
| | **Exact** | **Exact+Partial** | **Exact** | **Exact+Partial** | **Exact** | **Exact+Partial** | **Exact** | **Exact+Partial** |
| Causator | **0.4625** | <u>0.7000</u> | 0.4100 | 0.6400 | 0.4146 | 0.6463 | 0.0235 | 0.0235 |
| Cause | 0.0000 | <u>0.6667</u> | **0.1818** | 0.1818 | 0.0000 | 0.4000 | 0.0000 | 0.4000 |
| Deliberative | 0.0000 | 0.8571 | **0.8000** | 0.8000 | 0.3333 | 0.5000 | 0.0000 | <u>1.0000</u> |
| Experiencer | **0.6636** | <u>0.7465</u> | 0.6204 | 0.7130 | 0.5381 | 0.5685 | 0.3543 | 0.3657 |
| Instrument | 0.4444 | 0.4444 | 0.4211 | 0.4211 | **0.7500** | <u>0.7500</u> | 0.2500 | 0.2500 |
| Object | **0.6889** | 0.8667 | **0.6882** | <u>0.8817</u> | 0.4595 | 0.7027 | 0.1026 | 0.1197 |
| Overall F1 | **0.5808** | <u>0.7485</u> | 0.5404 | 0.6949 | 0.4739 | 0.6174 | 0.1965 | 0.2267 |
| Overall Precision | **0.6068** | <u>0.7821</u> | 0.5087 | 0.6540 | 0.5317 | 0.6927 | 0.2746 | 0.3169 |
| Overall Recall | 0.5569 | 0.7176 | **0.5765** | <u>0.7412</u> | 0.4275 | 0.5569 | 0.1529 | 0.1765 |

Table 1 presents the evaluation results for semantic role labeling of Russian emotion predicates across four models. Claude 3.7 Sonnet achieved the highest overall performance with a 0.7485 F1 score on partial matching, followed by GPT-5 Mini (0.6949) and DeepSeek-V3 (0.6174), while the RuELECTRA-SRL baseline showed substantially lower performance (0.2267). The performance gap between exact and partial matching metrics shows that identifying precise argument boundaries is challenging for LLMs. Notably, all LLM-based approaches struggled with the rare Deliberative role (marked by the preposition *за*), though GPT-5 Mini achieved 0.8000 F1 on exact matching for this category. Per-role analysis shows considerable variance across semantic categories, with more frequent Experiencer and Object roles generally yielding higher scores across all models, while Instrument and Cause roles presented consistently lower performance. The results demonstrate that while LLMs substantially outperform traditional token classification approaches, significant room for improvement remains, particularly in handling less frequent semantic roles and accurately identifying argument boundaries.

## 3.2. General Semantic Role Labeling

Three LLMs were compared to the baseline model scores. Overall, they performed significantly better than the baseline model. Out of the three models we used, the one that showed the best results was GPT-5 Mini.

However, as one can see from Tab. 2, our few-shot approach performs substantially worse on some less common roles. Out of the 10000 test sentences, Gemini got only 61.1% completely

**Table 2.** Performance of different models on semantic role labeling for general predicates. For specific roles, the F1 scores are provided. The best results in each category are in **bold**

| Role / Metric | DeepSeek-V3 | GPT-5 Mini | Gemini 2.5 Flash | Baseline |
|---|---|---|---|---|
| agent (11.7%) | 82.2 | **87.3** | 82.9 | 79.5 |
| patient (10.2%) | 86.6 | **88.4** | 85.7 | 86.9 |
| theme (6.9%) | 83.6 | 86.5 | **86.9** | 77.6 |
| sbj of psychol. state (6.2%) | 89.8 | **92.4** | 90.4 | 85.2 |
| goer (5.7%) | 87.9 | **91.4** | 89.0 | 85.9 |
| cause (4.7%) | 86.1 | 85.1 | **88.3** | 87.4 |
| speaker (4.5%) | 86.2 | **89.7** | 86.7 | 75.8 |
| location (4.1%) | 82.8 | 82.4 | 82.1 | **84.9** |
| content of action (3.6%) | 83.0 | 85.3 | 82.4 | **86.3** |
| content of thought (3.4%) | 86.8 | **88.1** | 84.8 | 77.0 |
| content of speech (3.4%) | 80.8 | **82.1** | 79.3 | 72.6 |
| final destination (3.4%) | 87.0 | **87.1** | 86.1 | 59.8 |
| result (2.8%) | 87.0 | 85.6 | **91.1** | 58.4 |
| patient of motion (2.6%) | 83.5 | **86.6** | 83.1 | 84.4 |
| stimulus (2.4%) | **87.2** | 86.1 | 85.3 | 78.1 |
| cognizer (2.3%) | 81.9 | **85.7** | 83.2 | 80.8 |
| addressee (1.8%) | 86.0 | **86.5** | 85.9 | 77.4 |
| perceiver (1.7%) | 88.4 | **93.5** | 91.7 | 84.3 |
| counteragent (1.6%) | 87.0 | **87.8** | 86.7 | 60.9 |
| effector (1.4%) | 66.1 | 65.8 | 67.6 | **78.9** |
| subject of social attitude (1.1%) | 77.9 | 81.0 | **82.4** | 80.8 |
| initial point (1.1%) | 88.3 | 84.5 | **88.4** | 78.1 |
| topic of speech (1.0%) | **82.5** | 81.7 | 78.3 | 68.0 |
| manner (1.0%) | 64.0 | 54.5 | 65.7 | **76.0** |
| recipient (1.0%) | 81.0 | **86.2** | 79.8 | 74.5 |
| goal (0.9%) | **76.0** | 75.5 | 75.8 | 73.3 |
| field (0.7%) | 78.8 | 80.0 | 74.5 | **91.3** |
| attribute (0.7%) | 73.1 | 70.3 | 67.0 | **82.5** |
| source of sound (0.7%) | 86.4 | **92.2** | 87.2 | 71.6 |
| behaver (0.6%) | 78.2 | **86.5** | 77.4 | 84.6 |
| situation in focus (0.6%) | 75.1 | 73.4 | 65.6 | **88.2** |
| counteragent of social attitude (0.6%) | **84.6** | 80.2 | 83.9 | 65.5 |
| sbj of physiol. reaction (0.6%) | 89.0 | **91.4** | 88.3 | 80.4 |
| topic of thought (0.6%) | 70.8 | 67.1 | 70.5 | **92.2** |
| potential patient (0.5%) | 88.3 | 89.8 | 89.3 | **90.1** |
| status (0.5%) | **87.8** | 86.1 | 78.1 | 83.3 |
| patient of social attitude (0.5%) | 58.4 | 58.2 | 60.7 | **80.8** |
| standard (0.5%) | 86.4 | **88.0** | 85.6 | 82.7 |
| term (0.5%) | 91.2 | 86.6 | 90.4 | **86.6** |
| attribute of action (0.5%) | 79.5 | 74.5 | 74.8 | **80.4** |
| causer (0.4%) | 54.3 | **69.3** | 68.6 | 68.7 |
| initial possessor (0.4%) | 76.1 | **81.0** | 63.3 | 78.3 |
| potential threat (0.4%) | 84.1 | **87.0** | 78.6 | 77.9 |
| path (2.3%) | 65.7 | 67.5 | 61.8 | **84.9** |
| Argument Extraction F1 | 87.6 | **88.4** | 86.4 | 79.4 |
| Argument Extraction Precision | **87.7** | 84.2 | 84.9 | 74.5 |
| Argument Extraction Recall | 87.5 | **93.2** | 86.4 | 85.1 |
| Role Labeling F1 | 83.3 | **85.0** | 83.1 | 80.1 |

right, and both DeepSeek and GPT-5 got 60.5%. Curiously, there is a substantial overlap of examples that the models got wrong. All three LLMs gave wrong answers on the same 20% of the data, indicating that some sentences may be inherently "counterintuitive" for LLMs (see also [9], where mistakes made by LLMs were compared to mistakes by human non-experts).

### 3.3. Analysis of Exact vs. Partial Matching Disparities

The performance patterns in Tab. 1 reveal a notable discrepancy across models in their ability to identify precise argument boundaries. Claude 3.7 Sonnet, despite achieving the highest overall F1 score (0.7485) on partial matching, recorded zero scores for exact matching in two semantic categories (Cause and Deliberative), while maintaining substantial partial matching scores for these same roles (0.6667 and 0.8571 respectively). In contrast, GPT-5 Mini demonstrated non-zero exact matching performance across all categories, including a particularly strong 0.8000 F1 score for the Deliberative role. DeepSeek V3 exhibited zero exact matching only for the Cause role.

This pattern correlates with our multiword argument analysis, where Claude achieved the highest success rate (70%) in correctly extracting and labeling multiword arguments, compared to GPT-5 Mini (65%) and DeepSeek V3 (63.33%). The data suggest that Claude adopts a more expansive approach to argument boundary identification, consistently capturing the semantic core of arguments while frequently deviating from expert-annotated boundaries. This is partic-

ularly evident in rare semantic roles: the Deliberative role appears only in 7 instances (4.1% of arguments), and Cause in 5 instances (3.0%). For such infrequent categories, Claude tendency to extract semantically appropriate but boundary-imprecise arguments results in complete exact matching failure, yet high partial matching success.

GPT-5 Mini consistent non-zero exact matching across all roles reflects a more conservative extraction strategy that better aligns with annotator boundary conventions, albeit at the cost of missing some semantically relevant material. The model reasoning component (which does not operate in our experiments) does not contribute to boundary-aware predictions. DeepSeek V3 intermediate behavior, with exact matching failure only for Cause, suggests it falls between these two strategies, as does reasoning, but lacks in depth understanding.

## 4. Discussion

The LLM-based approach allows practitioners to trade increased computational resources for reduced data annotation costs while also substantially improving the method robustness. In some instances, such as with SRL for the Russian language, this tradeoff is particularly effective, which stems from the natural complexity of semantic linguistic annotation. Indeed, the LLM-based approach, in our case, requires only minimal human annotation to briefly cover the predicate groups we are focusing on. Training a regular deep learning model, such as those we demonstrate as baselines in the section above, would require 30-40x more annotated examples to achieve the desired level of quality, not accounting for complex cases, such as multi-word arguments. This computational intensity makes this approach well-suited for HPC environments, where accelerated processing capabilities can efficiently handle the increased resource demands of LLM inference at scale.

Linguistically, the well-known difficulties of text annotation and semantic roles prediction in Russian, as a morphologically rich language, mentioned in the Introduction, should be extended by a number of additional problems. Generally speaking, they concern the problems of grammar of constructions, on one hand, and text production and genre, on the other hand.

Firstly, in the Russian grammar, there exist constructions with the personal and impersonal subjects of the sentence (Doer, Experiencer, etc.) omitted, which are generally referred to as "syntactic zeros" [15] of complete sentences and "zero pronouns" [6]. The syntactic zeros (∅) can take such meanings as "others (somebody excluding the speaker)" – *Уводили тебя на рассвете* (Akhmatova, ∅ took you away at the sunrise), "everybody including the speaker" – *Как потопаешь, так и полопаешь* (Proverb, As ∅ sow so shall ∅ mow) [5]. There are also 1st and 2nd person forms of verbs that definitely indicate the subject of a sentence: the Speaker (I – *Люблю грозу в начале мая* (Tyutchev, ∅ love the thunderstorm in the beginning of May)) and Listener (you – *Послушайте! Ведь, если звезды зажигают...* (Mayakovsky, ∅ Listen! In fact, if somebody lights the stars...)).

Secondly, text production supposes that the subject of the sentence can be replaced by a pronoun or omitted for the reasons of economy and coherence. Such modified subjects refer to the anaphora and can be represented by 3d person pronouns or syntactic zero. For both, the referent can be typically found in the left context of the text. See also summarizing presentation in [20].

Thirdly, we also encountered some unexpected difficulties with genre as we analyzed social media discussion content. This type of text belongs to written speech genres and consists of dialogue and, to some extent, monologue fragments. They are characterised by spoken and spon-

taneous features, and contain interrupted elements within them, as well as irregular usages. The latter, for example, is expressed in irregular cases of Causator: *ужаснулся* **над дырявостью наших законов**. Some syntactic fragments are brief and incomplete, and they frequently appear in dialogues. Some fragments are extensive, for instance in monologues. Structurally, this can result in a distant position between the emotion predicate and its Experiencer argument: **Водитель** *не видел что в другом ряду Жигули остановились, явно пропуская пешехода, куда все несутся, время экономит,* **боится** *лишние секунды потерять, хорошо не задел её.*

The monologue parts structurally and semantically are similar to egocentric 'I'('Я')-texts (such as narrative memoirs, diaries) and therefore contain 'I' (1st person pronoun) omitted constructions. Likely for all actual (present tense) emotion expressions, syntactic zeros of 'I' Experiencer are very specific for emotion verb constructions (compare: *Жаль; Обидно; Грустно; Пугает, что. . .*). Likely in 'I'-texts, constructions of emotion verbs denoting emotions of 'I' Experiencer often include corresponding syntactic zero of 1st person pronoun even in the past tense: *Проезжал мимо и просто ужаснулся,что же за гений архитектор это нарисовал!!!*(∅ was passing by and purely got frightened what a genius architect painted all this <...>). An additional argument for 'I' Experiencer reconstruction in this sentence is that the position of the Causator is filled with a direct speech clause (the problem of direct/indirect speech in complement clauses is discussed in [21, 26]).

Evidently, we could have never expected to extract arguments represented by such syntactic phenomena as noun phrases, clauses, and syntactic zeros, even with the help of LLM. Meanwhile, the expert analysis of what LLM recognizes as arguments of verbs of emotions convinced us that its process of thinking can be very productive and cover complicated and unsolvable cases for other methods, and carefully identify their semantic roles.

See examples of LLM prediction of Causator and Object arguments expressed by a single noun, noun phrase, and clause given below (Tab. 3). Our method is able to analyse and identify constructions with more than one pretender to be the argument of a verb, i.e. independent nouns and noun phrases connected by coordinating conjunction: *Толкучки и общественные места любите?* At the same time it in an arbitrary way can incorrectly shorten a long noun phrase: *Мне не понравилось [отсутствие реакции]#Object руководства больницы на мое предложение поставить там кофейный автомат*; *Особенно понравилось [исполнение]#Object второй песни* – instead of *Мне не понравилось [отсутствие реакции руководства больницы на мое предложение поставить там кофейный автомат]#Object*; *Особенно понравилось [исполнение второй песни]#Object*.

Additionally, our method recognizes not only Causators based on dependent clauses (followed by conjunction), but also independent clauses (no conjunction) ones: *Боюсь, что [в 24 году не до паркунов будет] - Если можно было бы взять его 4-ым, а бы взяла, но боюсь [тогда меня из дома выгонят вместе с ним].* It generally takes place in contexts of internal state verbs in the Present tense with 'I' (speaker) subjects. The meaning of 'I'-subject here becomes closer to Thinker than Experiencer, which reflects in the replacement of indirect (followed by a conjunction) by direct (no conjunction) speech, as we noticed above.

See Tab. 4 for examples of LLM predictions for syntactic zeros, both anaphoric (the meaning derived from the text) and paradigmatic ones (the meaning derived from the form of the verb or construction). It is interesting that in contexts of V3pl construction with a Nominative noun omitted and a present locative component, the latter is predicted as the subject of the sentence, i.e., the Experiencer, which entirely aligns with the common idea of grammatical interpretation.

**Table 3.** Examples of LLM Method recognizing Causator and Object arguments
expressed by a single noun, noun phrase, clause, etc.

| Syntactic status of argument | Text | Fragment in focus translated |
|---|---|---|
| Single noun or pronoun | **Женщин** люблю; **никого** не боялся | (I) like **women**; (I was) afraid of **nobody** |
| Noun phrase | Ольга, я с вами согласна, я тоже люблю **больших собак !** | I like **big dogs** |
| Independent elements | **Толкучки** и **общественные места** любите?  ; Люблю **больших** И **коренастых!** | (Do you) like **crowds** and **public places**; I like **big** and **stocky** (ones) |
| Infinive or infinitive phrase | Татьяна, я живу на 5 этаже и у меня за окном такие сосульки висят, но самостоятельно, я лично, **сбивать** боюсь.; Наталья, а Вы не думаете, что персонал боится **заразиться через детей и их родителей**??? | I am afraid **to break down icicles**; <...> the staff are afraid **to get infected through children and their parents** |
| Dependent clause | Не боитесь что **поклоники Высоцкого вас побьют**? | Aren't you afraid **to be beaten up by (his) fans**? |
| Pronoun followed by dependent clause | Дмитрий, а тех **кто летает по дорогам, как в** <...> **жаленый, не глядя на пешеходные переходы,** ты любишь? | <...> and do you like **those who rush along the roads** <...> |

The gerund constructions are also known as those that carry a syntactic zero of the subject of the action (coreferent to the subject of the main predicate) [24], and LLM can establish such a zero subject in our texts.

We have to consider particular cases that LLM performs incorrectly. It concerns 1st and 2nd person pronouns and syntactic zeros, which are rather deictic than anaphoric, i.e., they do not need to be replaced by nouns from the text, as their referents are the Speaker and the Listener, marked by corresponding personal pronouns. See LLM failures in Tab. 5.

On the other hand, LLM does not operate with classic anaphoric 3rd person pronouns, they are left without their context referents: *Катерина, ну да, мы* **его** *просто очень любим* (we love **him** very much); *Но мне* **это** *не нравится* (But I don't like **it**); *Натали,* **они** *как огня боятся областную жил инспекцию* (**they** are afraid of the local housing inspectorate). Probably, it happens for the reason that their referents are cut away in the process of text syntactic segmentation and therefore stay behind the frame of the sentence analysed. In general, as far as grammatically complicated cases are concerned, the advantages of the great performance of our method are complemented by the disadvantages of irregularity and instability of performance. Someone can still not entirely rely on the decisions of LLM. At the same time, in comparison to all other existing methods for automatic SRL for the Russian Language at scale, LLMs demonstrate the highest accuracy and robustness to linguistic phenomena.

**Table 4.** LLM prediction of Experiencer for syntactic zeros

| ∅ | Predicted from | Meaning | Text | Translation |
|---|---|---|---|---|
| Syntactic zero paradigmatic | Form of the verb, $V_{2s}$ | Personal pronoun 2 listener Вы Вы Ты | Ksu, ну хорошо что **любите** **Не страшитесь**, мне тоже нечупно и душевно и сопли от дыма уже пошли. Татьяна, **ужасайся** дальше | It is good ∅ **love** ∅ **Don't frighten** <...> ∅ **be horrified** again (= don't stop being horrified) |
| Syntactic zero paradigmatic | Form of the verb, $V_{1s}$ | Personal pronoun 1 speaker Я | Александр, только что сказали, что в Москве закрыты аэропорты больничный без посещения поли-ки, а мне позвонили и сказали, что надо идти в регистратуру и забирать его, вот тоже этого **страшусь** уже. | ∅ also **am afraid** already |
| | | Я | Валера, уже напротяжении 16 лет **люблю** творчество граффити, это позерство | ∅ have been **loving** graffiti art for 16 years <...> |
| Syntactic zero anaphoric | $V_{ps}$ chain of predicates | Personal pronoun 1 speaker Я | Приезжала в город в январские каникулы, **ужаснулась** | I went to the town for January vacation, ∅ **got frightened** |
| Syntactic zero anaphoric | Pronoun, $V_{ppl}$ chain of predicates | Все | А так, все увидели и **ужаснулись**. | Everybody saw and ∅ **got frightened** |
| Syntactic zero paradigmatic | $V_{3pl}$ construction | У нас | Или опять пешеход виноват, как у нас **любят** говорить? | <...> as ∅ **like** saying at our's (=People say that..., It is generally said that) |
| Syntactic zero anaphoric | Gerund | - | И хватит уже оглядываться на всяких либералов, **опасаясь** задеть их | ∅ Stop taking into account liberals ∅ **being afraid** to hurt them |
| Syntactic zero anaphoric | $V_{ps}$ chain of predicates, *для него* pronoun | Он | Решил развернуться, затем для него внезапно появилась машина (которая его почти догнала), в итоге **испугался**, потерял контроль над управлением. | ∅ **Decided** to turn around, then suddenly for him a car appeared, finally ∅ **got afraid** |
| Syntactic zero anaphoric | $V_{3s}$ chain of predicates | Марс | Марс ) самый милый пухляш на земле, супер ласковый и **любит** сидеть на ручках | Mars the nicest little fat kid on the Earth, very sweet and ∅ **loves** sitting in the arms |

**Table 5.** LLM predictions for 1st and 2nd person pronouns

| Predicted pronoun | Predicted from | Functioning | Text | Translation |
|---|---|---|---|---|
| 2$^{\text{nd}}$ person pronoun *вам* | Владимир | Wrong | Владимир, а вам **нравится**, когда вас с кем сравнивают? | Vladimir, and do you **like** $<\ldots>$ |
| 1$^{\text{st}}$ person pronoun *мне* | Александр | Wrong | Александр, соты мне **понравились** | Alexander, I **like** honeycombs |
| 1$^{\text{st}}$ person pronoun *мы* | мы | Right | Сынок, мы тебя очень **любим**!!! | Sonny, we **love** you very much |
| Syntactic zero (*Я*) | Валера | Wrong | Валера, уже напротяжении 16 лет **люблю** творчество граффити, а это позерство | Valera, ∅ have been **loving** graffiti art for 16 years |

## Conclusion

We presented a novel approach to semantic role labeling for Russian emotion predicates using large language models with few-shot learning. Our method demonstrates that LLMs can achieve significantly better performance than traditional supervised approaches, with Claude 3.7 achieving 74.85% F1 score on partial matching compared to 22.67% for the RuELECTRA baseline. For general predicates on FrameBank, GPT-5 Mini reached 85.0% F1, substantially outperforming the previous state-of-the-art of 80.1%.

The method successfully handles complex linguistic phenomena, both specific to Russian and natural to various languages in general, including syntactic zeros, anaphoric references (less than others), multi-word arguments, and clausal structures. LLMs demonstrate remarkable capability in identifying emotion arguments even in challenging social media texts with interrupted elements and irregular constructions. However, the approach shows limitations with certain anaphoric 3rd person pronouns and occasionally produces arbitrary segmentation of long noun phrases.

Future work should address the stability of predictions across different predicate types and explore hybrid approaches combining LLM reasoning with structured linguistic knowledge. The development of larger annotated corpora for Russian emotion predicates remains critical for advancing the field. Additionally, the study of cross-lingual transfer learning could leverage resources from morphologically similar languages to improve coverage of rare semantic roles.

## Acknowledgements

# References

1. Claude 3.7 Sonnet System Card. `https://api.semanticscholar.org/CorpusID:276612236`

2. Arkhipov, M., Trofimova, M., Kuratov, Y., Sorokin, A.: Tuning Multilingual Transformers for Language-Specific Named Entity Recognition. In: Erjavec, T., Marcińczuk, M., Nakov, P., *et al.* (eds.) Proceedings of the 7th Workshop on Balto-Slavic Natural Language Processing. pp. 89–93. Association for Computational Linguistics, Florence, Italy (Aug 2019). `https://doi.org/10.18653/v1/W19-3712`

3. Baker, C.F., Fillmore, C.J., Lowe, J.B.: The Berkeley FrameNet project. In: 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 1. pp. 86–90. Association for Computational Linguistics, Montreal, Quebec, Canada (Aug 1998). `https://doi.org/10.3115/980845.980860`

4. Bakker, R., Schoevers, A., van Drie, R., *et al.*: Semantic role extraction in law texts: a comparative analysis of language models for legal information extraction. Artificial Intelligence and Law. P. 1–35 (03 2025). `https://doi.org/10.1007/s10506-025-09437-x`

5. Bulygina, T.V.: Ya (I), ty (you) and others in the Russian grammar. Res philologica. Philological researches pp. 111–126 (1990)

6. Bulygina, T.V., Shmelev, A.D.: Language conceptualization of the world (based on the material of Russian grammar). Shkola "Jazyki russkoj kul'tury", Moscow (1997)

7. Campagnano, C., Conia, S., Navigli, R.: SRL4E – Semantic Role Labeling for Emotions: A unified evaluation framework. In: Muresan, S., Nakov, P., Villavicencio, A. (eds.) Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 4586–4601. Association for Computational Linguistics, Dublin, Ireland (May 2022). `https://doi.org/10.18653/v1/2022.acl-long.314`

8. Chen, H., Zhang, M., Li, J., *et al.*: Semantic role labeling: A systematical survey (2025). `https://doi.org/10.48550/arXiv.2502.08660`

9. Cheng, N., Yan, Z., Wang, Z., *et al.*: Potential and Limitations of LLMs in Capturing Structured Semantics: A Case Study on SRL. In: Advanced Intelligent Computing Technology and Applications: 20th International Conference, ICIC 2024, Tianjin, China, August 5–8, 2024, Proceedings, Part I. pp. 50–61. Springer-Verlag (2024). `https://doi.org/10.1007/978-981-97-5663-6_5`

10. Kuznetsov, I.: Semantic Role Labeling for Russian Language Based on Russian FrameBank. In: Khachay, M.Y., Konstantinova, N., Panchenko, A., *et al.* (eds.) Analysis of Images, Social Networks and Texts. pp. 333–338. Springer International Publishing, Cham (2015). `https://doi.org/10.1007/978-3-319-26123-2_32`

11. Larionov, D., Shelmanov, A., Chistova, E., Smirnov, I.: Semantic Role Labeling with Pre-trained Language Models for Known and Unknown Predicates. In: Mitkov, R., Angelova, G. (eds.) Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019). pp. 619–628. INCOMA Ltd., Varna, Bulgaria (Sep 2019). `https://doi.org/10.26615/978-954-452-056-4_073`

12. Li, X., Chen, H., Liu, C., *et al.*: LLMs Can Also Do Well! Breaking Barriers in Semantic Role Labeling via Large Language Models. In: Che, W., Nabende, J., Shutova, E., Pilehvar, M.T. (eds.) Findings of the Association for Computational Linguistics: ACL 2025. pp. 23162–23180. Association for Computational Linguistics, Vienna, Austria (Jul 2025). `https://doi.org/10.18653/v1/2025.findings-acl.1189`

13. Liu, A., Feng, B., Xue, B., *et al.*: Deepseek-v3 technical report. arXiv preprint arXiv:2412.19437 (2024). `https://doi.org/10.48550/arXiv.2412.19437`

14. Lyashevskaya, O., Kashkin, E.: FrameBank: a database of Russian lexical constructions. In: International Conference on Analysis of Images, Social Networks and Texts. Communications in Computer and Information Science, vol. 542, pp. 350–360. Springer (2015). `https://doi.org/10.1007/978-3-319-26123-2_34`

15. Mel'chuk, I.A.: About the syntactic zero. Typology of passive constructions, diatheses and voices pp. 343–360 (1974)

16. Nikitina, E.N., Onipenko, N.K.: Semantics and pragmatics of statements with psych verbs. Siberian Journal of Philology (2), 271–285 (2022). `https://doi.org/10.17223/18137083/79/19`

17. Nikitina, E.N., Smirnov, I.V.: Predicate-argument structure for intelligent text analysis of social media content. Speech Technology (1-2) (2020), `https://api.semanticscholar.org/CorpusID:256224132`

18. Oberländer, L.A.M., Klinger, R.: Token sequence labeling vs. clause classification for English emotion stimulus detection. In: Gurevych, I., Apidianaki, M., Faruqui, M. (eds.) Proceedings of the Ninth Joint Conference on Lexical and Computational Semantics. pp. 58–70. Association for Computational Linguistics, Barcelona, Spain (Online) (Dec 2020), `https://aclanthology.org/2020.starsem-1.7/`

19. Oberländer, L.A.M., Reich, K., Klinger, R.: Experiencers, stimuli, or targets: Which semantic roles enable machine learning to infer the emotions? In: Proceedings of the Third Workshop on Computational Modeling of People's Opinions, Personality, and Emotion's in Social Media. pp. 119–128 (2020)

20. Onipenko, N.K.: The model of the perspective of subjects (persons) and the problem of classification of egocentric elements. In: The problems of functional grammar: The principle of natural classification, pp. 92–121. Jazyki slavyanskoy kul'tury, Moscow (2013)

21. Paducheva, E.V.: Egocentric units of language and the modes of interpretation. In: Computational linguistics and intellectual technologies. Papers from the Annual International Conference "Dialogue". vol. 1, pp. 486–503. RGGU, Moscow (2013)

22. Senator, F., Lakhfif, A., Zenbout, I., *et al.*: Leveraging ChatGPT for Enhancing Arabic NLP: Application for Semantic Role Labeling and Cross-Lingual Annotation Projection. IEEE Access 13, 3707–3725 (2025). `https://doi.org/10.1109/ACCESS.2025.3525493`

23. Shelmanov, A., Devyatkin, D.: Semantic role labeling with neural networks for texts in Russian. In: Computational Linguistics and Intellectual Technologies. Papers from the Annual International Conference "Dialogue". vol. 1, pp. 245–256. RGGU, Moscow (2017)

24. Testelec, Ja. G.: Introduction to General Syntax. RGGU, Moscow (2001)

25. Troiano, E., Klinger, R., Padó, S.: On the relationship between frames and emotionality in text. Northern European Journal of Language Technology 9(1) (2023). `https://doi.org/10.3384/nejlt.2000-1533.2023.4361`

26. Voloshinov, V.N.: Marxism and the philosophy of language. Priboy, Leningrad (1930)

27. Zhou, Y., Fan, J., Zhang, Q., *et al.*: Modeling semantic-aware prompt-based argument extractor in documents. Applied Sciences 15(10) (2025). `https://doi.org/10.3390/app15105279`

# Document-Level Approach to Extracting Argumentation Structures from the Russian Texts of Scientific Communication

*Elena A. Sidorova*[1] iD *, Irina R. Akhmadeeva*[1] iD *, Daria V. Ilina*[1] iD *,
Irina S. Kononenko*[1] iD *, Alexey S. Sery*[1] iD *, Yury A. Zagorulko*[1] iD

The study addresses the problem of automatic extraction of argumentative structures in scientific communication texts in Russian. Such texts are characterized by a branched logical structure, including distant references and interrelations. To address these complexities, recent methodological advances attempt to leverage the text itself as a contextual foundation for extracting connections. This study presents a generative approach for extracting argumentative relations, reframing the prediction task as a problem of generating marked-up text and making it an end-to-end approach, rather than the traditional pipeline. Two Russian-language corpora were used in the experiments: the translated corpus of microtexts ruMTC and the annotated corpus of scientific communication texts ArgNetSC. A comparative analysis was conducted to evaluate the performance of T5 architecture models trained with supervised fine-tuning (SFT) and Large Language Models on various Russian-language datasets. To facilitate the analysis of long texts, a text segmentation method using a sliding window was proposed. The evaluation revealed that the highest performance in argumentative relation extraction was consistently achieved on the corpus of microtexts. Notably, the smaller models fine-tuned using the SFT method and large language models that were prompted to generate marked texts demonstrated comparable performance ($F_1 \sim 0.32 - 0.37$). For larger texts, however, this trend did not persist, as the FRED-T5 model outperformed all other models with $F_1 \sim 0.23$ on texts related to the genre of scientific articles.

*Keywords: argument mining, document-level argument relation prediction, long-range argumentative relation, text2text generative language model, scientific communication.*

## Introduction

One of the important areas of research in scientific communication, represented by scientific and popular science texts, is the study of the logical organization of reasoning that presents and substantiates the author's position from various points of view. In the process of such reasoning, in order to convince the audience, arguments formulated in the form of premises and conclusions are given in favor of or against the thesis under consideration. The author not only proves some positions through logical reasoning, but also mentally debates with an opponent, modeling possible counterarguments. In this case, a separate argument can act as an initial premise for constructing a new argument, and its conclusion is often used as a justification for another statement. In addition, different arguments can have common premises or conclusions. As a result, they all turn out to be interconnected and form a holistic system that can be represented as an argumentation graph.

Solving problems in the field of argumentation mining requires text corpora with annotated argumentative structures. In recent years, there has been an increase in the volume and diversity of annotated data, but most works are limited to using a few of the most well-known and widely used corpora, according to [11]. Less popular datasets, unfortunately, are often ignored. This is due, first of all, to the desire to compare new methods with existing ones based on uniform benchmarks. However, as the authors of the study note, such a practice is often criticized, since the benchmark data does not always reflect the features of real texts in terms of topic and genre.

[1]A.P. Ershov Institute of Informatics Systems, Novosibirsk, Russian Federation

Not only the language and topic of the text, but also its volume can have a significant impact on the process of extracting argumentation, as the length of the possible argumentative connection increases, hence the number of pairs of statements that can potentially be related.

The largest amount of argumentatively annotated data is available for the English language, while datasets for intellectual analysis of argumentation in Russian are extremely scarce. The following datasets are known for the Russian language:

– Argumentative Microtext Corpus (ruMTC) – a corpus of argumentative essays translated into Russian, with the original argumentative annotation automatically transferred from the original texts [10];

– RuArg-2022 – a corpus of comments from users of the VKontakte social network on news texts about COVID-19 [14]; the annotation model belongs to the APE (Argument Pair Extraction) class, where for a given thesis, supporting and attacking statements are found in different texts; in this corpus, a set of statements is specified, each of which is marked as «for», «against» the thesis, or has a neutral status;

– ArgNetSC – a corpus of scientific communication texts annotated on the ArgNetBankStudio resource based on D. Walton's model [30] being the traditional model of argumentative markup.

While the first two corpora contain rather short texts with contact or short-distance relations between statements, the third corpus contains longer texts as well, that are distinguished by greater structural-content complexity determined by the organization and logical connection of their parts. Scientific communication texts are characterized by a branched logical structure, with the presence of distant references and long-distance relations between content elements. At the same time, argumentative relations are implemented at the level of the entire text, and not only within a sentence, adjacent sentences or paragraphs. To take into account such long-distance relations, Document-Level approaches have recently been actively developing, which use the entire text as a context for finding connections.

In this paper, we propose to use a generative approach to solve the Document-Level Argument relation prediction problem, which, firstly, solves the relation prediction problem not as a classification problem, but as a problem of generating marked-up text, and secondly, uses an end-to-end approach to Argument Mining (E2E-AM) instead of the traditional pipeline [16], in which argument analysis is divided into separate modules trained and applied sequentially. Unlike pipeline frameworks, end-to-end frameworks jointly optimize all subtasks by studying global characteristics and dependencies, which allows us to obtain a holistic view of argument structures. In this paper, we focus on the following research questions.

RQ1. What is the quality of the solution of the End-to-End Argument Mining problem using the Document-Level approach implemented as a text generation task?

RQ2. How does the genre and volume of the text affect the quality of argumentation extraction?

We conducted comparative experiments on two Russian-language datasets: a) the Argumentative Microtext Corpus in Russian (hereinafter – ruMTC), obtained by manually translating the first part of the corpus of the same name from English [10], and b) the ArgNetSC corpus of scientific communication in Russian [30].

The article has the following structure. Section 1 is devoted to a review of the scholarly literature on the research problem. Section 2 describes the datasets used in the study and the data preparation. Section 3 presents the experiment and its result and provides an analysis of

common errors. In Section 4, the results are discussed. Conclusion summarizes the study and points directions for future work.

## 1. Related Work

The main task of argumentation analysis is to extract argumentatively related statements from the text based on formal models. The formal argument model proposed by Toulmin in his work [31] includes 6 components. But in practice, simplified representations of the argument structure are used for data annotating, including 2 components – premises and conclusions. Thus, the argumentative structure can be represented as a binary relation linking a pair of statements, one of which (*premise*) supports or refutes the second (*conclusion*).

There are many datasets, in which such relations were annotated: IAC [32], NoDE (Natural language arguments in online DEbates) [5], UKP-PE [27, 28], RuArg-2022 [14], etc. The Argumentative Microtext Corpus [24] and its Russian-language version [10] were annotated following more complex schemes, reducing, however, the complex arguments to sets of binary relations.

Typically, each corpus consists of texts of a specific genre. For example, the CDCP corpus [22] is used to analyze legal documents, the AbstRCT corpus [21] – for medical-related research, and the DrInventor [15] and SciDTB [1] corpora are used to analyze scientific publications and abstracts, respectively. The UKP-PE corpus of short essays [28] is widely known.

The standard solution of argumentation analysis is to build a pipeline that sequentially solves the following problems: identifying argumentative segments (ADUs), establishing the ADU type, determining the argument type and establishing relations between ADUs [16]. To solve these problems, BERT and BERT-based models are traditionally used [26, 33]. When extracting argumentative relations at short distances (if the premise and conclusion are within the same sentence or in adjacent sentences), the use of discourse markers and argumentation indicators [25], rhetorical relations [2] contributes to improving the results, but they are of little help when extracting long-distance relations.

Recent document-level approaches fall into several categories: sequence labeling methods (e.g., BIOLabel [29]); global context methods that capture document-wide information through question-answering frameworks (DocMRC [19]) or memory mechanisms (MemNet [9]); generative methods that use sequence-to-sequence models (e.g., BART-Gen [17]) for argument extraction, etc.

In general, the pipeline approach has been criticized: pipeline experiments, in general, suffer from the fact that error propagation occurs not only within each step, but also from one to another; the inflexibility of the models used is also noted [35]. In this regard, alternative approaches have recently been developed: end-to-end argument extraction methods based on a network architecture built on a biaffine parser [8, 35] and the use of Text2Text generative models [13]. In [7], an approach based on a biaffine dependency parser was applied to Russian-language texts, which also used rhetorical trees to clarify the boundaries of ADUs.

The idea of considering the AM task as a text generation task arose from related areas of NLP. Thus, the Translation between Augmented Natural Languages methodology [3, 12, 23] uses a pre-trained T5 encoder-decoder model, which has proven its effectiveness in the tasks of extracting relationships between entities, resolving coreference, constructing RST structure, etc.

In recent years, with the growth of pre-training methods, the development of a unified generative structure for solving a variety of tasks within a specific field has attracted increasing attention [6, 18, 20, 34]: solving various subtasks of named entity recognition, information

extraction, tonality analysis and other areas, such as understanding dialogue and multimodal referencing. The paper [4] presents a unified generative platform (UniASA) adapted for various tasks of structured argument analysis: a) E2E-AM Task, b) Argument Pair Extraction (the task is designed to extract pairs of arguments discussing the same point from two interrelated documents), c) Argument Quadruplet Extraction (sentence-level, four-component argument structure used mainly for discussion analysis).

Our work extends a similar approach to E2E-AM by applying it to a Russian-language dataset, characterized by a more complex conceptual and argumentative structure. This adaptation necessitated several key modifications: the development of a novel text annotation scheme, an investigation into the significance of argument sequencing – a problem salient in longer texts that remains unaddressed in prior literature – and the segmentation of lengthy texts into chunks.

## 2. Datasets

Two annotated text corpora were used in the study: 1) **ruMTC** – the first part of the English language Argumentative Microtext Corpus translated into Russian [10] and 2) **ArgNetSC** – the annotated corpus of scientific communication texts.

The corpus of microtexts is widely known and is often mentioned in studies on automatic argumentation analysis. It includes 112 texts (576 sentences) on various topics up to 10 sentences long. Each ADU (in this dataset, each sentence is an ADU) is labeled as supporting or disputing the main thesis of the text; statements are organized into a graph with the following relations: «support», «rebuttal» (attack to an ADU), «undercut» (attack to a relation between statements), «additional» (for combining multiple premises) and «example» (support by example) [24].

A subset of 160 texts was selected from the **ArgNetSC** corpus – a Russian-language corpus with complex argumentative markup. This dataset comprised short and medium-length texts from three subgenres: 30 popular science news texts (**News**), 30 academic paper reviews (**Reviews**), and 100 full-length academic papers (**Articles**). The inclusion of texts with considerable non-argumentative content was found to introduce noise, as they were causing argument-free chunks. Consequently, such texts were systematically filtered out during the dataset compilation process. The length of the texts and the specifics of the genre and topic determine the presence of long-distance argumentative relations, i.e. links between statements that are at least one paragraph apart. The average span of these relations was 330, 502, and 793 characters for the three subcorpora, respectively, notably exceeding their average paragraph lengths (188, 240, and 301 characters). Identifying such relations presents a significant NLP challenge, as the candidate pair space grows combinatorially and grammatical mechanisms become less effective over long distances. Russian-language corpora with argumentative markup, presented in Tab. 1, were used to prepare the data.

**Table 1.** Data statistics

| Corpus | Number of texts | Mean text length (in symbols) | Mean number of words | Mean number of sentences | Mean number of arguments |
|---|---|---|---|---|---|
| ruMTC | 112 | 446 | 61 | 4 | 4 |
| Reviews | 30 | 2860 | 356 | 19 | 20 |
| News | 30 | 4105 | 521 | 28 | 34 |
| Articles | 100 | 9431 | 1165 | 63 | 54 |

Argumentative relations were simplified in the manner shown in Fig. 1 for training data preparation. Each relation $R$ with more than one premise was decomposed into simple binary relations linking each of the premises $p_1, \ldots, p_n$ with a conclusion $C$ of the same type $R$. Each relation between a statement and another argument was replaced by a relation between the statement and the premises of that argument. A simplified argument structure was adopted to establish a baseline, reducing model complexity. While such an approach may lead to a risk of overlooking certain high-order relations, we believe it provided a necessary foundation for future work.



**Figure 1.** Transformation of a complex argument structure into binary relations;
*P – premise, R – relation, C – conclusion*

When making datasets, we developed a specialized annotation scheme, enclosing the structural elements of arguments with special tokens (Fig. 2).



**Figure 2.** An example of a text annotated using special tokens

ADU boundaries were marked with the tags `<adu>` and `</adu>`. The type of argumentative relation was marked with the symbols | (pipe). Two types of argumentative relations were considered: *support* and *attack*. The order of arguments in reference markups is of significant importance. It was observed that if the order of arguments in the markup differed from their order in the text, the model did not learn well. In the current study, arguments in the marked-up text are arranged in the order in which the premises of these arguments appear in the text.

## 3. Experiments

The experiments were conducted using the generative approach, allowing various argument mining tasks to be viewed as the task of generating a set of arguments. The study included a

comparative evaluation of small T5-based models trained using SFT, as well as Large Language Models (LLMs) from the GPT family, across different datasets.

## 3.1. Implementation

Pre-trained language models supporting Russian language were used in the experimental study.

1. **mT5-base** (Multilingual T5, **580M**) is a multilingual model based on the T5 architecture;
2. **ByT5-base** (**582M**) is a modification of T5 that does not use a tokenizer and works directly with UTF-8 bytes; this model can handle any language, is more robust to noise (e.g., typos), and is easier to use because it does not require additional preprocessing;
3. **FRED-T5-large** (**820M**) is a model for Russian language based on the T5 architecture;
4. **Gemma 3 (27B)** is a multilingual and multimodal LLM that supports long context and vision inputs;
5. **gpt-oss: 20b/gpt-oss:120b** are LLMs, which are considered good for performance-efficiency trade-off; particularly the gpt-oss-120b model, which performs comparably to OpenAI's proprietary o4-mini on many benchmarks, while the smaller gpt-oss-20b competes with o3-mini.

We employed SFT to train models of the **T5** architecture. The training was conducted over 20 to 50 epochs, utilizing a starting learning rate of $5 \times 10^{-4}$. To process lengthy documents from the **News**, **Reviews** and **Articles** datasets, a sliding window algorithm was employed for segmentation. This method generated overlapping chunks without regard for inherent textual units (e.g., sentences or paragraphs). The annotated text's length, averaging twice that of the source, dictated a maximum chunk size of half the model's context window. Therefore, the mT5-base model that is pretrained on 512-token sequences, was fed 256-token chunks, and the FRED-T5 model (4096-token context) received 512-token chunks. A larger chunk size for FRED-T5 was not possible due to limited computational resources. Annotations for a chunk were derived only from arguments fully contained within it, which led to the expected loss of long-range or oversized relations. This loss was measured at 21% for the 256-token chunking strategy and 7% for the 512-token strategy.

For the experiments with LLMs, texts were also divided into chunks, which, unlike the previous experiment, were aligned along sentence (paragraph) boundaries. This allowed to exclude incomplete contexts that may introduce noise from the source data for prompt. Note that both LLMs were employed in a zero-shot setting. The experiments were conducted locally using Ollama, the temperature was set to 1.0. Figure **??** shows the prompt template.

The prompt included the role specification that limited the subject area, problem statement and detailed description of the structure of the expected response. When analyzing the results of LLMs, a number of features were identified that allowed us to adjust the prompt.

– The model tended to paraphrase the original text. To counter such behavior we added the following requirement to the description of each component of the markup structure: *Must be an exact quote from the text. Do not paraphrase!!!*.
– The model tended to pay attention only to the main thesis and directly related arguments, so we also added the following requirement: *You must mark up all the text. There should be no fragments of the text that are not present in the argumentation graph.*

| Role specification | You are an expert in analyzing argumentation |
|---|---|
| Task Description | You extract the argumentation in the following CORRECT FORMAT:<br><br><arg><aid>IDENTIFIER</aid> <adu>EVIDENCE</adu> \|TYPE_OF_RELATION\| <adu>CLAIM</adu></arg> |
| Description of the format and argument model | Where:<br><aid> is an unique identifier, for example r0, r1,etc.;<br><adu> are argumentative units (parts of sentences, clauses or whole sentences) which must be an exact (!!!) quotation from the source text, without paraphrasing or changing words.<br>**EVIDENCE** is a statement that serves as the basis for an argument, contains facts, observations, data, rules, or principles.<br><u>It must be an exact quote from the text. Do not paraphrase!!!</u><br>**CLAIM** is a thesis or conclusion that logically follows from a premise. It is the result of reasoning.<br><u>It must be an exact quote from the text. Do not paraphrase!!!</u><br>Not only the main thesis of the document, but also any intermediate conclusions or statements based on **EVIDENCE** can be selected as **CLAIM**.<br>**TYPE_OF_RELATION** — support or attack.<br>Multiple arguments are separated by [SEP]. |
| Result description | The result should be a coherent argument graph where all premises and theses are logically related to each other. |
| Additional notes | The entire text must be marked up. There should be no text fragments left that are not present in the argument graph.<br>**EVIDENCE** must not be the same as **CLAIM** in one argument.<br>**CLAIM** can act as **EVIDENCE** in other arguments.<br>Be careful, **CLAIM** and **EVIDENCE** can be in different parts of the text. |

**Figure 3.** The prompt template for argument extraction

## 3.2. Results

Table 2 summarizes the results of the conducted experiments. We utilized Precision ($P_{adu}$), Recall ($R_{adu}$), and $F1$ score ($F1_{adu}$) as evaluation metrics for ADU Extraction part and $F1$-score for extraction of unlabeled ($F1_{urel}$) and labeled (support, attack) relations ($F1_{rel}$). ADUs were compared at the character level using the Dice coefficient (equivalent to the $F_1$ score), and partial matches above a predefined threshold (equal to 0.8) were considered correct.

For the SFT experiments, the data were randomly partitioned into training, validation, and test sets, comprising 72%, 8%, and 20% of the total data, respectively. To ensure statistical robustness, this procedure was repeated across 10 distinct stratified splits (folds) for each dataset, and a full training-validation-testing cycle was conducted on each fold. LLM-based experiments were conducted on the test sets of each fold. The results, reported in Tab. 2, are presented as the mean performance across all folds with a 95% confidence interval.

The **ByT5** model was applied exclusively to the **ruMTC**-based dataset. Operating directly on UTF-8 bytes without a tokenizer caused the attention tensor to expand rapidly, which – given available computational resources – precluded its application to longer texts from ArgNetSC.

As it can be seen from the results, the highest result of the SFT approach was obtained on microtexts of the **ruMTC** corpus. Small sizes of both source and marked-up texts allowed to fit them completely into the context window of the model without information loss. This can also explain the lack of quality improvement for this dataset in argument extraction when moving to a larger model (**mT5** vs. **FRED-T5**).

LLMs applied to the **ruMTC** data demonstrated performance comparable to that of fine-tuned models. The highest performance for extracting argument relations (without type classification) was achieved by the **gpt-oss-120b** model ($F1_{urel} = 0.42$). When relation types were considered, the **gpt-oss-120b** model performed best on the same data ($F1_{rel} = 0.37$).

On texts of other genres, the fine-tuned models performed better than LLMs, and in total the results were expectedly lower. This may be due to both the presence of non-argumentative

**Table 2.** Experimental results

| Dataset | | Model | ADU Extraction | | | $F1_{urel}$ | $F1_{rel}$ |
|---|---|---|---|---|---|---|---|
| | | | $P_{adu}$ | $R_{adu}$ | $F1_{adu}$ | | |
| ruMTC | | mT5-Base | $0.86 \pm 0.02$ | $0.78 \pm 0.03$ | $0.81 \pm 0.02$ | $0.39 \pm 0.03$ | $0.32 \pm 0.03$ |
| | | ByT5-Base | $0.85 \pm 0.01$ | $0.64 \pm 0.03$ | $0.72 \pm 0.02$ | $0.24 \pm 0.02$ | $0.17 \pm 0.02$ |
| | | FRED-T5-Large | $\mathbf{0.92 \pm 0.03}$ | $\mathbf{0.91 \pm 0.02}$ | $\mathbf{0.91 \pm 0.03}$ | $0.38 \pm 0.06$ | $0.31 \pm 0.05$ |
| | | gpt-oss:20b | $0.77 \pm 0.05$ | $0.75 \pm 0.04$ | $0.74 \pm 0.04$ | $0.33 \pm 0.05$ | $0.29 \pm 0.04$ |
| | | gemma3:27b | $0.86 \pm 0.03$ | $0.82 \pm 0.03$ | $0.84 \pm 0.03$ | $0.34 \pm 0.04$ | $0.30 \pm 0.03$ |
| | | gpt-oss:120b | $0.80 \pm 0.03$ | $0.82 \pm 0.02$ | $0.81 \pm 0.03$ | $\mathbf{0.42 \pm 0.04}$ | $\mathbf{0.37 \pm 0.03}$ |
| ArgNetSC | Reviews | mT5-Base | $\mathbf{0.71 \pm 0.04}$ | $0.52 \pm 0.05$ | $0.58 \pm 0.04$ | $0.10 \pm 0.03$ | $0.09 \pm 0.03$ |
| | | FRED-T5-Large | $0.70 \pm 0.04$ | $0.64 \pm 0.05$ | $0.65 \pm 0.04$ | $\mathbf{0.19 \pm 0.03}$ | $\mathbf{0.18 \pm 0.03}$ |
| | | gpt-oss:20b | $0.60 \pm 0.03$ | $0.56 \pm 0.05$ | $0.56 \pm 0.04$ | $0.06 \pm 0.02$ | $0.04 \pm 0.02$ |
| | | gemma3:27b | $0.62 \pm 0.02$ | $\mathbf{0.78 \pm 0.03}$ | $\mathbf{0.68 \pm 0.02}$ | $0.13 \pm 0.01$ | $0.12 \pm 0.01$ |
| | | gpt-oss:120b | $0.61 \pm 0.03$ | $0.68 \pm 0.06$ | $0.63 \pm 0.04$ | $0.15 \pm 0.02$ | $0.13 \pm 0.02$ |
| | News | mT5-Base | $0.58 \pm 0.05$ | $0.55 \pm 0.03$ | $0.55 \pm 0.03$ | $0.10 \pm 0.02$ | $0.10 \pm 0.02$ |
| | | FRED-T5-Large | $0.64 \pm 0.05$ | $0.57 \pm 0.03$ | $\mathbf{0.59 \pm 0.03}$ | $\mathbf{0.14 \pm 0.01}$ | $\mathbf{0.13 \pm 0.01}$ |
| | | gpt-oss:20b | $0.49 \pm 0.03$ | $0.45 \pm 0.03$ | $0.46 \pm 0.02$ | $0.07 \pm 0.02$ | $0.06 \pm 0.01$ |
| | | gemma3:27b | $\mathbf{0.66 \pm 0.03}$ | $0.49 \pm 0.01$ | $0.55 \pm 0.02$ | $0.11 \pm 0.03$ | $0.11 \pm 0.03$ |
| | | gpt-oss:120b | $0.59 \pm 0.03$ | $\mathbf{0.59 \pm 0.02}$ | $0.58 \pm 0.02$ | $0.13 \pm 0.03$ | $\mathbf{0.13 \pm 0.03}$ |
| | Articles | mT5-Base | $\mathbf{0.70 \pm 0.04}$ | $0.43 \pm 0.06$ | $0.51 \pm 0.05$ | $0.11 \pm 0.007$ | $0.10 \pm 0.006$ |
| | | FRED-T5-Large | $0.70 \pm 0.04$ | $\mathbf{0.83 \pm 0.01}$ | $\mathbf{0.74 \pm 0.02}$ | $\mathbf{0.25 \pm 0.01}$ | $\mathbf{0.23 \pm 0.01}$ |
| | | gpt-oss:20b | $0.47 \pm 0.01$ | $0.62 \pm 0.02$ | $0.52 \pm 0.01$ | $0.08 \pm 0.01$ | $0.08 \pm 0.005$ |
| | | gemma3:27b | $\mathbf{0.73 \pm 0.04}$ | $0.61 \pm 0.07$ | $0.66 \pm 0.05$ | $0.22 \pm 0.03$ | $0.17 \pm 0.03$ |
| | | gpt-oss:120b | $0.54 \pm 0.02$ | $\mathbf{0.83 \pm 0.02}$ | $0.65 \pm 0.02$ | $0.12 \pm 0.01$ | $0.11 \pm 0.01$ |

zones and the loss of long-range relations and the fact that chunks often did not correspond to the logical structure of the text: paragraphs, sections or sentences.

The recall of ADU detection on texts in the ArgNetSC corpus improved significantly when utilizing the larger **FRED-T5 model**, especially on the **Articles** (0.42 vs. 0.82). The average length of an ADUs in **Articles** is greater than in **News** and **Reviews**, so expanding the context window was particularly critical for processing articles. Due to the smaller context size seen during training, some ADUs were «lost» by the **mT5** model.

## 3.3. Analysis of Argumentation Extraction Errors

Error analysis was conducted by experts who participated in the annotation of scientific communication texts. Tables of comparison of arguments found by models with reference ones (type I errors) and tables of comparison of reference arguments with arguments found by the model (type II errors) were considered separately.

### 3.3.1. Segmentation errors

Argumentative analysis involves identifying text fragments that make sense from the point of view of argumentation – ADUs. In general, they are statements based on propositions. Whole sentences are most often identified automatically. However, the analysis of manual segmentation by human annotators shows that these can be both smaller and larger fragments of text, depending on its length and genre.

Scientific and popular science texts (i.e., in this case, the corpus of scientific communication as a whole) are characterized by a high density of argumentation, which contributes to the per-

suasiveness and effectiveness of the impact on the reader. However, in the experiment for various subgenres and models, low values of recall for segmentation can be noted, i.e. identification of a lower density of argumentation compared to human annotation. This corresponds, first of all, to a smaller volume of predicted ADUs as shown in the list below and in Tab. 3.

1. Complex ADUs. Boundaries are set incorrectly in complex ADUs consisting of several sentences.

2. ADUs representing the source of information. Segmentation errors are observed in multi-component structures denoting someone's speech or opinion: direct speech with inserted segment corresponding to the act of speaking *X explains*; reported speech or opinion *According to X*, etc. Separation of the source indicator is an absolute rule for annotators, due to the peculiarities of the analysis of argumentation from the source (expert, witness, etc.), while models either combine speech indicator with the main proposition, or do not isolate such a segment at all.

3. Subordinate clauses. Subordinate clauses connected to the main clause by means of a subordinate conjunction or a conjunction word *where, what, because, as a result of what, since,* etc. are not distinguished into independent ADUs – in this case, the correctness of the segmentation is determined by the semantics of the conjunction/conjunction word (usually a cause-and-effect relationship).

4. Nested phrases. Independent ADUs do not include embedded phrases that represent a comparison or exemplification of the situation in the main sentence.

5. Incomplete Propositions. Collapsed propositions represented by prepositional constructions with a substantive predicate (verbal noun) are not distinguished as independent ADUs.

6. Discontinuous structures. There is a lack of identification of discontinuous structures, the necessity of which is demonstrated by the example in Tab. 3 (it also presents the above-mentioned errors related to exemplification and collapsed proposition).

**Table 3.** Examples of segmentation errors.
Labels (a), (b), etc., denote the ADUs comprising the fragment

| Error type | Expert segmentation | Model segmentation |
|---|---|---|
| Complex ADUs | *There are many methods for detecting a mask on the face, and most of them are a combination of other methods. But they can all be divided into two categories* | *But they can all be divided into two categories.* |
| ADUs representing the source of information | (a) *This is a compelling, evidence-based case for freshwater fishing at the end of the last ice age,* (b) *– Potter noted* | *This is a compelling, evidence-based case for freshwater fishing at the end of the last ice age* |
| Subordinate clauses | (a) *AR devices are the future of surgery,* (b) *because they can significantly reduce the number of medical errors,* (c) *and they can also be used to teach surgery* | *AR devices are the future of surgery, because they can significantly reduce the number of medical errors, and they can also be used to teach surgery* |
| Nested phrases | (a) *and they can also be used to teach surgery,* (b) *as the sports medicine specialist from Aglaya, Christopher Heading, did.* | *and they can also be used to teach surgery, as the sports medicine specialist from Aglaya, Christopher Heading, did.* |
| Incomplete Propositions | (a) *Mennonite dialects are generally recognized as German,* (b) *however, due to their constant migration to regions with other languages or German dialects,* (c) *a simple mention of this is not enough.* | (a) *Mennonite dialects are generally recognized as German,* (b) *however, due to their constant migration to regions with other languages or German dialects, a simple mention of this is not enough.* |
| Discontinuous structures | (a) *the ancestors of the indigenous people of this region, many of whom still depend on freshwater fish* (b) *(salmon, for example)* (a) *may have started subsistence fishing* (c) *in response to declining food resources during long-term climate change.* | *the ancestors of the indigenous people of this region, many of whom still depend on freshwater fish (salmon, for example), may have started subsistence fishing in response to declining food resources during long-term climate change.* |

### 3.3.2. Errors of argument relations extraction

The analysis of the false positive and false negative responses of the models showed the following reasons for the incorrectly extracted relations.

**Common causes.** Expectedly, the facts explained by the causes of false results common to AM tasks were found: proximity in the text of clauses and sentences combined into one pair (one or neighboring sentences); lexical and semantic similarity (presence of the same words and synonyms in the statements of a pair); presence in a pair of relations similar to argumentative ones but not being them.

**Technical reasons for false positive responses.** The use of a generative approach led to the appearance of false positive results due to incorrect segmentation; this type of results also included pairs of statements that, in the expert annotation, are connected by an argumentative relation, but indirectly, through other statements.

**Peculiarities** of recognition of typical reasoning models (schemes of argumentation). The analysis of true- and false-negative responses revealed that the SFT approach has a greater coverage of schemes that are recognized with a quality greater than 10% than the LLM approach – 12 vs. 9. In general, the pairs realizing the Example and Cause to Effect relations are well recognized by the largest number of models on the largest number of subcases.

The relations Expert Opinion, Part to Whole and Sign are recognized in some cases most well (up to 28%), in other cases very poorly (up to complete absence of true positives).

***Expert Opinion***. This scheme is well recognized on the **News** subcorpus by SFT models; both LLM models for the same subcorpus and **mT5** models for the **Articles** subcorpus performed poorly. The low results are unexpected, since our previous experiments, which solved the problem of binary classification of pairs of statements, showed high results for this scheme (up to $F1_{urel}$=0.89). The analysis showed that the errors are related either to insufficient explication of the argument components or to insufficiently detailed segmentation performed by the models.

***Part to Whole***. This relation is well recognized by the **FRED-T5** model, but the other models perform poorly. The scheme contains two premises and a conclusion: *m is a species (part) of n; m has property G → n has property G.* The errors can be partially explained by the fact that experts use the ***Part to Whole*** scheme as a transitive rather than argumentative scheme when marking up the text, in order to preserve the integrity of the graph. Therefore, an isolated pair of sentences out of context may not contain argumentation, and the models naturally fail to detect it.

***Sign***. The ***Sign*** relation is well recognized by LLM models and poorly by SFT models, including **FRED-T5**. This reasoning model also has three components: *B is generally indicated as true when its sign, A, is true; A (a finding) is true in this situation → B is true in this situation.* False negative examples of this scheme, as a rule, had no explicit relation indicators, or the indicators are polysemous (may indicate not only argumentative relation: e.g., parenthetical explanations, markers *associated with*, *since*), or the segmentation of annotators is too fractional, and the selected segments are not informative enough for models. Another possible cause of errors is the simplification of arguments during data preparation, which made an isolated premise statement insufficient to support the conclusion.

***Negative Consequences***. In examples implementing this reasoning model, the models most frequently make errors in determining the type of relation (support or attack). This is because in the AIF ontology this scheme is supporting (*If A is implemented, bad consequences*

*are likely to occur → A should not be implemented*), but due to the implicitness of some reasoning components, this scheme is more often implemented as attacking.

Table 4 demonstrates examples of the above-described errors made by the models in extracting argument relations.

**Table 4.** Examples of Argument Relation Extraction errors.
FP denotes False Positives, FN denotes False Negatives

| Error type | Errors in Argument Relation Prediction | Comments |
|---|---|---|
| Incorrect segmentation (**FP**) | (a) *There are inaccuracies in the article,* \|**attack**\| (b) *5) There are inaccuracies in the article,* | Pairs of identical statements have been merged, either in full or truncated form. |
| Incorrect segmentation (**FP**) | (a) *5) The article* \|**support**\| (b) *The material of the article raises a number of questions:* | The statements correspond to pairs from the reference dataset, but one of the statements is truncated. |
| Mediated argument relation (**FP**) | \<arg\>\<aid\>r6\</aid\>\<adu\> *Since the main goal of our research was to compare the obtained results with the data presented in the «Russian Associative Dictionary»*\</adu\>\|**support**\|\<adu\>*In this study, the free association experiment was used as the main method, in which the respondent is required to give an unrestricted response to a stimulus word in the form of a response word or phrase.*\</adu\>\</arg\> | In the expert annotation, the pairs of statements are mediated by others; in this example, the statement reports an intermediate research objective: *The focus of the research interest was on identifying, within the mini-group, matches with the most frequent responses recorded in the lexicographic source.* |
| Relation similar to argumentative ones but not being them (**FP**) | (a) *The author has collected interesting material (the corpus of examples found in each translation is given in the article), but limits himself exclusively to its quantitative analysis, presented in two tables.* \|**support**\| (b) *The reviewed article is devoted to the study of archaisms and historicisms in five Russian translations of The Song of Roland.* | The candidate premise is a detail of the candidate conclusion |
| Peculiarities of ***Part to Whole*** reasoning model (transitive scheme, **FN**) | (a) *The ATT&CK Matrix for Enterprise information security threat assessment methodology has the following merits* \<...\> \|**support**\| (b) *It helps to understand what tools attackers use, to familiarize with their techniques and tactics.* \|**support**\| (c) *This knowledge allows predicting the likely point of entry into organizations.* | Statement (b) is not a premise to statement (a) because it clarifies rather than proves it, but since statement (c) supports (b), the link between (a) and (b) is necessary to demonstrate the support for statement (a) by statement (c) in the graph. |
| Peculiarities of **Sign** reasoning model (**FN**) | (a) *No clarity in the interpretation of the term «language».* (b) *In the first sentence of the abstract we read: «one of the Germanic languages», which is known «in the world» as «Lower Germanic language (???) (dialect???) of the Mennonites».* (c) *Below it is labeled as «vernacular» for both Siberian and Canadian Mennonites.* | Thesis (a) is supported by two premises (b) and (c). Together (b) and (c) are able to prove the thesis, but separately they are insufficient to support (a). |
| Peculiarities of ***Negative Consequences*** reasoning model (Incorrect Relation type) | (a) *However, when dealing with large amounts of data, training diffusion models can be time-consuming and require large computational resources* \|**attack**\| (b) *Generalizing, diffusion models allow generating an image from a textual description by sequentially varying the noise in pixel space.* | If we explicate all the statements, the first fragment should be divided into two: (1) However, when dealing with large amounts of data, training diffusion models may not be appropriate (conflicts logically with the second statement) and (2) It can be time-consuming and computationally intensive (supporting premise to statement (1)). Thus, the model has generated a markup with a correct label, but it contradicts the «reference»markup. |

## 4. Discussion

**RQ1** discussion. In most cases, the generative approach to argument extraction shows rather low results on Russian-language texts, which is generally consistent with the work of other

researchers. For example, in [4] on the English-language **MTC** corpus the $F1_{rel}$ score is 0.35, while on the Russian-language **ruMTC** corpus in our experiment $F1_{rel} = 0.37$.

On the positive side, a comprehensive analysis of the whole text through the chunking mechanism allowed us to include about 98% of all argumentative relations. Some problems arise at the boundaries of the chunks due to the arbitrariness of text partitioning, and aligning the chunk boundaries to the boundaries of sentences/paragraphs leads to an unstable context window size in the training sample, which has a bad effect on model training.

**RQ2** discussion. The analysis of the dependence of the quality of argument extraction on the text length shows that such correlation takes place only for ultra-small texts, when it is possible to fully place the text in the context window (**ruMTC** vs. **ArgNetSC**). On texts larger than one chunk this ceases to play a significant role. For larger texts, genre features seem to play an important role, in particular, the complexity and type of argumentation used, the coverage of the text with argumentation, the size of the argumentative statement, etc. Thus, the experimental results show better performance when analyzing scientific articles than on other subgenres of the **ArgNetSC** corpus, despite their larger size (see Tab. 1). This can be explained by the fact that scientific articles have a stricter organization of the presentation of the material. In addition, there is less complex argumentation in this sub-corpus (in the **Articles**, 21% of binary relations are obtained as a result of simplification, in **Reviews** – 25%, and in **News** – 26%), so errors arising from the simplification of such argumentation are less frequent in this corpus.

## Conclusion

The main goal of the study was to test new document-level generative approaches for solving the problem of argumentation analysis and long-range argumentative relation extraction.

Experiments conducted on Russian-language data showed that the highest results were achieved on microtexts, with small models fine-tuned with SFT showing approximately the same quality as large language models running on a specialized prompt. However, for large texts of scientific genres, this trend does not hold and the best results are obtained with the trained **FRED-T5** model.

To analyze long texts, a technique was proposed to segment them into chunks by a sliding window. The size of the chunk depended on the context window of the model used. This approach guarantees consideration of relations «fitting» into such a window.

The main types of errors that reduce argument extraction performance were identified: errors related to segmentation, to establishing an argument relation between two ADUs, and to determining the type of relation. The models often fail to recognize ADUs consisting of more than one sentence, presenting subordinate clauses, comparative turns, explanatory turns with examples, discontinuous structures, and indications of the source of information. The errors in establishing an argumentative relation between statements occurred both due to segmentation errors and general quality-reducing factors (lexico-semantic similarity, contact of statements between which a relation is falsely established, etc.), peculiarities of individual argumentation schemes. In addition, some of the erroneously identified relations in the expert markup are argumentatively related through other statements, but such relations could not be reflected in the dataset due to the peculiarities of the experiment. The most frequent errors in identifying the type of relation are explained by the implicitness of part of the reasoning elements, due to which supporting schemes are sometimes used as attacking schemes.

In further development of the approach, a stage for identifying the main thesis will be added, which is likely to adapt the document-level argument mining task for cases of long-range links over two paragraphs.

We have published our code at the GitHub[2]. All datasets used in this article are publicly available from each distributor.

## Acknowledgements

## References

1. Accuosto, P., Saggion, H.: Mining arguments in scientific abstracts with discourse-level embeddings. Data & Knowledge Engineering 129, 101840 (2020). `https://doi.org/10.1016/j.datak.2020.101840`

2. Akhmadeeva, I.R., Kononenko, I., Sidorova, E., Shestakov, V.: Using rhetorical structures to analyze argumentation in scientific communication texts. Computational Linguistics and Intellectual Technologies (2025), `https://api.semanticscholar.org/CorpusID:280935169`

3. Bao, J., He, Y., Sun, Y., *et al.*: A generative model for end-to-end argument mining with reconstructed positional encoding and constrained pointer mechanism. In: Goldberg, Y., Kozareva, Z., Zhang, Y. (eds.) Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. pp. 10437–10449. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates (dec 2022). `https://doi.org/10.18653/v1/2022.emnlp-main.713`

4. Bao, J., Jing, M., Dong, K., *et al.*: UniASA: A Unified Generative Framework for Argument Structure Analysis. Computational Linguistics 51(3), 739–784 (09 2025). `https://doi.org/10.1162/coli_a_00553`

5. Cabrio, E., Villata, S.: Node: A benchmark of natural language arguments. In: Computational Models of Argument, pp. 449–450. IOS Press (2014). `https://doi.org/10.3233/978-1-61499-436-7-449`

6. Chen, Z., Chen, L., Chen, B., *et al.*: UniDU: Towards a unified generative dialogue understanding framework. In: Lemon, O., Hakkani-Tur, D., Li, J.J., *et al.* (eds.) Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue. pp. 442–455. Association for Computational Linguistics, Edinburgh, UK (sep 2022). `https://doi.org/10.18653/v1/2022.sigdial-1.43`

7. Chistova, E.: End-to-end argument mining over varying rhetorical structures. In: Rogers, A., Boyd-Graber, J., Okazaki, N. (eds.) Findings of the Association for Computational Linguistics: ACL 2023. pp. 3376–3391. Association for Computational Linguistics, Toronto, Canada (jul 2023). `https://doi.org/10.18653/v1/2023.findings-acl.209`

---

[2]`https://github.com/Inscriptor/doc-level-approach-arg-extraction.git`

8. Dozat, T., Manning, C.D.: Deep biaffine attention for neural dependency parsing. In: International Conference on Learning Representations. Toulon, France (apr 2017), `https://openreview.net/forum?id=Hk95PK9le`

9. Du, X., Li, S., Ji, H.: Dynamic global memory for document-level argument extraction. In: Muresan, S., Nakov, P., Villavicencio, A. (eds.) Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 5264–5275. Association for Computational Linguistics, Dublin, Ireland (may 2022). `https://doi.org/10.18653/v1/2022.acl-long.361`

10. Fishcheva, I., Kotelnikov, E.: Cross-Lingual Argumentation Mining for Russian Texts. In: van der Aalst, W.M.P., Batagelj, V., Ignatov, D.I., *et al.* (eds.) Analysis of Images, Social Networks and Texts. pp. 134–144. Springer International Publishing, Cham (2019). `https://doi.org/10.1007/978-3-030-37334-4_12`

11. Galassi, A., Lippi, M., Torroni, P.: Multi-task attentive residual networks for argument mining. IEEE/ACM Transactions on Audio, Speech, and Language Processing 31, 1877–1892 (2023). `https://doi.org/10.1109/TASLP.2023.3275040`

12. Hu, X., Wan, X.: RST Discourse Parsing as Text-to-Text Generation. IEEE/ACM Transactions on Audio, Speech, and Language Processing 31, 3278–3289 (2023). `https://doi.org/10.1109/TASLP.2023.3306710`

13. Kawarada, M., Hirao, T., Uchida, W., Nagata, M.: Argument mining as a text-to-text generation task. In: Graham, Y., Purver, M. (eds.) Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 2002–2014. Association for Computational Linguistics, St. Julian's, Malta (mar 2024). `https://doi.org/10.18653/v1/2024.eacl-long.121`

14. Kotelnikov, E., Loukachevitch, N., Nikishina, I., Panchenko, A.: RuArg-2022: Argument Mining Evaluation. pp. 333–348 (06 2022). `https://doi.org/10.28995/2075-7182-2022-21-333-348`

15. Lauscher, A., Glavaš, G., Ponzetto, S.P.: An argument-annotated corpus of scientific publications. In: Slonim, N., Aharonov, R. (eds.) Proceedings of the 5th Workshop on Argument Mining. pp. 40–46. Association for Computational Linguistics, Brussels, Belgium (nov 2018). `https://doi.org/10.18653/v1/W18-5206`

16. Lawrence, J., Reed, C.: Argument mining: A survey. Computational Linguistics 45(4), 765–818 (01 2020). `https://doi.org/10.1162/coli_a_00364`

17. Li, S., Ji, H., Han, J.: Document-level event argument extraction by conditional generation. In: Toutanova, K., Rumshisky, A., Zettlemoyer, L., *et al.* (eds.) Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 894–908. Association for Computational Linguistics, Online (jun 2021). `https://doi.org/10.18653/v1/2021.naacl-main.69`

18. Li, Z., Lin, T.E., Wu, Y., *et al.*: UniSA: Unified Generative Framework for Sentiment Analysis. In: Proceedings of the 31st ACM International Conference on Multimedia. pp. 6132–6142. MM '23, Association for Computing Machinery, New York, NY, USA (2023). `https://doi.org/10.1145/3581783.3612336`

19. Liu, J., Chen, Y., Xu, J.: Machine reading comprehension as data augmentation: A case study on implicit event argument extraction. In: Moens, M.F., Huang, X., Specia, L., Yih, S.W.t. (eds.) Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. pp. 2716–2725. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic (nov 2021). `https://doi.org/10.18653/v1/2021.emnlp-main.214`

20. Lu, Y., Liu, Q., Dai, D., *et al.*: Unified structure generation for universal information extraction. In: Muresan, S., Nakov, P., Villavicencio, A. (eds.) Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 5755–5772. Association for Computational Linguistics, Dublin, Ireland (may 2022). `https://doi.org/10.18653/v1/2022.acl-long.395`

21. Mayer, T., Marro, S., Cabrio, E., Villata, S.: Enhancing evidence-based medicine with natural language argumentative analysis of clinical trials. Artificial Intelligence in Medicine 118, 102098 (2021). `https://doi.org/10.1016/j.artmed.2021.102098`

22. Niculae, V., Park, J., Cardie, C.: Argument mining with structured SVMs and RNNs. In: Barzilay, R., Kan, M.Y. (eds.) Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 985–995. Association for Computational Linguistics, Vancouver, Canada (jul 2017). `https://doi.org/10.18653/v1/P17-1091`

23. Paolini, G., Athiwaratkun, B., Krone, J., *et al.*: Structured prediction as translation between augmented natural languages. In: International Conference on Learning Representations (2021), `https://openreview.net/forum?id=US-TP-xnXI`

24. Peldszus, A., Stede, M.: An annotated corpus of argumentative microtexts. In: Argumentation and Reasoned Action: Proceedings of the 1st European Conference on Argumentation, Lisbon 2015. vol. 2, pp. 801–815. College Publications, London (2016)

25. Sidorova, E., Akhmadeeva, I., Zagorulko, Y., *et al.*: An integrated approach to the analysis of argumentative relationships in scientific communication texts. Ontology of Designing 13(4), 562–579 (12 2023). `https://doi.org/10.18287/2223-9537-2023-13-4-562-579`, (in Russian)

26. Srivastava, P., Bhatnagar, P., Goel, A.: Argument Mining using BERT and Self-Attention based Embeddings. In: 2022 4th International Conference on Advances in Computing, Communication Control and Networking (ICAC3N). pp. 1536–1540 (2022). `https://doi.org/10.1109/ICAC3N56670.2022.10074559`

27. Stab, C., Gurevych, I.: Annotating argument components and relations in persuasive essays. In: Tsujii, J., Hajic, J. (eds.) Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers. pp. 1501–1510. Dublin City University and Association for Computational Linguistics, Dublin, Ireland (aug 2014), `https://aclanthology.org/C14-1142/`

28. Stab, C., Gurevych, I.: Parsing argumentation structures in persuasive essays. Computational Linguistics 43(3), 619–659 (09 2017). `https://doi.org/10.1162/COLI_a_00295`

29. Strubell, E., Verga, P., Andor, D., *et al.*: Linguistically-informed self-attention for semantic role labeling. In: Riloff, E., Chiang, D., Hockenmaier, J., Tsujii, J. (eds.) Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. pp. 5027–5038. Association for Computational Linguistics, Brussels, Belgium (oct – nov 2018). `https://doi.org/10.18653/v1/D18-1548`

30. Timofeeva, M., Ilina, D., Kononenko, I.: Argumentative annotation of the scientific Internet-communication corpus: Genre analysis and study of typical reasoning models based on the ArgNetBank Studio platform. NSU Vestnik 22(1), 27–49 (2024). `https://doi.org/10.25205/1818-7935-2024-22-1-27-49`, (in Russian)

31. Toulmin, S.E.: The Uses of Argument. Cambridge University Press, 2 edn. (2003). `https://doi.org/10.1017/CBO9780511840005`

32. Walker, M., Tree, J.F., Anand, P., *et al.*: A corpus for research on deliberation and debate. In: Calzolari, N., Choukri, K., Declerck, T., *et al.* (eds.) Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12). pp. 812–817. European Language Resources Association (ELRA), Istanbul, Turkey (may 2012), `https://aclanthology.org/L12-1643/`

33. Xu, H., Ashley, K.: Multi-granularity argument mining in legal texts. In: Legal Knowledge and Information Systems, pp. 261–266. Frontiers in Artificial Intelligence and Applications, IOS Press (2022). `https://doi.org/10.3233/FAIA220477`

34. Yan, H., Gui, T., Dai, J., *et al.*: A unified generative framework for various NER subtasks. In: Zong, C., Xia, F., Li, W., Navigli, R. (eds.) Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 5808–5822. Association for Computational Linguistics, Online (aug 2021). `https://doi.org/10.18653/v1/2021.acl-long.451`

35. Ye, Y., Teufel, S.: End-to-end argument mining as biaffine dependency parsing. In: Merlo, P., Tiedemann, J., Tsarfaty, R. (eds.) Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume. pp. 669–678. Association for Computational Linguistics, Online (apr 2021). `https://doi.org/10.18653/v1/2021.eacl-main.55`

# Can LLMs Get to the Roots? Evaluating Russian Morpheme Segmentation Capabilities in Large Language Models

*Dmitry A. Morozov*[1] (iD)*, Anna V. Glazkova*[2] (iD)*, Boris L. Iomdin*[3] (iD)

Automatic morpheme segmentation, a crucial task for morphologically rich languages like Russian, is persistently hindered by a significant drop in performance on words containing out-of-vocabulary (OOV) roots. This issue affects even state-of-the-art models, such as fine-tuned BERT models. This study investigates the potential of modern Large Language Models (LLMs) to address this challenge, focusing on the specific task of root identification in Russian. We evaluate a diverse set of eight state-of-the-art LLMs, including proprietary and open-weight models, using a prompt-based, few-shot learning approach. The models' performance is benchmarked against strong baselines – a fine-tuned RuRoberta model and a CNN ensemble – on a 500-word test set. Our results demonstrate that one model, Gemini 2.5 Pro, surpasses both baselines by approximately 5 percentage points in root identification accuracy. An examination of the model's reasoning capabilities shows that while it can produce logically sound, etymologically-informed analyses, it is also highly prone to factual hallucinations. This work highlights that while LLMs show significant promise in overcoming the OOV root problem, the inconsistency of their reasoning presents a significant obstacle to their direct application, underscoring the need for further research into improving their factuality and consistency.

*Keywords: morpheme segmentation, tokenizers, large language models, Russian language.*

## Introduction

Word segmentation into minimal meaningful substrings, morphemes, is important for morphologically rich languages. The construction of morpheme segmentations can be applied in language learning, for building hypotheses about possible word etymology, or as a subword tokenizer. In the latter capacity, morpheme-oriented tokenizers can be used for low-resource languages as an alternative to common BPE tokenizers [2, 9, 11, 14].

The need for automation in this task arises from the incompleteness of existing morpheme dictionaries. For instance, in Russian, the largest dictionaries contain no more than half of the words found in the Russian National Corpus [5]. At the same time, algorithmic approaches today have an accuracy that, on average, is not inferior to expert annotation [13]. However, existing algorithms have a significant drawback: a low quality of performance with roots not encountered in the training sample [7, 13]. This problem can be partially solved by using pre-trained BERT-like models, but the quality of annotation for words containing out-of-vocabulary (OOV) roots is still significantly lower than the average, and the vast majority of annotation errors are specifically related to identifying root boundaries in such words [12].

A potential solution to the problem of identifying OOV roots could be the use of large language models (LLMs) for word segmentation and root finding [1, 17]. However, the applicability of LLMs to morpheme segmentation remains insufficiently investigated.

In this work, we focus on Russian, the language that is both well-represented in the training corpora of state-of-the-art LLMs and well-studied within the context of morpheme segmentation. Since the primary challenge for existing methods is the identification of OOV roots, we decided to concentrate on the applicability of LLMs specifically to the task of word root identification.

---

[1]Novosibirsk State University, Novosibirsk, Russian Federation
[2]University of Tyumen, Tyumen, Russian Federation
[3]Käthe-Kollwitz-Gymnasium, Berlin, Germany

The rationale is that if the root can be successfully identified, the remainder of the word can be segmented with very high accuracy using established algorithms.

Our main contributions are as follows:

- LLMs can outperform other approaches in identifying word roots, yet the quality of segmentation using them is still far from ideal;
- examining the content of the reasoning fields from a linguistic perspective showed that in a number of cases, the segmentation variant proposed by the model is logically justified and more suitable than the dictionary-based one.

The rest of the paper is structured in the following way. Section 1 contains a brief review of related work on automatic morpheme segmentation. Section 2 describes the dataset and the models utilized for the experiments. Section 3 presents the experimental results and discussion. Section 3 concludes this paper, summarizing the study and pointing directions for further work.

## 1.  Related Work

Automatic morpheme segmentation comprises two main paradigms: surface segmentation, where a word is segmented into its constituent substrings (e.g., funniest → funn-i-est), and canonical segmentation, which aims to restore the underlying forms of the morphemes (e.g., funniest → fun-y-est) [6]. For both tasks, machine learning-based algorithms have demonstrated the highest efficacy.

The task of surface segmentation has been extensively studied for the Russian language. The problem is typically framed as a character-level classification task, where each character is assigned a two-part label. The first part of the label indicates the character's position within a morpheme, while the second specifies the morpheme type. Among traditional approaches not leveraging pre-trained models, strong results have been achieved using Convolutional Neural Networks (CNNs) [20] and Long Short-Term Memory (LSTM) networks [4], with the performance of CNNs on random word samples being comparable to expert-level annotation. The use of convolutional networks has also shown strong results for other languages [15, 19].

Fine-tuning BERT-like models on Russian data has further improved performance, achieving a character-level accuracy exceeding 97% and a perfect-segmentation rate for words over 92% on random samples [12]. This approach has also yielded strong results for other Slavic languages, including Belarusian and Czech. Furthermore, using pre-trained models partially addresses the key limitation of CNN and LSTM architectures: the sharp decline in performance on words containing roots not seen in the training data. Nevertheless, the challenge of OOV roots remains critical, with a performance drop of over 15% in word-level accuracy for such cases [13].

For canonical segmentation, state-of-the-art algorithms were benchmarked in the SIGMOR-PHON 2022 shared task [3]. The top-performing systems were developed by the DeepSPIN team using LSTM and Transformer-based models [16], closely followed by the CLUZH team with models based on neural transducers [22]. During testing on nine languages (English, Spanish, Hungarian, French, Italian, Russian, Czech, Latin, and Mongolian), participants achieved morpheme-level F1-scores exceeding 93.5 for all languages, with scores surpassing 99 for three languages, including Russian. However, subsequent testing of the DeepSPIN approach confirmed that, as with surface segmentation, OOV roots remain the primary challenge [12].

Notable existing methods utilizing LLMs include the LLMSegm algorithm [17]. In this framework, a pre-trained Glot500 model [8] performs binary classification for each potential boundary position within a word. Despite its significant computational complexity (requiring $N-1$ LLM

inferences for a word of length $N$), this method has outperformed previous techniques for several low-resource South African languages. Conversely, end-to-end generation of segmentations for another low-resource language, Bribri, using an LLM (Claude Sonnet 3.7), while more efficient on average, proved to be less effective for multi-morphemic words than a non-pretrained algorithm [1]. Therefore, the application of LLM-based approaches to morpheme segmentation is currently under-explored.

## 2. Data and Models

### 2.1. Data

The task of morpheme segmentation for the Russian language is complicated by the absence of a unified approach among linguists for defining what constitutes a morphemic analysis. In practice, the choice of a dataset for experiments determines the segmentation paradigm, and an algorithm trained on one dataset will inherently produce segmentations that are incorrect from the perspective of another. At the same time, previous studies have shown that the performance of algorithms is generally similar across different datasets [7, 13]. The largest machine-readable datasets for Russian morpheme segmentation currently available are Morphodict-K (based on the "Dictionary of Morphemes of the Russian Language" (ed. A. I. Kuznetsova and T. F. Efremova) [10]), Morphodict-T (based on the "Word Formation Dictionary of the Russian Language" (ed. A. N. Tikhonov) [21]), and the Russian dataset from the SIGMORPHON 2022 Shared Task on Morpheme Segmentation [3]. In the present study, we chose to use the Morphodict-K dataset. This decision was based on two reasons:

- The Morphodict-K and the Morphodict-T datasets use the same notation: words are segmented into morphemes with an indication of the type of each of them, in total, 7 types of morphemes are used: PREF (prefixes), ROOT, SUFF (suffixes), END (endings), LINK (connecting vowels), POST (postfixes), HYPH (hyphens). This makes it easy to move from experiments with one dataset to experiments with another. However, the segmentation in Morphodict-K is based on etymology and features a high degree of morpheme granularity. Although this paradigm is not strictly formalized, it allows for significantly less subjectivity than the paradigm underlying the Morphodict-T dataset. The latter employs a non-obvious criterion of the transparency of derivational chains in modern Russian. This leads to the identification of different roots in words, which relatedness is obvious to native speakers (e.g., in *dobro* 'goodness', the root is identified as *-dobr-*, whereas in *odobrenie* 'approval', the root is identified as *-odobr-*). The use of such a criterion reduces the internal consistency of the annotation, which is also reflected in the lower performance of automated methods on this dataset.
- Upon closer inspection, the SIGMORPHON dataset is annotated in a highly inconsistent manner. Specifically, there is no uniform approach to segmenting the infinitive suffix *-t'*- or the adjectival ending *-yy-*; in about half of the cases, they are merged with the preceding morpheme. This internal inconsistency, along with the absence of documented segmentation rules, makes this dataset a poor candidate for our research.

Therefore, within the scope of our study, we aimed to identify the etymological root specifically. This approach may be more valuable for creating morpheme-oriented tokenizers for training language models, as it offers more fine-grained segmentation and results in a smaller token vocabulary.

The utilized dataset was split by word roots into training and test sets at an approximate 4:1 ratio. Words containing multiple roots were preliminarily removed from the dataset. However, due to financial constraints, the LLM testing was not performed on the entire test set, but on a random sample of 500 words from it. For baseline models, we calculated the quality both on the entire test sample and on the selected 500 words. A summary of the characteristics of the dataset and the samples is provided in Tab. 1.

**Table 1.** Brief characteristics of the datasets

| Dataset | Morphodict-K | Train set | Test set | 500 test words |
|---|---|---|---|---|
| Unique words | 75 649 | 51 620 | 12 401 | 500 |
| Unique morphemes | 8 079 | 5 606 | 1 662 | 434 |
| Unique roots | 7 148 | 4 768 | 1 192 | 275 |
| Avg morphemes per word | 4.12 | 3.95 | 3.98 | 4.07 |
| Avg morpheme occurrence | 38.56 | 36.36 | 29.71 | 4.69 |
| Avg root occurrence | 12.24 | 10.83 | 10.40 | 1.82 |
| Avg characters in root | 4.62 | 3.64 | 3.74 | 3.73 |

## 2.2. Prompt-based models

To ensure the representativeness of the study, we used a variety of state-of-the-art multilingual general-purpose LLMs developed and trained by different teams. The access to the models was provided via the OpenRouter API.

Some of the models used are proprietary, accessible only through API, and architecture and training details are unknown. Using such models inevitably leads to a worse understanding of the reasons for the effectiveness of a particular model, but the high performance of these models in independent benchmarks necessitates their inclusion for a comprehensive evaluation. These models include:

1. Claude Sonnet 4[4];
2. Gemini 2.5 Pro and Gemini 2.5 Flash Lite[5];
3. Mistral Medium 3.1[6];
4. GPT 5 Chat[7] (we used this model instead of the base GPT 5 because GPT 5 was still not available in the OpenRouter API at the time of our experiments.).

We compared these models to the models whose weights are available on HuggingFace:

1. Llama 4 Maverick[8], an auto-regressive language model that uses a mixture-of-experts (MoE) architecture with 128 experts and has 17B activated parameters (400B parameters in total);
2. gpt-oss-120b[9] with 117B parameters and 5.1B active parameters;
3. Qwen3-235B-A22B[10], MoE model with 128 experts (8 activated experts) and 235B parameters in total (22B parameters activated).

---

[4]https://www.anthropic.com/claude/sonnet
[5]https://deepmind.google/models/gemini/
[6]https://mistral.ai/news/mistral-medium-3
[7]https://openai.com/index/introducing-gpt-5/
[8]https://huggingface.co/meta-llama/Llama-4-Maverick-17B-128E-Instruct
[9]https://huggingface.co/openai/gpt-oss-120b
[10]https://huggingface.co/Qwen/Qwen3-235B-A22B

It should also be noted that the models used differ significantly in the costs of API requests. This is due to three reasons: first, different tokenization of requests and responses, second, different token prices, and third, the presence of reasoning, which significantly affected the volume of generated text. For the Claude and Gemini models, we forced the reasoning mode to be enabled. The other models used do not have such a setting via the OpenRouter API, but we recorded non-empty reasoning in the response in the case of gpt-oss-120b (for all 500 words) and Qwen3-235B-A22B (for 35 words). The cost of a request varied from $10^{-4}$ to $10^{-2}$ USD. In total, less than 25 USD were spent on experiments (including preliminary ones). A temperature value of 0 was used during generation for all models.

## 2.3. Prompt

To construct the prompt, we used the segmentation paradigm description from the original dataset[11] and a series of examples from the training set.

During preliminary experiments with models Gemini 2.5 Flash Lite, Mistral Medium 3.1, and gpt-oss-120b, we tested different strategies for selecting examples presented to the model in the prompt. In particular, we attempted to use randomly selected examples from the training set, as well as words similar to the analyzed word by a substring at the beginning or end of the word. We also tried different numbers of examples in the prompt (from 10 to 100). These strategies turned out to be insufficiently effective. Increasing the number of examples with a random selection strategy improved the quality only slightly, and when adding words with a similar substring, the models tended to overfit.

The best result was achieved by adding a small number of examples to the prompt, illustrating various features of the segmentation paradigm: consideration of etymology and a high degree of morpheme granularity (in comparison, for example, with the approach of the "Word Formation Dictionary of the Russian Language" (ed. A. N. Tikhonov) [21]. Examples were selected iteratively; the final prompt included 13 words from the training sample: *"nasekomoe"* ('insect', root *-sek-*), *"ulybat'sya"* ('to smile', root *-lyb-*), *"revolyutsiya"* ('revolution', root *-volyuts-*), *"vostochnyy"* ('eastern', root *-toch-*), *"obidet'sya"* ('to be offended', root *-obid-*), *"pozlashchat'"* ('to gild', root *-zlashch-*), *"obratit'"* ('to turn', root *-obrat-*), *"bytovoy"* ('domestic, household', root *-by-*), *"nenastnyy"* ('inclement', root *-nast-*), *"izverzhenie"* ('eruption', root *-verzh-*), *"uproshchat'sya"* ('to be simplified', root *-proshch-*), *"truzhenichestvo"* ('hard work', *-truzh-*), *"annulirovanie"* ('cancellation', root *-nul-*). An example of word formatting is shown in Fig. 1.

Additionally, the prompt included requirements describing the surface segmentation procedure: the concatenation of morphemes must exactly match the original word, and in cases of alternations, the root must be extracted in the exact form in which it appears in the word. Finally, during preliminary experiments, we determined that the segmentation quality improves when the response includes a complete segmentation. The final prompt was written in Russian. In Fig. 2 we provide a translation of the original prompt into English

## 2.4. Baselines

As baselines we considered two approaches, originally developed for constructing surface segmentation with prediction of the type of each morpheme. The task of constructing morpheme

---

[11]https://ruscorpora.ru/en/page/instruction-derivation

```
- Source word: "nasekomoe"
  JSON output:
  {
    {
      "original_word": "nasekomoe",
      "etymological_root": "sek",
      "morphemic_analysis": "na-sek-om-oe"
    }
  }
```

**Figure 1.** An example of word formatting

segmentation is considered as a task of character-level classification. The choice of these models was based on their high quality in previous studies [12, 13, 20]. In our study, we also trained models on the surface segmentation task, and then extracted the morpheme marked by the model as the root of the word. The baselines are:

1. **Convolutional neural network ensemble** [20]. The ensemble consists of three identical convolutional networks trained independently. We used the original implementation of the algorithm. The number of convolutional layers in each convolutional network was 3, the number of filters was 192. Each of the models was trained for 25 epochs on an AMD Ryzen 5 5600X CPU.

2. **Fine-tuned RuRoberta model** [12]. We fine-tuned `ruRoberta-large`[12] (355M parameters) [23] for token-classification task. The input sequence consisted of the lemma itself and sequence of the separated word letters. The output sequence is '0' for the lemma and letter class for each letter. For implementation we used the `simpletransformers`[13] framework [18]. The batch size during training was set to 16, and the learning rate was set to 4e-6. The values of the remaining parameters were set to default. We fine-tuned the model for 30 epochs on an Nvidia RTX 4090 GPU.

## 3. Results and Discussion

The results we obtained are presented in Tab. 2 below, where the LLMs are ordered by the decreasing number of correctly identified roots. The performance of the baseline models is also included for comparison.

Notably, Gemini 2.5 Pro was the only model to outperform the baselines, achieving a root accuracy nearly 5% higher. The remaining LLMs underperformed compared to the baseline models. Furthermore, we observed that proprietary models generally surpassed open-weight models.

Our analysis of LLM errors revealed no single, consistent pattern of incorrect root identification; the models tended to err on different words rather than the same ones (Fig. 3). For instance, only 185 words were correctly segmented by all eight models, while 474 words were segmented correctly at least once (meaning 26 words were never segmented correctly by any model).

---

[12]https://huggingface.co/ai-forever/ruRoberta-large
[13]https://simpletransformers.ai/

You are a linguistic expert specializing in morphemic and etymological
    analysis of the Russian language. Your task is to conduct morphemic
    analysis of words, strictly following the principles of the "Dictionary
    of Morphemes of the Russian Language" by A. I. Kuznetsova and T. F.
    Efremova. Your answer should always be in JSON format.

Your task is to perform a morphemic analysis of a word and identify its
    etymological root. Follow these rules:

0. **Exact match:** The morpheme segmentation of the word MUST be a letter-
    for-letter match with the original word and must not change any letters.
1. **Principle of granularity:** Divide the word into morphemes (prefixes,
    root, suffixes, endings) in as much detail as possible.
2. **Historical root:** Identify historical and etymological roots, even if
    their connection to the modern meaning is not obvious to a native
    speaker.
3. **Structural correlation:** Identify morphemes (including the root) if
    there are other words in the language with a similar structure or
    morphemes, even if the word is not used without them (e.g., "u-lyb-at'
    sya" by analogy with "u-smekh-at'sya").
4. **Analysis of loanwords:** Segment borrowed stems if there are other
    lexemes in the Russian language with similar structural elements (e.g.,
    "re-volyuts-iya" and "e-volyuts-i-ya").
5. **Handling of the soft sign:** If a soft sign follows the root, include
    it in the root (e.g., "kol'-ts-o").
6. **Alternations:** If an alternation is allowed in the root, you must
    choose the spelling that is present in the word (e.g., in the word "
    pozlashchat'" the root is "zlashch").

Here are some reference examples:

[EXAMPLES]

Now, please process the following word:

Source word: "{word_to_analyze}"

Provide the answer strictly in the following JSON format, without including
    any other explanations. The etymological root MUST be part of the full
    morpheme segmentation.
{{
  "original_word": "<the word being processed>",
  "etymological_root": "<the etymological root>",
  "morphemic_analysis": "<the full morpheme segmentation with hyphens>"
}}

**Figure 2.** Translation of the used prompt into English

**Table 2.** Experiment results

| Model | Correct roots | Root-level accuracy | Fully correct segmentations | Word-level accuracy |
|---|---|---|---|---|
| Gemini 2.5 Pro | 430 | 0.86 | 392 | 0.78 |
| Mistral Medium 3.1 | 391 | 0.78 | 349 | 0.69 |
| Claude Sonnet 4 | 386 | 0.77 | 294 | 0.59 |
| Gemini 2.5 Flash Lite | 355 | 0.71 | 279 | 0.56 |
| GPT 5 Chat | 343 | 0.69 | 251 | 0.50 |
| Llama 4 Maverick | 341 | 0.68 | 305 | 0.61 |
| Qwen3 235B A22B | 335 | 0.67 | 242 | 0.48 |
| gpt-oss-120b | 334 | 0.67 | 228 | 0.46 |
| fine-tuned ruRoberta-large | 401 | 0.80 | 387 | 0.77 |
| CNN ensemble | 406 | 0.81 | 358 | 0.72 |



**Figure 3.** Number of words for which the root was correctly identified by $N$ models

The primary challenge for the LLMs was the root *-sta-* in words such as *zastava* 'outpost', *nastavlyat'* 'to instruct', and *predstavitel'skiy* 'representative'. In these cases, all LLMs incorrectly identified the root as *-stav-*. However, the correct root *-sta-* was successfully identified by at least some models in words like *perestavat'* 'to cease' and *ostanovit'sya* 'to stop'. The current experimental design does not allow us to draw general conclusions about the relationship between specific root features and the quality of their identification. We plan to replicate this experiment with a larger dataset in the future to investigate this issue further.

Interestingly, even the top-performing model, Gemini 2.5 Pro, incorrectly processed some words for which the other seven models provided the correct answer, such as *"vaflya"* ('waffle') and *"prorab"* ('foreman'). However, an analysis of the responses showed that once the model had violated the response format: for *"vaflya"*, it identified the root *-vafl'-*, where the soft sign might be included from a phonetic standpoint but is incorrect within a surface segmentation paradigm.

Since Gemini 2.5 Pro was the only model to surpass the baseline approach, we decided to focus our error analysis on this model as the most promising. A preliminary analysis revealed that in several cases, a discrepancy between the predicted root and the reference root might not indicate a model error but rather an inaccuracy in the dictionary or the possibility of a different interpretation. Consequently, we conducted a detailed manual analysis of 70 words for which the root identified by the model differed from the reference. We evaluated the model's responses based on two criteria:

1. **Is the model's answer more suitable than the reference?** A score of 2 indicated the model's answer was better, 1 meant it was difficult to choose the better option, and 0 meant the reference was better.

2. **Is the model's reasoning factually and logically sound?** A score of 2 meant the reasoning represented a well-conducted word-formation analysis containing only correct etymological facts; 1 indicated that the reasoning contained a mix of correct statements and hallucinations; and 0 meant the reasoning consisted entirely of hallucinations.

The analysis showed that the root predicted by the model was more suitable than the reference in 13 out of 70 cases, and in another 11 cases, choosing the better option was difficult. Cases where the model's prediction was more accurate included words such as:

- *"anonimnost'"* 'anonymity', reference root *-nim-*, predicted *-onim-*: in this case, -onim- is better suited as a root, since words such as *"anonim"* 'anonymous' and *"sinonim"* 'synonymous' contain the borrowed ancient Greek root $-ονομα-$, that is, *-o-* cannot be considered either a prefix or a linking vowel;

- *"oblechennyy"* 'endowed', reference root *-oblech-*, predicted *-lech-*: from the point of view of the segmentation paradigm used in the Morphodict-K dataset, it is reasonable to isolate the historical prefix *-ob-*;

- *"podkrylok"* 'fender liner', reference root *-kry-*, predicted *-kryl-*: in the Proto-Slavic language, this root was supposedly contained in a form containing *-l-* (*\*kridlo*, *\*skrdlo*), so its isolation here is incorrect.

The model provided adequate reasoning in 12 cases and partially correct reasoning in 29. However, it is important to note that correct reasoning coincided with a better-identified root in only 6 of these cases (Tab. 3).

**Table 3.** Analysis of predicted roots different from reference ones for the Gemini 2.5 Pro model

| Root correctness \ Reasoning validity | 0 | 1 | 2 |
|---|---|---|---|
| 0 | 27 | 16 | 3 |
| 1 | 1 | 7 | 3 |
| 2 | 1 | 6 | 6 |

On three occasions, correct reasoning led to a root that was less suitable than the reference:

- *"ukomplektovyvat'sya"* 'to be staffed/equipped' and *"doukomplektovat'sya"* 'to become fully staffed/equipped', reference root *-komplekt-*, predicted *-plekt-*: despite the fact that etymologically this substring contains the prefix *-kom-* (in Latin *"completus"* 'complete'), this root was borrowed into the Russian language in its current form through Polish and,

before that, German, and there are no other Russian words with substring *-plekt-*, so the rule about lexemes with similar structural elements does not apply here;

- *"sovetizirovat'sya"* 'to become sovietized', reference root *-vet-*, predicted *-sovet-*: from the point of view of the segmentation paradigm used in the Morphodict-K dataset, it is reasonable to isolate the historical prefix *-so-*.

In some instances, the model's reasoning arrived at the correct root, yet a different option was chosen for the final answer (e.g., for *"rabota"* 'work' the reasoning pointed to the root *-rab-*, but the answer given was *-rabot-*).

The majority of Gemini 2.5 Pro's errors relate to the incorrect handling of root alternations (e.g., *-treb-/-trebl-*, *-yav-/-yavl-*), insufficient consideration of etymology (for instance, the model failed to connect the words *"stat'"* ('to become/stand') and *"stavit'"* ('to put/place')), or the excessive segmentation of loanwords in cases where there is no basis for segmentation in the Russian language (*interes → inter-es*, *krendel → krend-el*, or the aforementioned *komplekt → kom-plekt*). Furthermore, the reasoning often contains fabricated facts or flawed logical transitions (even when the identified root is correct), which should be considered a significant drawback for the potential integration of the model into lexicography for automating dictionary creation.

In addition to analyzing the model's errors, we conducted a zero-shot evaluation using Gemini 2.5 Pro to evaluate the example selection strategy. For this, we removed the examples from the prompt, leaving the rest of the instructional text unchanged. We then generated roots for the test set. The model produced 420 correct roots and 387 fully correct analyses. This resulted in a correct root rate of 84%, compared to 86% achieved with the few-shot approach. However, we observed an increased proportion of incorrectly formatted responses among the model's errors in the zero-shot setting. The model often appended hyphens or extraneous explanations to the root, and in two instances, it generated an empty response. Therefore, despite the small difference in metrics, the few-shot approach enhances the stability of the response format, which can be critical for practical applications.

## Conclusion

A key challenge for state-of-the-art automatic morpheme segmentation algorithms is their poor performance on words containing roots that were not present in the training data. This paper presents an investigation into the potential of using Large Language Models (LLMs) to overcome this limitation. For our experiments, we utilized the Russian-language Morphodict-K dataset and a range of multilingual, general-purpose LLMs, including the most current proprietary models. The Russian language was selected because it is, on the one hand, well-represented in the training corpora of these LLMs and, on the other, a well-studied language for the morpheme segmentation task.

We compared the efficacy of LLMs against two strong baselines – a fine-tuned BERT-like model and an ensemble of convolutional neural networkson the specific task of word root identification. Using the Gemini 2.5 Pro model, we successfully surpassed the baselines by 5 percentage points in accuracy. A subsequent linguistic analysis of this model's errors revealed that in several instances, the root predicted by the LLM was more suitable than the reference one from the dataset. An examination of the model's reasoning fields showed that it is sometimes capable of justifying its choices with factual evidence. However, it frequently generates reasoning that

contains hallucinations and fabricated facts. This should be considered a significant drawback for the potential integration of the model into lexicography for automating dictionary creation.

The limitations of this study include the relatively small size of the test set (500 words), the focus on a single target language, and the limited number of prompting strategies explored. In addition, a significant limitation of the overall approach lies in its speed and cost: although we do not query the LLM for every possible boundary position in a word, each word is processed using a separate query to the model via the API. Despite these constraints, our approach managed to outperform state-of-the-art baselines, which underscores the need for more extensive and larger-scale research in this domain.

# References

1. Anderson, C., Nguyen, M., Coto-Solano, R.: Unsupervised, semi-supervised and LLM-based morphological segmentation for Bribri. In: Mager, M., Ebrahimi, A., Pugh, R., *et al.* (eds.) Proceedings of the Fifth Workshop on NLP for Indigenous Languages of the Americas (AmericasNLP). pp. 63–76. Association for Computational Linguistics, Albuquerque, New Mexico (May 2025). `https://doi.org/10.18653/v1/2025.americasnlp-1.7`

2. Asgari, E., Kheir, Y.E., Javaheri, M.A.S.: MorphBPE: A Morpho-Aware Tokenizer Bridging Linguistic Complexity for Efficient LLM Training Across Morphologies (2025), `https://arxiv.org/abs/2502.00894`

3. Batsuren, K., Bella, G., Arora, A., *et al.*: The SIGMORPHON 2022 shared task on morpheme segmentation. In: Nicolai, G., Chodroff, E. (eds.) Proceedings of the 19th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology. pp. 103–116. Association for Computational Linguistics, Seattle, Washington (Jul 2022). `https://doi.org/10.18653/v1/2022.sigmorphon-1.11`

4. Bolshakova, E., Sapin, A.: Bi-LSTM model for morpheme segmentation of Russian words. In: Ustalov, D., Filchenkov, A., Pivovarova, L. (eds.) Artificial Intelligence and Natural Language. pp. 151–160. Springer International Publishing, Cham (2019). `https://doi.org/10.1007/978-3-030-34518-1_11`

5. Bonch-Osmolovskaya, A., Gladilin, S., Kozerenko, A., *et al.*: Russian National Corpus 2.0: corpus platform, analysis tools, neural network models of data markup. In: Computational Linguistics and Intellectual Technologies. Papers from the Annual International Conference "Dialogue" (01 2025). `https://doi.org/10.28995/2075-7182-2025-23-57-73`

6. Cotterell, R., Vieira, T., Schütze, H.: A joint model of orthography and morphological segmentation. In: Knight, K., Nenkova, A., Rambow, O. (eds.) Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 664–669. Association for Computational Linguistics, San Diego, California (Jun 2016). `https://doi.org/10.18653/v1/N16-1080`

7. Garipov, T., Morozov, D., Glazkova, A.: Generalization ability of CNN-based Morpheme Segmentation. In: 2023 Ivannikov Ispras Open Conference (ISPRAS). pp. 58–62 (2024). `https://doi.org/10.1109/ISPRAS60948.2023.10508171`

8. Imani, A., Lin, P., Kargaran, A.H., *et al.*: Glot500: Scaling multilingual corpora and language models to 500 languages. In: Rogers, A., Boyd-Graber, J., Okazaki, N. (eds.) Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 1082–1117. Association for Computational Linguistics, Toronto, Canada (Jul 2023). `https://doi.org/10.18653/v1/2023.acl-long.61`

9. Kildeberg, M.W., Schledermann, E.A., Larsen, N., van der Goot, R.: From Smør-re-brød to Subwords: Training LLMs on Danish, One Morpheme at a Time (2025), `https://arxiv.org/abs/2504.01540`

10. Kuznetsova, A.I., Efremova, T.F.: Dictionary of Morphemes of the Russian Language. Russkii yazyk, Moscow (1986)

11. Matthews, A., Neubig, G., Dyer, C.: Using morphological knowledge in open-vocabulary neural language models. In: Walker, M., Ji, H., Stent, A. (eds.) Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers). pp. 1435–1445. Association for Computational Linguistics, New Orleans, Louisiana (Jun 2018). `https://doi.org/10.18653/v1/N18-1130`

12. Morozov, D., Astapenka, L., Glazkova, A., Garipov, T., Lyashevskaya, O.: BERT-like models for Slavic morpheme segmentation. In: Che, W., Nabende, J., Shutova, E., Pilehvar, M.T. (eds.) Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 6795–6815. Association for Computational Linguistics, Vienna, Austria (Jul 2025). `https://doi.org/10.18653/v1/2025.acl-long.337`

13. Morozov, D., Garipov, T., Lyashevskaya, O., *et al.*: Automatic morpheme segmentation for Russian: Can an algorithm replace experts? Journal of Language and Education 10(4), 71–84 (Dec 2024). `https://doi.org/10.17323/jle.2024.22237`

14. Nzeyimana, A., Niyongabo Rubungo, A.: KinyaBERT: a morphology-aware Kinyarwanda language model. In: Muresan, S., Nakov, P., Villavicencio, A. (eds.) Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 5347–5363. Association for Computational Linguistics, Dublin, Ireland (May 2022). `https://doi.org/10.18653/v1/2022.acl-long.367`

15. Olbrich, M., Žabokrtský, Z.: Morphological segmentation with neural networks: Performance effects of architecture, data size, and cross-lingual transfer in seven languages. In: Ekštein, K., Konopík, M., Pražák, O., Pártl, F. (eds.) Text, Speech, and Dialogue. pp. 275–286. Springer Nature Switzerland, Cham (2026). `https://doi.org/10.1007/978-3-032-02551-7_24`

16. Peters, B., Martins, A.F.T.: Beyond characters: Subword-level morpheme segmentation. In: Nicolai, G., Chodroff, E. (eds.) Proceedings of the 19th SIGMORPHON Workshop on

Computational Research in Phonetics, Phonology, and Morphology. pp. 131–138. Association for Computational Linguistics, Seattle, Washington (Jul 2022). `https://doi.org/10.18653/v1/2022.sigmorphon-1.14`

17. Pranjić, M., Robnik-Šikonja, M., Pollak, S.: LLMSegm: Surface-level morphological segmentation using large language model. In: Calzolari, N., Kan, M.Y., Hoste, V., *et al.* (eds.) Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024). pp. 10665–10674. ELRA and ICCL, Torino, Italia (May 2024), `https://aclanthology.org/2024.lrec-main.933/`

18. Rajapakse, T.C., Yates, A., de Rijke, M.: Simple transformers: Open-source for all. In: Proceedings of the 2024 Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region. pp. 209–215. SIGIR-AP 2024 (2024). `https://doi.org/10.1145/3673791.3698412`

19. Sorokin, A.: Improving Morpheme Segmentation Using BERT Embeddings. In: Burnaev, E., Ignatov, D.I., Ivanov, S., *et al.* (eds.) Analysis of Images, Social Networks and Texts. pp. 148–161. Springer International Publishing, Cham (2022). `https://doi.org/10.1007/978-3-031-16500-9_13`

20. Sorokin, A., Kravtsova, A.: Deep convolutional networks for supervised morpheme segmentation of Russian language. In: Ustalov, D., Filchenkov, A., Pivovarova, L., Žižka, J. (eds.) Artificial Intelligence and Natural Language. pp. 3–10. Springer International Publishing, Cham (2018). `https://doi.org/10.1007/978-3-030-01204-5_1`

21. Tikhonov, A.N.: Word Formation Dictionary of the Russian language [Slovoobrazovatelnyi slovar russkogo yazyka]. Russkiy yazyk, Moscow (1990)

22. Wehrli, S., Clematide, S., Makarov, P.: CLUZH at SIGMORPHON 2022 shared tasks on morpheme segmentation and inflection generation. In: Nicolai, G., Chodroff, E. (eds.) Proceedings of the 19th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology. pp. 212–219. Association for Computational Linguistics, Seattle, Washington (Jul 2022). `https://doi.org/10.18653/v1/2022.sigmorphon-1.21`

23. Zmitrovich, D., Abramov, A., Kalmykov, A., *et al.*: A family of pretrained transformer language models for Russian. In: Calzolari, N., Kan, M.Y., Hoste, V., *et al.* (eds.) Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024). pp. 507–524. ELRA and ICCL, Torino, Italia (May 2024), `https://aclanthology.org/2024.lrec-main.45/`

# RuBookSum: Dataset for Russian Literature Abstractive Summarization

*Denis A. Grigoriev*[1] , *Daniil V. Khudiakov*[1], *Daniil I. Chernyshev*[1]

The majority of existing Russian document summarization datasets focus on short-form source documents which does not require complex causal analysis or coreference resolutions. Furthermore, processing longer multi-page texts poses a serious challenge to current generation of language models as the limited context window complicates response generation by demanding additional task partitioning. To lay the groundwork for future research of the problem, we introduce RuBookSum, an abstractive summarization dataset for Russian long-form narrative summarization. Our dataset covers documents from various literature domains, including fiction, classic, children books and popular science, and includes high-quality human-written summaries. To establish a baseline, we evaluate popular open-source large language models and provide comprehensive analysis on their performance. Additionally, we propose optimized algorithms for long-document summarization, which enable up to 300% summary generation speed up without significant drops in quality.

*Keywords: large language model, summarization, literature, books.*

## Introduction

Automatic text summarization is one of the key tasks in natural language processing. The goal is to create an informative synopsis of the source text while preserving its main meaning. In recent years, with the advent of large language models (LLMs), interest in automating summarization has increased across many genres, including fiction. Unlike scientific, news, or technical texts, fiction is characterized by high stylistic and semantic complexity. Non-linear storytelling, imagery, metaphor, and stylistic devices make synopsis writing especially challenging. The limited context window of modern models further complicates processing long texts as it imposes additional text generation constraints and demands additional task partitioning.

At present moment there are not many datasets focusing specifically on summarizing fiction, and available collections focus on non-Russian material. BookSum [3] is one of the first and best-known English-language datasets for abstractive summarization of narrative works. It contains books, plays, and short stories paired with summaries of varying granularity (paragraph level, chapter level, book level). Echoes from Alexandria [8] is a multilingual corpus of fiction, including five languages: English, German, French, Italian, and Spanish. FABLES [2] is a hand-curated corpus designed to evaluate factual faithfulness of summaries for book-length fiction. It includes 3 158 claims extracted from LLM-generated summaries for 26 books. Each claim is evaluated across model outputs by experts. According to FABLES, even advanced models (e.g., Claude) commit 20–30% factual errors, including distorted causal relations, incorrect characterization of protagonists, and overemphasis on minor details, judged by three criteria: agreement with original events, logical correctness, and absence of distortions.

To study the specifics of Russian long-narrative automatic summarization, we introduce RuBookSum, a dataset for Russian literature abstractive summarization. Our dataset contains high-quality human-written summaries for documents of different domains including fiction, popular science, children books and classical literature. To demonstrate the issues of multi-page text summarization, we conduct an extensive evaluation of popular open-source large language

---

[1]Lomonosov Moscow State University, Research Computing Center, Moscow, Russia

models. Our analysis finds that only largest models exceeding 100 billion parameters are able to fully comprehend long-range causal relations, while smaller models only capture general semantics. Additionally, to adapt the models for the task, we propose new abstractive summarization algorithms optimized for long-document processing. Compared to existing approaches new methods achieve up to 300% summary generation speed up while retaining the same level of quality.

The article is organized as follows. Section 1 describes the RuBookSum dataset. In Section 2 we present the summarization methods, including the hierarchical approach with node filtering and the blueprint-based approach with question clustering. Section 3 defines the evaluation metrics. Section 4 contains the experimental setup. Section 5 reports and analyzes the results. Conclusion summarizes the study and points directions for further work.

Code and data are publicly available[2].

## 1. Dataset

At the moment of the study, there were no publicly available corpora designed specifically for Russian literature summarization. To address the lack of resources, a new dataset was created, using "Narodny Briefly" platform [7] where users publish summaries for popular books. The summaries vary in length (from a few sentences to several paragraphs) and in style: some reproduce key phrases verbatim, while others use free-form narration. Some cover the whole work, others split content by chapter. Usually they contain the main facts and conclusions from the source text, but may include the author's commentary.

To collect the respective summary sources, we leveraged Librusec digital library [4], one of the largest Russian-language online book collections. Each text underwent automatic preprocessing: meta-information (e.g., titles, chapter descriptions, technical inserts) were removed, then the text was formatted into a unified, standardized form suitable for use with models.

To better link books with their summaries, cosine similarity was used: the author name text from Briefly [7] and from Librusec [4] was embedded via SentenceTransformers with the model[3] and compared using cosine similarity. The summaries were automatically cleaned of HTML tags, comments, and service markers using LLM Meta-Llama 3-70B-Instruct. Then Librusec was searched and a collection of "book text — summary" pairs was formed.

The resulting dataset includes:
- 600+ cleaned user summaries from "Narodny Briefly" [7];
- 40+ different genres;
- source works from the Librusec digital library [4].

**Table 1.** Dataset overview

| Dataset | Number of documents | Avg. document length (# words) | Avg. summary length (# words) | Compression ratio (summary length / text length) |
|---|---|---|---|---|
| **RuBookSum** | 634 | 35 052.64 | 700.77 | 8.43% |
| BookSum | 405 | 112 885.15 | 1 167.20 | 0.79% |
| Gazeta | 60 964 | 632.77 | 41.94 | 6.99% |

---

[2] https://github.com/Nejimaki-Tori/RuBookSum
[3] https://huggingface.co/deepvk/USER-bge-m3

**Figure 1.** Distribution of texts by genres (top 10; "Other" aggregates all remaining genres)

The genre distribution in the collection is shown in Fig. 1 and Tab. 1 gives dataset statistics versus analogs.

## 2. Methodology

### 2.1. Hierarchical Method

Most common method (Algorithm 1) for long-document summarization [11] splits the text into chunks and generates a local summary for each chunk. First, the document is split into chunks and each chunk is summarized, yielding the list $S_0$. Then, at level $\ell$, GROUPSUMMARIES takes the current list $S_\ell$ and partitions it into non-overlapping groups. For each group, MERGEGROUP produces a single higher-level summary, forming $S_{\ell+1}$. Because at least some groups contain two or more elements, $|S_{\ell+1}| < |S_\ell|$, and the process terminates when a single root summary remains, which we report as the document-level summary.

### 2.2. Node Filtering Optimization

The classical hierarchical method constructs the final summary through a multi-layered combination of intermediate summaries derived from individual text chunks. However, literature often contains chunks that have little impact on plot development or contains redundant information without any additional details. During the generation of the final summary, these chunks can shift the narrative towards repetitive content, thus reducing the overall informativeness.

To address this issue, a node filtering mechanism based on cosine similarity was implemented (Algorithm 2). To eliminate low-informative or redundant chunks, we evaluate cosine similarity between all intermediate summaries at each hierarchy level. Chunks that are close in cosine similarity to previous ones are considered redundant and are excluded from compilation of the summary at the current level. We compute a pairwise similarity matrix PAIRWISESIMILARITY$(S_\ell)$ using cosine similarity of summary embeddings (the first element is always retained). This modification aims to accelerate generation by removing potentially superfluous parts of information, thereby increasing the salient detail density in the final summaries.

**Algorithm 1** Hierarchical method

**Input:** $W$ – model context window, $D$ – input text of length $L \gg W$, $p_\theta$ – model, $C$ – chunk length

  Split $D$ into chunks $c_1 \ldots c_{\lceil \frac{L}{C} \rceil}$

  $\ell \leftarrow 0$

  $S_\ell \leftarrow \{c_1 \ldots c_{\lceil \frac{L}{C} \rceil}\}$

  **repeat**

    $Groups \leftarrow GroupSummaries(S_\ell)$

    $\ell \leftarrow \ell + 1$

    $S_\ell \leftarrow \{\}$

    **for** $g \in Groups$ **do**

      $S_\ell \leftarrow S_\ell \cup \{MergeGroup(p_\theta, g)\}$

    **end for**

  **until** $|S_\ell| = 1$

  **return** $S_\ell[0]$

---

**Algorithm 2** Hierarchical method with node filtering

**Input:** $W$ – model context window, $D$ – input text of length $L \gg W$, $p_\theta$ – model, $\theta$ – similarity threshold, $C$ – chunk length

  Split $D$ into chunks $c_1 \ldots c_{\lceil \frac{L}{C} \rceil}$

  $\ell \leftarrow 0$

  $S_\ell \leftarrow \{c_1 \ldots c_{\lceil \frac{L}{C} \rceil}\}$

  **repeat**

    $M \leftarrow PairWiseSimilarity(S_\ell)$

    $S_\ell \leftarrow \{s_i : \ s_i \in S_\ell \wedge$

        $(\max_{j<i} M_{ij} < \theta \vee i = 0)\}$

    $Groups \leftarrow GroupSummaries(S_\ell)$

    $\ell \leftarrow \ell + 1$

    $S_\ell \leftarrow \{\}$

    **for** $g \in Groups$ **do**

      $S_\ell \leftarrow S_\ell \cup \{MergeGroup(p_\theta, g)\}$

    **end for**

  **until** $|S_\ell| = 1$

  **return** $S_\ell[0]$

## 2.3. Text-Blueprint

This method [1] is essentially a modification of the hierarchical method that improves summary robustness by building an intermediate outline before text generation (Algorithm 3). The outline is formed as a set of question-answer pairs, which enhances the controllability of the generation process and ensures the structured nature of the result. First the model creates a list of questions reflecting key events, themes, and characters. Then short answers are automatically generated for each question. Given a group $g$ at merge level $\ell$, GENERATEBLUEPRINT produces a set of question–answer pairs. This structure serves as a blueprint used to produce the final summary. The function SUMWITHBP uses generated blueprint and given chunks to produce a higher-level summary. Full prompt templates and additional Q/A examples are provided in Appendix A.

## 2.4. Question Clustering Optimization

The baseline blueprint implementation generates a question-answer outline for each chunk and at each merge level. With fiction, however, questions produced for different chunks may overlap and yield conflicting answers, which in turn corrupts merging process, making the summary less structured and complete. Moreover, generating an outline at every step slows the method down and consumes extra computational time. To address the issue, we add additional question clustering step aimed at reducing merge level content overlap (Algorithm 4). The obtained question clusters are generalized using the same summary generation LLM to produce universal question-answer outline. We add the following functions to modify blueprint method: EXTRACTQUESTIONS builds $Q'$, CLUSTERIZE forms clusters from $Q'$ and GENERALIZE maps each cluster to one generalized question.

---

**Algorithm 3** Blueprint method

**Input:** $W$ – model context window, $D$ – input text of length $L \gg W$, $p_\theta$ – model, $C$ – chunk length, $R$ – length limit

Split $D$ into chunks $c_1 \ldots c_{\lceil \frac{L}{C} \rceil}$

$\ell \leftarrow 0$

$S_\ell \leftarrow \{c_1 \ldots c_{\lceil \frac{L}{C} \rceil}\}$

**repeat**            ▷ Merging summaries

    $Groups \leftarrow GroupSummaries(S_\ell)$

    $\ell \leftarrow \ell + 1$

    $S_\ell \leftarrow \{\}$

    **for** $g \in Groups$ **do**

       **if** $Length(g) > R$ **then**

          $b_i \leftarrow GenerateBlueprint(p_\theta, g)$

          $S_\ell \leftarrow S_\ell \cup \{SumWithBp(p_\theta, b_i, g)\}$

       **else**

          $S_\ell \leftarrow S_\ell \cup \{g\}$

       **end if**

    **end for**

**until** $|S_\ell| = 1$

**return** $S_\ell[0]$

---

**Algorithm 4** Blueprint method with clustering

**Input:** $W$ – model context window, $D$ – input text of length $L \gg W$, $p_\theta$ – model, $C$ – chunk length, $R$ – length limit

Split $D$ into chunks $c_1 \ldots c_{\lceil \frac{L}{C} \rceil}$

$\ell \leftarrow 0$

$S_\ell \leftarrow \{c_1 \ldots c_{\lceil \frac{L}{C} \rceil}\}$

**for** $c \in S_\ell$ **do**

    $b \leftarrow GenerateBlueprint(p_\theta, c)$

    $Q' \leftarrow \{ExtractQuestions(p_\theta, b)\}$

**end for**

$K \leftarrow Clusterize(Q')$

**for** $k_i \in K$ **do**

    $Q \leftarrow Q \cup \{Generalize(p_\theta, k_i)\}$

**end for**

**repeat**            ▷ Merging summaries

    $Groups \leftarrow GroupSummaries(S_\ell)$

    $\ell \leftarrow \ell + 1$

    $S_\ell \leftarrow \{\}$

    **for** $g \in Groups$ **do**

       $S_\ell \leftarrow S_\ell \cup \{SumWithBp(p_\theta, Q, g)\}$

    **end for**

**until** $|S_\ell| = 1$

**return** $S_\ell[0]$

---

## 3. Metrics

For an objective comparison of the described methods and models in the task of literature summarization, four metrics were considered.

**ROUGE-L** [5] is based on the length of the longest common subsequence (LCS) between the generated summary $S$ and the reference $R$.

$$\text{Precision} = \frac{\text{LCS}(S, R)}{|S|}, \tag{1}$$

$$\text{Recall} = \frac{\text{LCS}(S, R)}{|R|}, \tag{2}$$

$$\text{ROUGE-L} = \frac{2\,\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}. \tag{3}$$

**BERTScore** [14]. For every token pair from prediction and reference we compute cosine similarity of their embeddings. Then:

$$P = \frac{1}{|S|} \sum_{t \in S} \max_{u \in R} \text{sim}(e_t, e_u), \tag{4}$$

$$R = \frac{1}{|R|} \sum_{u \in R} \max_{t \in S} \text{sim}(e_u, e_t), \tag{5}$$

$$\text{BERTScore} = \frac{2\,P \cdot R}{P + R}, \tag{6}$$

---

where $R$ is the reference text and $S$ is the generated one.

**Key question coverage (Coverage)** is the proportion of questions that are answered in reference that are covered in generated summary:

$$\text{Coverage} = \frac{\#\{q_i \colon P(\text{yes} \mid q_i, S) > 0.75\}}{N}, \tag{7}$$

where $N$ is the total number of questions and $P(\text{yes} \mid q_i, S)$ is the probability that the answer to $q_i$ is present in $S$, obtained with an LLM.

**Factual agreement (Agreement)** is the average cosine similarity between answers $a_i^{\text{pred}}$ generated based on predicted summary and answers $a_i^{\text{ref}}$ obtained from reference summary to the same questions from Coverage metric:

$$\text{Agreement} = \frac{1}{N} \sum_{i=1}^{N} \text{sim}(a_i^{\text{pred}}, a_i^{\text{ref}}), \tag{8}$$

where `sim` is cosine similarity of embeddings.

## 4. Experimental Setup

All measurements were performed on the test split of the dataset, comprising one sixth ($\approx 16.7\%$) of the corpus. For all methods the summaries were limited to 500 words maximum. The input text was split into fixed-size chunks of 2 000 tokens which were obtained by `AutoTokenizer` from `DeepPavlov/rubert-base-cased`[4] with default settings. To ensure reproducibility, the random seed is fixed ($random\_seed = 42$). To obtain embeddings, we used USER-bge-m3 model[3]. To generate questions and answers for Coverage and Agreement metrics, we use Qwen3-235B-A22B [13] model.

In the **hierarchical method with node filtering**, cosine similarity threshold is set at $\theta = 0.85$: if for a summary $S_j$ there existed a previous summary $S_i$ with a cosine similarity above this threshold, then $S_j$ is discarded as redundant. This choice of threshold provides a compromise between preserving meaningful information and eliminating duplication, which empirically led to a noticeable reduction in the volume of intermediate representations without significant degradation in quality.

In the **blueprint method** with question clustering, we utilize KMeans clustering algorithm. The number of clusters is chosen using a heuristically derived rule:

$$n_{\text{clusters}} = \max\left(2, \left\lceil \sqrt{N_{\text{questions}}} \right\rceil\right), \tag{9}$$

where $N_{\text{questions}}$ is the total number of questions generated across all chunks before clustering.

Runtime was measured as the average value (in seconds) of the generation time per book for each method across 100 books.

The experiments used the following large language models: RuadaptQwen2.5-7B-Lite-Beta [10], RuadaptQwen3-32BInstruct-v2 [10], DeepSeek V3 [6], Qwen3-235B-A22B [13], tpro [9] and yagpt5lite [12]. We selected these models to cover three comparable size tiers and to ensure strong Russian capability. At the large tier, DeepSeek V3 [6] and Qwen3-235B-A22B [13] serve as high-capacity baselines for long-document reasoning. At the mid tier, RuadaptQwen3-32BInstruct-v2 [10] and tpro [9] represent instruction-tuned models

---

[4] `https://huggingface.co/DeepPavlov/rubert-base-cased`

with robust Russian coverage. At the lightweight tier, RuadaptQwen2.5-7B-Lite-Beta [10] and yagpt5lite [12] provide cost-efficient options suitable for constrained inference. This pairing by parameter scale lets us compare quality–speed trade-offs under similar computational budgets, while the RuAdapt, tpro and yagpt5lite were chosen specifically for their reported performance on Russian text. Availability on our compute infrastructure and reproducibility considerations also guided the final choice.

## 5. Experimental Results

In all tables models are grouped by parameter count: within each size group, the best result is underlined, and the overall best is in bold type. Values are mean $\pm$ SD across test documents. Metrics of automatic book summarization across models and methods are shown in Tab. 2.

**Table 2.** Main evaluation results

| Model | Metrics | Hierarchical | Blueprint | Hierarchical with node filtering | Blueprint with clustering |
|---|---|---|---|---|---|
| DeepSeek V3 | bertscore | 60.0 $\pm$ 3.1 | 58.0 $\pm$ 4.0 | 60.0 $\pm$ 2.9 | 58.4 $\pm$ 3.6 |
| | rouge-l | 13.7 $\pm$ 3.9 | 12.6 $\pm$ 4.6 | 13.5 $\pm$ 3.7 | 11.2 $\pm$ 3.9 |
| | coverage | **53.57 $\pm$ 21.66** | 40.19 $\pm$ 23.68 | **45.00 $\pm$ 23.03** | **34.68 $\pm$ 23.77** |
| | agreement | 42.38 $\pm$ 17.73 | 32.31 $\pm$ 19.33 | 35.64 $\pm$ 18.88 | 27.76 $\pm$ 19.75 |
| | time | 196.77 $\pm$ 187.85 | 315.67 $\pm$ 321.89 | 147.21 $\pm$ 146.4 | 132.60 $\pm$ 197.25 |
| Qwen3-235B-A22B | bertscore | 61.2 $\pm$ 3.0 | 61.6 $\pm$ 3.3 | 60.9 $\pm$ 2.7 | 59.3 $\pm$ 3.4 |
| | rouge-l | 14.9 $\pm$ 4.0 | 15.8 $\pm$ 4.5 | 14.8 $\pm$ 3.7 | 12.2 $\pm$ 3.6 |
| | coverage | 52.48 $\pm$ 20.79 | **54.78 $\pm$ 21.16** | 44.54 $\pm$ 23.03 | 30.19 $\pm$ 21.96 |
| | agreement | 41.68 $\pm$ 17.18 | 43.99 $\pm$ 17.54 | 35.67 $\pm$ 18.87 | 24.10 $\pm$ 17.62 |
| | time | 103.49 $\pm$ 97.30 | 230.35 $\pm$ 271.03 | 83.06 $\pm$ 102.05 | 158.30 $\pm$ 196.35 |
| RuadaptQwen3-32B Instruct-v2 | bertscore | 57.3 $\pm$ 2.9 | 58.9 $\pm$ 3.6 | 57.7 $\pm$ 3.3 | 55.3 $\pm$ 3.3 |
| | rouge-l | 11.0 $\pm$ 2.4 | 10.6 $\pm$ 3.2 | 10.7 $\pm$ 2.4 | 7.8 $\pm$ 2.1 |
| | coverage | 33.12 $\pm$ 21.50 | 33.18 $\pm$ 22.83 | 32.19 $\pm$ 22.52 | 17.72 $\pm$ 15.23 |
| | agreement | 25.25 $\pm$ 16.94 | 26.21 $\pm$ 18.22 | 24.82 $\pm$ 17.74 | 13.97 $\pm$ 12.39 |
| | time | 218.30 $\pm$ 195.16 | 379.24 $\pm$ 500.40 | 166.79 $\pm$ 164.61 | 286.35 $\pm$ 395.97 |
| tpro | bertscore | 59.4 $\pm$ 3.0 | 59.0 $\pm$ 4.9 | 59.5 $\pm$ 3.3 | 58.2 $\pm$ 3.7 |
| | rouge-l | 13.8 $\pm$ 3.1 | 14.7 $\pm$ 4.9 | 13.5 $\pm$ 3.0 | 11.8 $\pm$ 3.9 |
| | coverage | 40.27 $\pm$ 20.23 | 40.83 $\pm$ 22.42 | 37.13 $\pm$ 20.72 | 26.03 $\pm$ 18.44 |
| | agreement | 31.77 $\pm$ 16.63 | 32.60 $\pm$ 18.57 | 29.44 $\pm$ 16.83 | 20.83 $\pm$ 15.26 |
| | time | 367.32 $\pm$ 324.49 | 592.39 $\pm$ 772.19 | 267.73 $\pm$ 253.34 | 247.59 $\pm$ 361.20 |
| RuadaptQwen2.5-7B Lite-Beta | bertscore | 55.4 $\pm$ 2.9 | 56.1 $\pm$ 4.9 | 55.8 $\pm$ 2.9 | 54.0 $\pm$ 4.0 |
| | rouge-l | 8.6 $\pm$ 2.5 | 10.1 $\pm$ 3.9 | 8.7 $\pm$ 2.5 | 7.7 $\pm$ 2.8 |
| | coverage | 19.66 $\pm$ 17.77 | 24.94 $\pm$ 21.08 | 20.31 $\pm$ 17.95 | 15.51 $\pm$ 14.83 |
| | agreement | 15.16 $\pm$ 14.11 | 20.03 $\pm$ 17.50 | 15.94 $\pm$ 14.39 | 12.23 $\pm$ 12.30 |
| | time | 68.86 $\pm$ 64.85 | 126.84 $\pm$ 145.74 | 53.59 $\pm$ 47.28 | 76.66 $\pm$ 91.78 |
| yagpt5lite | bertscore | 62.5 $\pm$ 3.5 | 61.1 $\pm$ 3.8 | 62.1 $\pm$ 3.2 | 61.5 $\pm$ 3.3 |
| | rouge-l | 16.9 $\pm$ 5.1 | 15.8 $\pm$ 5.1 | 16.4 $\pm$ 4.7 | 14.3 $\pm$ 4.4 |
| | coverage | 36.85 $\pm$ 19.40 | 33.17 $\pm$ 21.58 | 31.75 $\pm$ 20.06 | 24.28 $\pm$ 16.95 |
| | agreement | 29.69 $\pm$ 16.43 | 26.58 $\pm$ 18.13 | 25.60 $\pm$ 16.85 | 19.70 $\pm$ 14.29 |
| | time | 31.02 $\pm$ 28.51 | 113.34 $\pm$ 123.78 | 27.39 $\pm$ 28.05 | 42.15 $\pm$ 56.50 |

In terms of exact matching (ROUGE-L) and semantic replication all models exhibit similar behavior. Low ROUGE-L scores can be explained by high sensitivity to word permutation which

are common in Russian paraphrasing. While BERTScore, also vulnerable to this kind of text perturbations comparing its values to our established similarity threshold (0.85), indicates a major semantic dissimilarity with reference summary. Our question-based metrics (Coverage and Agreement) also confirm frequent summary content deviation. However, these metrics seem to be much more efficient at distinguishing real storytelling errors as they demonstrate a considerably wider value range at the same BERTScore levels.

The best overall performance was achieved by Qwen3-235B-A22B: it delivered the highest coverage and answer agreement. At the same time, the hierarchical method with node filtering offered the best quality-time trade-off. It significantly sped up processing (e.g., almost 2× faster for DeepSeek V3), with comparable quality to the blueprint method which on average achieved the best metrics. The exception was Qwen3-235B-A22B, which achieved its top results with the baseline blueprint. Experiments show that the hierarchical method with node filtering provides the best compromise between speed and quality.

**Table 3.** Comparison of the best and worst english-translated generated summaries

| Title | Text |
|---|---|
| **A Sound of Thunder** | . . .**The main character, Eckels**, a thrill-seeking and overconfident hunter, **pays** a huge sum of money **for the chance to travel 60 million years back in time to kill a Tyrannosaurus rex**. Before the journey, the guide Travis strictly warns him about the rules: **under no circumstances should anyone step off the anti-gravity Path** or interfere with the natural course of events, as even the slightest violation could catastrophically change the future. . . . Travis explains **the fragility of the temporal balance: even the death of a single mouse could wipe out entire species**, and thus alter human history. **The group tracks down a Tyrannosaurus**, marked with red paint — a sign that **its death will not affect the future**. However, at the sight of the giant predator **Eckels panics, steps off the Path, and accidentally crushes a butterfly**. . . Upon returning to 2055 . . .**the world has changed beyond recognition: the language is coarse, the atmosphere oppressive, and instead of the moderate President Keith, a cruel dictator, Deutscher, is in power.** Eckels realizes that **his carelessness triggered the "butterfly effect"** — the crushed insect set off a chain of events that distorted history. In despair, he begs to undo the mistake, but **Travis, understanding the irreversibility of the consequences, raises his rifle.** . . . |
| **Kastrjuk** | . . . The story takes place in a Russian village in early spring, where nature awakens, but people's lives remain harsh and monotonous. **The main character — an old man named Semyon, nicknamed Kastrjuk**, is spending his final days in loneliness, **tormented by memories of his former strength** and regrets over his present frailty. **Once he was known as the best worker in the district,** but now, frail and forgotten, he is forced to stand aside while his fellow villagers work in the fields. . . .**Only in the evening, persuading his son to let him go to the night watch (to herd horses), does Kastrjuk find brief happiness**. Out in the open, among the children and under the starry sky, he feels almost young again. By the pond, a mare drinks water reflecting the sunset, while the **old man, looking at the Milky Way, whispers a prayer — as if he is reconnecting with the world and regaining lost harmony**. But it is only a fleeting solace: tomorrow he will again face hopeless loneliness and the realization of his own uselessness. . . . |

The deviation of question-based metrics can be illustrated by results of hierarchical method obtained by DeepSeek V3. Two summaries were chosen for the analysis: "A Sound of Thunder" and "Kastrjuk". In the first case the model scored high, answering all but one question, but the other summary contained answers to only two questions out of eleven, leading to a low score. These two summaries are shown in Tab. 3. For brevity only the main points that affected the final metric were highlighted. The "Kastrjuk" summary contains many lyrical digressions and stylistic details, making it hard to capture the essence, so the model gets distracted from key facts, whereas in "A Sound of Thunder" events are presented sequentially and clearly, with core plot elements explicitly listed, simplifying retrieval of important information. In the texts, bold type marks plot-relevant fragments, while underlined type indicates content that could be omitted.

**Table 4.** Comparison of models in summary generation using the "Blueprint" method (english-translated)

| Model | Text |
|---|---|
| RuadaptQwen3 | "The company *Time Safari* organizes paid excursions into the past for dinosaur hunting, using time machines capable of moving between eras. Clients are required to follow strict rules: to stay on the metal Path . . . " |
| tpro | "In the text, the main character, Eckels, goes on a time safari in order to kill a Tyrannosaurus rex. The company that organizes the safari guarantees only dinosaurs and strictly forbids hunters from stepping off the Path . . . Mr. Travis, the safari guide, explains that even the destruction of a single mouse could lead to the extinction of all its descendants . . . " |
| DeepSeek V3 | "**Summary by outline:** 1. **Eckels** — the hunter . . . 2. **The company 'Time Safari'** organizes hunting in the past . . . 3. **Travis** — the guide supervising the expedition. . . . " |

Comparing model behavior, DeepSeek V3 generally outperforms smaller models; however, within the blueprint method, in 30% of cases RuadaptQwen3-32B-Instruct-v2 performs best, and tpro in 43%. For reference, consider the summary for "A Sound of Thunder" generated with the blueprint method, with small excerpts shown in Tab. 4. While the DeepSeek V3 summary resembles a numbered list of main events, the outputs from RuadaptQwen3-32B-Instruct-v2 and tpro are cohesive narratives that cover the key plot points.

**Table 5.** Comparison of hierarchical and blueprint methods

| Method | Text |
|---|---|
| **Hierarchical** | . . . Zhulka is a graceful, well-groomed **horse** that lives on the estate . . . |
| **Blueprint** | . . . Zhulka was a small black **dog** with yellow markings . . . |

Note that the best result overall was achieved by the blueprint method with the large model Qwen3-235B-A22B, as shown in Tab. 2. For comparison, on the story "Barbos and Zhulka", the hierarchical method with Qwen3-235B-A22B misclassified "Zhulka" as a horse rather than a dog as shown in Tab. 5. Also, DeepSeek V3 tends to strictly follow the blueprint template and produces a numbered list of key events and main characters, rather than a coherent summary, whereas Qwen3-235B-A22B writes plain text. Thus, the unmodified blueprint method delivered the best results when using the strongest available model – Qwen3-235B-A22B.

**Table 6.** Runtime (seconds) for a text of 81 049 characters (11 chunks)

| Model | Hierarchical | Hierarchical with node filtering | Blueprint | Blueprint with clustering |
|---|---|---|---|---|
| DeepSeek V3 | 237.83 | 72.42 | 292.80 | 268.75 |
| Qwen3-235B-A22B | 113.24 | 39.45 | 215.63 | 145.20 |
| RuadaptQwen3-32BInstruct-v2 | 218.23 | 72.54 | 227.7 | 203.30 |
| tpro | 472.23 | 127.38 | 391.29 | 185.94 |
| RuadaptQwen2.5-7B-Lite-Beta | 84.64 | 25.70 | 103.66 | 78.99 |
| yagpt5lite | 34.17 | 14.08 | 99.70 | 27.26 |

To confirm the efficiency of proposed algorithm modifications and to measure the actual speed up, we conducted an isolated test using text "1408" by Stephen King. The average results of three runs are provided in Tab. 6. Interestingly, larger models such as Qwen3-235B-A22B and DeepSeek V3 showed higher speed than some 32B models achieving almost a 300% speed up on some stories. A key reason is the Mixture-of-Experts (MoE) architecture: during generation only a subset of parameters is active (e.g., $\approx 30B$ out of $\approx 600B$), thus maintaining throughput of smaller models while having substantially higher level of knowledge and task solving skills. Moreover, both RuadaptQwen3-32BInstruct-v2 and tpro generate at least 1.5x more tokens, which noticeably increases the overall runtime.

## Conclusion

In this work we introduced RuBookSum, the first open dataset for Russian long-narrative summarization. To address high computational costs of LLM-based summary generation, we proposed two optimizations to existing approaches: hierarchical with node filtering and blueprint method with clustering. The hierarchical method with node filtering achieves up to 300% speed up with minimal quality loss, making it a perfect choice for long-document summarization under tight context window limits.

Our comparative analysis shows that larger models such as DeepSeek V3 and Qwen3-235B-A22B generally deliver higher key question coverage and factual agreement while having more complete summaries than smaller models, especially with hierarchical and blueprint methods. However, for certain text types and methods (e.g., baseline blueprint), more compact models such as RuadaptQwen3-32B-Instruct-v2 can be a competitive cost-efficient alternative. Qualitative analysis shows that models are better at summarizing linear texts with simple descriptions of events, while books with an abundance of lyrical digressions lead to models omitting key facts. In addition, while blueprint method in conjunction with strong model such as Qwen3-235B-A22B gives the best results, some of generated summaries may turn out to be similar to a enumeration of key events, rather than a coherent text. This implies that future research should consider more advanced quality metrics that would account for stylistic deviations.

## Limitations

Larger-scale study is required to establish statistical significance for the reported metrics and to make defensible claims that one method truly outperforms another. Outcomes can vary with many factors: the choice of models, the composition of the test set (which specific books, genres, and lengths are included in), etc. Accordingly, in this submission we limit ourselves to descriptive

reporting and refrain from significance claims, a full inferential analysis as well as small-scale human analysis is deferred to future work.

We fixed the node-filtering threshold at $\theta = 0.85$ as a pragmatic choice. If $\theta$ were too small, too many fragments would be filtered out, risking loss of useful content, and if too large, most fragments would be kept, defeating the purpose of filtering. A full sensitivity study of $\theta$ is left for future work.

## Acknowledgements

## References

1. Huot, F., Maynez, J., Narayan, S., *et al.*: Text-blueprint: An interactive platform for plan-based conditional generation. In: Croce, D., Soldaini, L. (eds.) Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations. pp. 105–116. Association for Computational Linguistics, Dubrovnik, Croatia (May 2023). `https://doi.org/10.18653/v1/2023.eacl-demo.13`

2. Kim, Y., Chang, Y., Karpinska, M., *et al.*: FABLES: Evaluating faithfulness and content selection in book-length summarization. In: First Conference on Language Modeling (2024)

3. Kryscinski, W., Rajani, N., Agarwal, D., *et al.*: BOOKSUM: A collection of datasets for long-form narrative summarization. In: Goldberg, Y., Kozareva, Z., Zhang, Y. (eds.) Findings of the Association for Computational Linguistics: EMNLP 2022. pp. 6536–6558. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates (Dec 2022). `https://doi.org/10.18653/v1/2022.findings-emnlp.488`

4. LibRusEc: Library of works of art. `https://librusec.org/` (2025), accessed: 2025-07-30

5. 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/`

6. Liu, A., *et al.*: DeepSeek-V3 Technical Report. CoRR (2024), `https://arxiv.org/abs/2412.19437`

7. Narodny Briefly: Digital library of short summaries of literary works. `https://wiki.briefly.ru/` (2025), accessed: 2025-07-30

8. Scirè, A., Conia, S., Ciciliano, S., Navigli, R.: Echoes from Alexandria: A Large Resource for Multilingual Book Summarization. In: Rogers, A., Boyd-Graber, J., Okazaki, N. (eds.) Findings of the Association for Computational Linguistics: ACL 2023. pp. 853–867.

Association for Computational Linguistics, Toronto, Canada (Jul 2023). `https://doi.org/10.18653/v1/2023.findings-acl.54`

9. 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

10. 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`

11. Wu, J., Ouyang, L., Ziegler, D.M., *et al.*: Recursively summarizing books with human feedback (2021), `https://arxiv.org/abs/2109.10862`

12. Yandex: YandexGPT 5 with reasoning mode. `https://ya.ru/ai/gpt` (2025), accessed: 2025-07-30

13. Yang, A., *et al.*: Qwen3 technical report (2025), `https://arxiv.org/abs/2505.09388`

14. Zhang, T., Kishore, V., Wu, F., *et al.*: BERTScore: Evaluating Text Generation with BERT. In: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net (2020), `https://openreview.net/forum?id=SkeHuCVFDr`

# Appendix A. Blueprinting Prompts and Examples

This appendix provides the full English translated prompt templates used in the blueprinting pipeline as well as illustrative English translations of the Q/A pairs.

## A.1. Prompt Templates

**Blueprint: Question Generation**

```
Create a plan for the following text as a list of key questions
that help capture the text's main elements.

Strictly follow the rules:
- Questions must reveal the main events, characters, conflicts,
  and important details.
- Do not create more than 15 questions; choose only the most essential ones.
- Avoid semantically duplicate questions.
- Output questions only, each on a new line, without numbers
  or any additional explanations.

Text:
---
{chunk}
---
```

## Blueprint: Answer Generation

Answer the question using exclusively the information from the provided text.

Strictly follow the rules:
- Be as precise and concise as possible.
- Do not add explanations, commentary, or analysis beyond the text.
- Preserve the original terminology and proper names from the text.
- Use only the information provided in the text.

Text:
---
{chunk}
---

Question:
**{question}**

## Question Generalization

Formulate one generalized key question that:
- Covers the common theme of all questions,
- Preserves their semantic essence,
- Eliminates redundancy and duplication.

Source questions:
---
{questions}
---

Output only the final generalized question without additional comments.

## Summary from Blueprint (with Q/A)

Using the following plan of questions and answers, create a concise
summary of the text presented below. Ensure the summary is logically
coherent and preserves the important elements of the original context.
Do not add anything extraneous.

Plan:
---
{blueprint}
---

Text:
---
{chunk}
---

**Summary from Blueprint (questions only)**

```
Using the following plan of questions, create a concise summary of
the text presented below. Ensure the summary is logically coherent
and preserves the important elements of the original context.
Do not add anything extraneous.


Plan:
---
{blueprint}
---


Text:
---
{chunk}
---
```

## A.2. Illustrative Q/A Examples

The examples below are translated into English and reflect typical outputs produced by the blueprinting stage on two different stories.

**Q1. Who is Mike Enslin and what motivates him?** *A.* Mike Enslin is a writer who investigates alleged paranormal phenomena. He is driven by professional curiosity and the need to gather material for his book, despite his underlying skepticism.

**Q2. What happened to the tower after Maisie was lowered?** *A.* The Wolverden Tower was destroyed: a lightning bolt split it, and part of the structure collapsed. This occurred about an hour after Maisie was lowered.

# Do Open Large Language Models Know What, Where, and When? A Case Study with Quiz-Style Questions

*Anna V. Kuznetsova*[1] (iD), *Viktor A. Byzov*[2] (iD), *Ilias V. Aslanov*[1] (iD),
*Evgeny V. Kotelnikov*[1] (iD)

Large language models (LLMs) are increasingly tested on reasoning-intensive benchmarks, yet their performance on complex quiz-style tasks remains underexplored. In this paper we evaluate modern open-source LLMs on the Russian intellectual game *What? Where? When?*, a challenging format requiring fact recall, associative reasoning, and interpretation of hidden clues. We introduce a new dataset of 2600 questions (2018–2025), enriched with empirical human team success rates and annotated with structural and thematic clusters. We benchmark 14 recent open models accessible via API using both automatic metrics (Exact Match, BLEU, ROUGE) and an LLM-as-a-Judge framework. The best system, `Qwen3-235B-A22B-Thinking`, achieved 32.4% accuracy, but still lagging behind the average human team success rate (45.8%). Large-scale reasoning-enabled models consistently outperformed non-reasoning or smaller counterparts, particularly in domains such as technology, ancient world, psychology, and nature. However, omission, wordplay, and proper-name questions remained difficult across all systems. Comparison with *CheGeKa* (MERA leaderboard) shows that our dataset is substantially harder: while leading proprietary and open models reach EM of 0.534–0.645 and 0.442 on *CheGeKa*, respectively, the strongest model in our benchmark achieves only 0.255 EM. Correlation analysis indicates that human and model perceptions of difficulty only weakly align, suggesting different problem-solving strategies. Qualitative case studies further show that models excel more in fact recall than in reconstructing hidden logic. Our findings highlight both the progress of open LLMs and their current limitations in quiz-style reasoning. The new dataset offers a complementary and more challenging benchmark for Russian-language evaluation.

*Keywords: Large Language Models (LLMs), question answering, reasoning, evaluation metrics, quiz datasets, LLM-as-a-judge, human-AI comparison.*

## Introduction

Recently, there has been a growing interest in evaluating large language models (LLMs) on tasks that require: (1) extensive and well-organized memory; (2) ability to hold multiple information units in working memory while tracking chains of thought; (3) logical-associative reasoning; and (4) fluid intelligence – the capacity to tackle novel problems beyond rote patterns [3, 11, 14, 17, 21]. These competencies are particularly well-tested in intellectual games such as *What? Where? When?* [8]. *What? Where? When?* (Russian: *Chto? Gde? Kogda?*) is a long-running Russian intellectual quiz game similar to *Jeopardy*, where teams of players are given one minute to answer complex, riddle-like questions that often involve hidden clues, wordplay, and multi-step reasoning [2]. Such games offer a challenging testbed to assess LLM reasoning ability.

While many widely used benchmarks (such as *BIG-Bench* [20], *MMLU* [10], and *GSM-8K* [7]) offer multi-step reasoning challenges, they are often too general or already included in model training data. Quiz-style datasets like *TriviaQA* [12], *QANTA (Quizbowl)* [19], and *HotpotQA* [23] provide more direct parallels to trivia competitions, yet still focus mainly on fact recall or retrieval. More recent resources, including *modeLing* (Linguistics Olympiad puzzles) [6], *TurnBench-MS* (multi-turn logic games) [24], and *QUENCH* (open-domain quizzing

---

[1]European University at St. Petersburg, St. Petersburg, Russian Federation
[2]Vyatka State University, Kirov, Russian Federation

with masked rationales) [13], better mirror the cognitive demands of *What? Where? When?* and allow for more robust assessment of reasoning and inference beyond standard static tasks.

However, existing studies on LLM performance in intellectual quiz-style tasks are limited. Prior work has often relied on outdated datasets, such as those likely already included in model pre-training corpora, and on models that no longer represent the state of the art. Moreover, there remains a lack of comparative benchmarks against human teams, and insufficient analysis of how performance varies with question themes and structural features.

In this paper, we address these gaps by presenting a refreshed and rigorous evaluation of modern, open LLMs in the context of *What? Where? When?* game. First, we construct a novel dataset of 2600 questions and answers from *What? Where? When?*, annotated with human team answer-rates. We then perform structural and thematic clustering of the dataset, enabling fine-grained analysis. Employing an "LLM-as-a-Judge" methodology, we evaluate model responses across question clusters for 14 open models accessible via API, identifying both strengths and systematic errors, and complement this analysis with automatic metrics such as Exact Match, BLEU, and ROUGE. Through this analysis, we demonstrate how LLM success correlates with question type and topic, and highlight key reasoning limitations. We also compare model performance to human teams using the recorded success rates and analyze how question difficulty aligns across humans and models. Finally, we present qualitative case studies – successes and failures – that illustrate typical reasoning patterns and types of errors.

Our contributions are as follows:
- We created a new dataset of 2600 *What? Where? When?* questions with human team success rates.
- We applied structural and thematic clustering to enable fine-grained analysis of reasoning requirements.
- We evaluated 14 recent open-access LLMs using LLM-as-a-Judge and automatic metrics, and compared their performance to human teams.
- We analyzed results across clusters, identifying strengths, recurring reasoning failures, and systematic challenges for LLMs.
- We complemented the quantitative results with qualitative case studies that illustrate characteristic successes and errors.

The remainder of this paper is organized as follows. Section 1 reviews previous work on quiz-style question answering and reasoning evaluation of large language models. Section 2 describes the construction and annotation of our What? Where? When? dataset. Section 3 details the methodology of our experiments, including model selection, evaluation metrics, and analysis procedures. Section 4 presents and discusses the results, highlighting model performance across structures, topics, and comparisons with human teams. Finally, the Conclusion summarizes the key findings and outlines directions for future research.

## 1. Previous Work

Research on LLMs in quiz-style question answering spans two main directions. On the one hand, traditional trivia corpora such as *TriviaQA* [12], *QANTA (Quizbowl)* [19], and *HotpotQA* [23] provide large-scale collections of factoid or multi-hop questions. These datasets are valuable for measuring knowledge coverage and retrieval ability, but they only partially reflect the associative reasoning and hidden-clue structure characteristic of intellectual games. On the other hand, a new generation of reasoning-oriented benchmarks (including *modeLing* [6],

*TurnBench-MS* [24], and *QUENCH* [13]) explicitly targets logical inference, multi-turn reasoning, and puzzle-like inference. These resources move closer to the spirit of *What? Where? When?*, though they remain English-centric and do not capture its cultural specificity.

In the Russian-language context, the most significant contribution to date is the *CheGeKa* dataset [17, 21], which contains nearly 29 375 Jeopardy-style questions annotated by topic and difficulty. *CheGeKa* distinguishes between factoid and reasoning tasks and introduced a scoring system adapted to gameplay. It has become the reference benchmark for Russian quiz QA, with a public leaderboard [1] where human teams still lead (token-wise F1=0.719, EM=0.645). Among proprietary models, `Gemini 1.5 Pro` (F1=0.630, EM=0.534) and `Claude 3.7 Sonnet` (F1=0.630, EM=0.526) perform best, while the strongest open model, `DeepSeek-V3-0324`, trails behind (F1=0.531, EM=0.442).

Other studies have explored model behavior on quiz questions in smaller settings. Lifar et al. [14] tested `LLaMA3-405B` on a 416-question *CheGeKa* sample and showed that multi-agent prompting strategies such as self-consistency and suggesterdiscriminator improved Exact Match by about 8 percentage points over single-agent baselines. Aßenmacher et al. [3] introduced *wwm-german-18k*, a German multiple-choice dataset, and found that accuracy remained high on easy questions but dropped to near-random on the hardest levels. Hu et al. [11] proposed a dynamic benchmark of interactive games (*Akinator*, *Taboo*, *Bluffing*) for testing deductive, abductive, and inductive reasoning, showing that different frontier models excel in different reasoning modes.

Our work extends this literature by introducing a new Russian dataset of 2600 *What? Where? When?* questions (2018–2025), enriched with empirical human success rates – an element absent in prior corpora. Unlike earlier studies, we combine structural and thematic clustering, evaluation of 14 recent open models, and comparison to human teams, complemented by qualitative case studies of successes and failures.

## 2. Dataset

The *What? Where? When?* quiz questions were collected from the IQ Game website[3]. The initial dataset contained 3526 entries, which were then preprocessed: blitz questions, multi-part questions with limited answering time, questions with accompanying materials, and rarely used questions were removed. Specifically, we excluded questions that had been played fewer than 100 times on the platform.

After filtering, the final dataset[4] included 2600 unique questions spanning 2018–2025, with an average length of 29.1 words. An example question is shown in Fig. 1.

Using regular expressions, we identified five structural clusters of questions (Tab. 1). Each question was assigned only to the first matching cluster in a priority hierarchy, where more specific patterns had higher precedence.

We employed a two-step procedure for thematic clustering:
1. generation of a list of topics using `BERTopic` [9], `HDBSCAN` algorithm [4], and embeddings from `FRIDA`[5] model;
2. assignment of questions to the identified topics using `Qwen3-235B-A22B` model.

To avoid bias from structural patterns, these were stripped from the questions before thematic clustering. Embedding dimensionality was reduced from 1536 to 50 using `UMAP` [16]. Op-

---

[3]`https://iqga.me`
[4]`https://github.com/kotelnikov-ev/quiz-dataset`
[5]`https://huggingface.co/ai-forever/FRIDA`; we used "`categorize_topic: `" prefix.

> **Example Question**
>
> **id:** 3453
>
> **Question:** According to the construction plan, the scene depicting *THIS* was supposed to be located high above. To better capture the perspective, Antonio asked to hoist a donkey. Answer in two words: what is *THIS*?
>
> **Answer:** Nativity of Christ
>
> **Accepted answers:** Birth of Jesus; Christmas Nativity
>
> **Commentary:** During the construction of the Sagrada Familia, Gaud asked to hoist a donkey by a winch to the height where the Nativity scene was planned to be placed.
>
> **Answer rate:** 82/188 (44%)
>
> **Season:** 2024–2025

**Figure 1.** Illustrative example of a quiz question

**Table 1.** Structural clusters of questions

| Cluster | Number | Share, % | Examples of questions (answer) |
|---|---|---|---|
| Word substitution (HE, SHE, X, SUCH, DOING THIS, ...) | 1363 | 52.4 | In a story by John Coetzee, a savage brought HIM to life with his breath. Name HIM (fire). |
| Answer format (answer in N words, consecutive letters, etc.) | 477 | 18.3 | A riddle of Turkic nomads: "I sit on a hill, stepping on copper bowls." Name these bowls in one word (stirrups). |
| Omission (missing word/letter, abbreviation) | 128 | 4.9 | A jewelry studio is called Room. Write the two Latin letters that we omitted in the name of this studio (au[room]). |
| Name (proper name required) | 85 | 3.3 | Name the person who, according to Michel Pastoureau, was often depicted with black lips (Judas). |
| Other | 547 | 21.0 | What did Tsvetan Angelov call the spears of the snow army? (icicles). |
| **Total** | **2600** | **100** | |

timal `UMAP` and `HDBSCAN` parameters were selected via the Tree-structured Parzen Estimator implemented in `optuna`[6], targeting silhouette maximization and noise minimization.

This process yielded 30 preliminary topic clusters (silhouette score: 0.331). To improve interpretability, semantically similar clusters were identified and merged with the assistance of the `Claude Sonnet 4` model, which compared the most frequent terms and representative questions for each cluster. The LLM was provided with 50 most frequent terms from each

---

[6]`https://optuna.org`

cluster along with 15 randomly selected questions. As a result, we obtained 16 topic clusters with automatically generated names. Subsequently, we assigned all questions to the identified topics using `Qwen3-235B-A22B` model (Tab. 2).

Although it is impossible to completely exclude the possibility of training data overlap, the overall performance of the evaluated models on our dataset (see Section 4.1) suggests that large-scale memorization of the questions is unlikely. If the dataset had been substantially included in pretraining corpora, accuracy levels would plausibly be much higher.

## 3. Methodology

Our study included the following main stages:
1. Selection of a judge model from several proprietary models.
2. Obtaining answers to questions from several open models accessible via API.
3. Analysis of the answers.

### 3.1. Judge Model Selection

To evaluate the quality of the answers, we employed both automatic evaluation metrics (such as Exact Match, BLEU [18], ROUGE-1 and ROUGE-L [15]) and the LLM-as-a-Judge approach [25]. BLEU was computed as a precision-oriented metric based on n-gram overlap, while ROUGE-1 and ROUGE-L were calculated as F1-scores combining precision and recall. Recent studies indicate that metrics based solely on surface overlap (n-grams, exact spans) are limited in their applicability to more complex QA tasks – for example, Chen et al. (2019) show that F1 and similar metrics may not capture answer quality beyond extraction or simple generation tasks [5]. Similarly, Xian et al. (2025) demonstrate that in long-form question answering the style, length and category of answers can heavily bias traditional automatic metrics, and LLM-based evaluators exhibit significantly higher consistency with human judgments [22]. We therefore adopt the LLM-as-a-Judge approach as a complementary method, while acknowledging its own limitations and the need for further review of semantic-based and embedding-based evaluation metrics.

For the LLM-as-a-Judge, a random sample of 10% of the dataset (260 questions) was selected to compare evaluations from candidate LLM judges against human annotators. For these questions, answers were obtained from five open models: `Gemma-3-27b-it`, `QwQ-32B`, `Phi-4-multimodal`, `Llama-4-Scout-17B-16E`, and `Qwen3-32B` (52 answers per model). Their correctness was independently evaluated by two human annotators as well as by several proprietary LLMs (available via API) considered as candidates for the role of judge. The human annotators first labeled the answers of the open models independently. The initial inter-annotator agreement was high, with only a few discrepancies that were discussed and reconciled to obtain a consensus gold-standard set of labels. We then measured which candidate judge model aligned best with this reconciled human annotation set, using Cohens kappa coefficient (Tab. 3). The complete prompt used to instruct the judge model during automatic evaluation is shown below.

**Table 2.** Topic clusters of questions

| Cluster | Number | Share, % | Examples of questions (answer) |
|---|---|---|---|
| Literature | 443 | 17.0% | In the "Aeneid" it is said that Styx forms THEM. Name THEM in two words (nine circles) |
| History | 337 | 13.0% | Interestingly, potatoes first appeared in China during the rule of... Which dynasty? (Ming) |
| Art | 258 | 9.9% | Who demanded to rename his painting to "Love in the Bin"? (Banksy) |
| Nature | 226 | 8.7% | In the illustration to the first chapter of Dr. Komarovskys parenting guide, a plant is shown. Which one? (cabbage) |
| Science | 222 | 8.5% | "First the sails, then the ships hull." Aristotle used this observation to prove... What? (sphericity of Earth) |
| Etymology | 189 | 7.3% | In Norway THIS is called krøllalfa–curly alpha. Name THIS (@) |
| Cinema | 160 | 6.2% | Maria Scholl writes that the Swiss remain true to tradition and still announce IT in local cinemas. Name IT (intermission) |
| Technology | 152 | 5.8% | In the late 1970s a famous company hired young people to stroll around Tokyo. What were they advertising? (Sony Walkman) |
| Ancient World | 100 | 3.8% | Whom did the Aztecs equate with warriors fallen in battle? (women who died in childbirth) |
| Games | 94 | 3.6% | At an event, a chess player from a TV series is asked to say "queen." Instead of which word? (cheese) |
| Sports | 88 | 3.4% | MMA fighter Diana Avsaragova threw in the towel already during IT. Name IT (weigh-in) |
| Geography | 69 | 2.7% | What name was given to the land where Indians wore moccasins of roughly tanned hides? (Patagonia) |
| Psychology | 50 | 1.9% | A wealthy patient of Sigmund Freud had a phobia because of which he literally... did what? (laundered money) |
| Design | 48 | 1.8% | One ATM in Vienna is stylized as HER. Name HER with a hyphenated word (piggy-bank) |
| Numismatics | 11 | 0.4% | The first German radio listener paid 350 marks for the right to use a receiver. In this question we omitted nine of THEM. Name THEM (zeros) |
| Other | 153 | 5.9% | A Russian tattoo salon is called "Yes and No." Which two rhyming words did we replace? (wants, hurts) |
| **Total** | **2600** | **100** | |

```
You are an expert evaluating answers in an intellectual quiz game.
Your task is to assess a list of answers to questions. Evaluate each
answer independently of the others. For each question, you are given:
- "id" - question identifier,
- "question" - the question text,
- "answer" - the answer to be evaluated,
- "correct_answer" - the correct answer,
- "variations" - acceptable alternative answers that should also be
considered correct.
Return only JSON, without any additional comments:  a list of evaluations,
where each item is a dictionary with the keys:
- "id" - question identifier,
- "is_correct" - a boolean value indicating whether the answer is correct.
```

**Table 3.** Agreement of proprietary models with annotators (Cohen's kappa).
Best value is in **bold**, second-best is <u>underlined</u>

| Judge model | Cohen's $\kappa$ |
|---|---|
| GPT-4.1 | **0.9370** |
| Gemini-2.5-flash | <u>0.9348</u> |
| Claude-sonnet-4 | 0.8984 |
| Gemini-2.0-flash-001 | 0.7444 |
| GPT-4.1-mini | 0.6907 |
| GPT-4o-mini | 0.6429 |
| Claude-3.5-haiku | 0.5448 |

The best results were demonstrated by GPT-4.1 and Gemini-2.5-flash. Since at the time of the study the API cost of GPT-4.1 was several times higher than that of Gemini-2.5-flash, the latter was chosen as the judge model, as the quality difference was negligible.

### 3.2. Answer Generation

We evaluated 14 open-source models, focusing on recently released models accessible via API:

- **DeepSeek family:**
  - DeepSeek-R1-0528: Mixture-of-Experts (MoE) architecture; reasoning-first RL model from the V3 family.
  - DeepSeek-V3-0324: MoE, 671B total / 37B active.
  - DeepSeek-V3.1: MoE, hybrid thinking / non-thinking variant.
- **Qwen family:**
  - Qwen3-235B-A22B-Thinking: MoE, 235B total / 22B active; reasoning-oriented.
  - Qwen3-235B-A22B: MoE, 235B total / 22B active; hybrid (instruction + reasoning).
  - Qwen3-30B-A3B: MoE, 30B total / 3B active; hybrid (instruction + reasoning).
  - Qwen3-32B: dense, 32B; hybrid (instruction + reasoning).
  - QwQ-32B: dense, 32B; reasoning-oriented.
- **Llama family:**

- – `Llama-4-Maverick-17B-128E`: MoE; 128 experts, long-context optimization.
  - – `Llama-4-Scout-17B-16E`: MoE; 16 experts, efficiency-focused.
- **Kimi family:**
  - – `Kimi-K2-Instruct`: MoE; ∼1T total / 32B active; long-context model.
- **GPT-OSS family:**
  - – `GPT-OSS-120B`: MoE; ∼117B total / 5.1B active; reasoning-tuned.
- **GLM family:**
  - – `GLM-4.5-Air`: MoE; ∼106B total / 12B active; hybrid reasoning.
- **Gemma family:**
  - – `Gemma-3-27B-it`: dense; 27B instruction-tuned model for dialogue and QA.

We focused on open-source models because they can, in principle, be reproduced or fine-tuned by the community, unlike proprietary counterparts. At the same time, many of the strongest open models require substantial computational resources to run locally. Since our access to hardware was limited, we relied on those open models that are available via the DeepInfra API[7]. This setup enabled systematic evaluation across families and scales at a fraction of the cost of operating large models in-house.

The same prompt was used for all models:

```
You are participating in an intellectual quiz game.
Please briefly reason about the following question and provide an answer.
Question: {question}.
Output your reasoning and answer in JSON format:
{ "reasoning": "your reasoning here",
"answer": "your answer here" }
```

All models were used in their default inference configuration as provided in the official model documentation. The temperature was set to zero to improve determinism. Each model was allowed up to five attempts per question, not to introduce variability, but to handle cases where a model produced no valid response due to output looping or incorrect formatting. If no answer was generated, the maximum output length was increased by 1000 tokens at each retry (starting from 2000 tokens). Despite these multiple attempts, in some cases models still failed to produce any valid answer; such cases were recorded as unanswered and treated as incorrect in subsequent evaluation.

### 3.3. Answer Analysis

Answer evaluation was carried out using two approaches: (1) automatic metrics (Exact Match, BLEU, ROUGE-1, ROUGE-L), and (2) evaluation by the judge model (`Gemini-2.5-flash`).

Before applying automatic metrics, answers were lemmatized, lowercased, and stripped of punctuation. Both the exact reference answers and acceptable variants provided in the dataset were considered correct. For the LLM-as-a-Judge evaluation, the judge model classified each response as correct or incorrect relative to the reference answers, yielding a binary decision. From these judgments we computed *Accuracy*, defined as the proportion of correctly classified responses.

---

[7] `https://deepinfra.com`

After preprocessing and evaluation with both automatic metrics and the judge model, we analyzed the results along several dimensions: overall model performance, variation across structural and thematic clusters, alignment with human team success rates, and representative qualitative examples.

# 4. Results and Discussion

## 4.1. Overall Model Performance

Table 4 summarizes the performance of 14 open-source LLMs on our benchmark, evaluated with both automatic metrics and an LLM-as-a-Judge approach. The results show a considerable variation across models, reflecting differences in reasoning capability, training paradigms, and model scale.

**Table 4.** Performance of open models

| Model | Accuracy | Reasoning | EM | BLEU | R-1 | R-L |
|---|---|---|---|---|---|---|
| DeepSeek-R1-0528 | 30.00 | ✓ | 0.223 | 0.255 | 0.290 | 0.289 |
| DeepSeek-V3-0324 | 29.00 | | 0.222 | 0.250 | 0.287 | 0.286 |
| DeepSeek-V3.1 | 29.65 | | 0.227 | 0.258 | 0.292 | 0.291 |
| Qwen3-235B-A22B-Thinking | **32.42** | ✓ | **0.255** | **0.290** | **0.320** | **0.319** |
| Qwen3-235B-A22B | 20.31 | | 0.156 | 0.181 | 0.208 | 0.207 |
| Qwen3-30B-A3B | 6.12 | | 0.047 | 0.052 | 0.065 | 0.065 |
| Qwen3-32B | 8.58 | | 0.062 | 0.072 | 0.087 | 0.086 |
| QwQ-32B | 12.62 | ✓ | 0.084 | 0.100 | 0.124 | 0.123 |
| Llama-4-Maverick-17B-128E | 21.77 | | 0.172 | 0.199 | 0.227 | 0.226 |
| Llama-4-Scout-17B-16E | 13.81 | | 0.100 | 0.123 | 0.148 | 0.147 |
| Kimi-K2-Instruct | 19.77 | | 0.136 | 0.158 | 0.189 | 0.188 |
| GPT-OSS-120b | 13.65 | ✓ | 0.095 | 0.107 | 0.128 | 0.127 |
| GLM-4.5-Air | 12.42 | ✓ | 0.091 | 0.108 | 0.129 | 0.129 |
| Gemma-3-27b-it | 12.23 | | 0.085 | 0.103 | 0.126 | 0.125 |

Among all evaluated systems, `Qwen3-235B-A22B-Thinking` achieved the best overall performance, with the highest accuracy (32.42%) and superior results on all automatic metrics. Importantly, this model explicitly incorporates a reasoning mode, which appears to contribute significantly to its advantage over the non-reasoning counterpart `Qwen3-235B-A22B`, which reached only 20.31% accuracy. This contrast highlights the effectiveness of explicit reasoning strategies for complex question answering tasks, where solutions often require multi-step inference rather than surface-level retrieval.

The `DeepSeek` family demonstrated relatively strong performance, with accuracies around 2930%. The reasoning-enabled `DeepSeek-R1-0528` slightly outperformed the non-reasoning `DeepSeek-V3` and `DeepSeek-V3.1` variants in terms of *Accuracy*, again underscoring the importance of reasoning traces. However, the performance gap between the reasoning and non-reasoning `DeepSeek` models was narrower than that observed in the `Qwen3` family, suggesting that other architectural or training factors may also play a role.

By contrast, smaller-scale models such as `Qwen3-30B-A3B`, `Qwen3-32B`, and `QwQ-32B` achieved substantially lower accuracies (613%), with weak scores on EM, BLEU, and ROUGE. `QwQ-32B`,

explicitly positioned as a reasoning-oriented variant, outperformed its standard dense counterpart `Qwen3-32B` (12.62% vs. 8.58%), showing that reasoning specialization can bring relative gains even at moderate scale. However, the absolute performance of both models remained low, far behind the large reasoning-enabled `Qwen3-235B-A22B-Thinking`. This suggests that while reasoning traces improve results, they cannot compensate for limited model size and knowledge coverage.

Other evaluated families, including `Llama-4`, `Kimi-K2`, `GPT-OSS`, `GLM-4.5-Air`, and `Gemma-3`, demonstrated moderate to low performance (1222% accuracy). While some of these models occasionally produced plausible answers, their overall metrics remained below those of the strongest `DeepSeek` and `Qwen` variants. Notably, both `GPT-OSS-120B` ($\sim$117B total / 5.1B active) and `GLM-4.5-Air` ($\sim$106B total / 12B active) are reasoning-enabled Mixture-of-Experts architectures, yet their accuracy (1213%) was far below that of the much larger `Qwen3-235B-A22B-Thinking` (235B total / 22B active, 32.42% accuracy). This contrast highlights that scale and effective integration of reasoning capabilities are critical: smaller MoE models with reasoning signals did not achieve competitive performance.

To illustrate typical failure patterns of lower-performing models, Tab. 5 shows two representative examples where all `DeepSeek` models and `Qwen3-235B-A22B-Thinking` produced correct answers, while some weaker models, for example `Kimi-K2-Instruct` and `Qwen3-32B`, failed. Both questions were among the easiest for human participants (answered correctly by over 96% of them), which highlights the qualitative gap between the higher- and lower-performing models.

**Table 5.** Examples of failure cases of lower-performing models

---

**Question:** The symbiotic relationship attributed to certain birds is merely a legend. In fact, these birds catch flies that appear in meat leftovers rather than DO THIS. What does DO THIS mean?

**Correct answer:** clean crocodiles teeth

**Kimi-K2-Instruct / Qwen3-32B:** remove parasites from an animal

---

**Question:** When composer John Tesh came up with a good melody, he was in a hotel and could not write the music down. To preserve the melody, he called his home number of... whom?

**Correct answer:** himself

**Kimi-K2-Instruct / Qwen3-32B:** his wife

---

## 4.2. Performance by Question Structure

Figure 2 presents model accuracies broken down by the main structural categories of questions: word substitution, answer format, omission, name, and other (see Tab. 1). The results reveal systematic differences in difficulty across categories, as well as clear trends in how reasoning-enabled models perform relative to their non-reasoning counterparts.

Word substitution questions (e.g., replacing pronouns or phrases) are the most frequent type and generally yielded the highest accuracies across models (except for the *Other* category). The best results were achieved by `Qwen3-235B-A22B-Thinking` (32.1%) and `DeepSeek-R1-0528` (31.2%), with other `DeepSeek` variants following closely. Even medium models such as `Llama-4-Maverick` (22.2%) and `Kimi-K2-Instruct` (19.1%) achieved moderate success here,
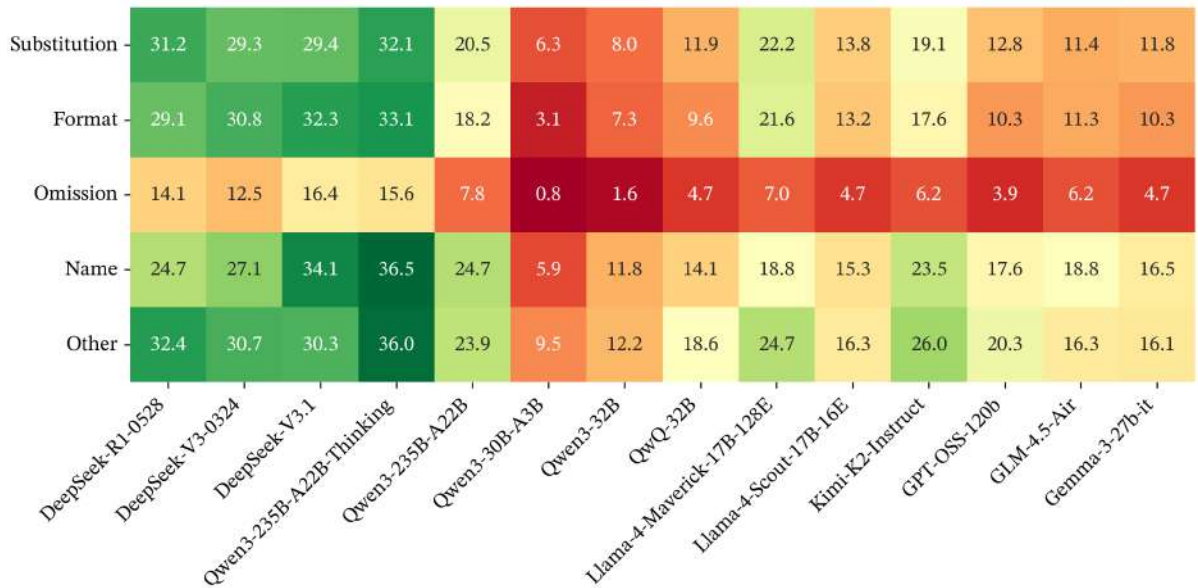
**Figure 2.** Accuracy of LLMs across question structures

indicating that substitution tasks benefit from lexical flexibility and do not always require deep multi-step reasoning.

Answer format questions (e.g., constrained by number of words or letter sequences) proved more challenging, and the strongest results were achieved by `Qwen3-235B-A22B-Thinking` (33.1%) and `DeepSeek-V3.1` (32.3%). The gap between reasoning and non-reasoning models is visible here: `Qwen3-235B-A22B` scored only 18.2%, suggesting that explicit reasoning helps models interpret and enforce output constraints.

Omission questions (requiring restoration of missing words, letters, or abbreviations) were the most difficult across all the models. Even the strongest models did not exceed 17% accuracy. This category appears particularly challenging because it requires precise contextual recall or cultural knowledge rather than general reasoning ability.

Name questions (requiring specific proper names) posed a considerable challenge. The reasoning-enabled `Qwen3-235B-A22B-Thinking` achieved the best result (36.5%), clearly outperforming both its non-reasoning counterpart `Qwen3-235B-A22B` (24.7%) and large non-reasoning models such as `DeepSeek-V3.1` (34.1%). This contrast highlights that while scale alone can bring solid performance, explicit reasoning traces provide an additional advantage for tasks demanding precise entity recall. Mid- and small-scale models struggled markedly (e.g., `Qwen3-30B-A3B` at 5.9%, `Qwen3-32B` at 11.8%), confirming that both scale and reasoning capabilities are crucial for reliably handling proper-name questions.

Finally, the Other category, which aggregates more heterogeneous question types, confirmed the general advantage of reasoning-enabled large models. `Qwen3-235B-A22B-Thinking` reached 36.0%, followed by `DeepSeek-R1-0528` at 32.4%, whereas most mid-scale models remained below 26%.

## 4.3. Performance by Question Topic

Figure 3 reports accuracies by thematic category. The analysis highlights both the relative difficulty of different knowledge areas and the advantage of reasoning-enabled models across most topics.



| | DeepSeek-R1-0528 | DeepSeek-V3-0324 | DeepSeek-V3.1 | Qwen3-235B-A22B-Thinking | Qwen3-235B-A22B | Qwen3-30B-A3B | Qwen3-32B | QwQ-32B | Llama-4-Maverick-17B-128E | Llama-4-Scout-17B-16E | Kimi-K2-Instruct | GPT-OSS-120b | GLM-4.5-Air | Gemma-3-27b-it |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Literature | 20.8 | 24.8 | 23.0 | 23.9 | 13.8 | 2.9 | 4.7 | 6.3 | 13.1 | 8.6 | 12.9 | 9.3 | 7.2 | 8.8 |
| History | 35.0 | 32.6 | 37.1 | 33.5 | 22.3 | 7.1 | 8.0 | 15.1 | 24.3 | 14.5 | 23.7 | 14.2 | 14.2 | 14.8 |
| Art | 27.1 | 29.5 | 26.4 | 29.1 | 19.4 | 5.4 | 7.0 | 10.5 | 21.3 | 10.5 | 17.8 | 11.6 | 10.5 | 11.2 |
| Nature | 38.5 | 32.7 | 31.9 | 41.2 | 27.9 | 6.6 | 8.4 | 12.8 | 24.8 | 15.0 | 18.6 | 12.8 | 14.2 | 12.4 |
| Science | 32.0 | 32.4 | 32.9 | 38.7 | 26.1 | 10.4 | 13.5 | 18.0 | 25.7 | 16.7 | 24.8 | 17.6 | 17.1 | 13.1 |
| Etymology | 25.4 | 23.8 | 23.8 | 24.3 | 15.3 | 5.3 | 6.3 | 9.5 | 18.5 | 13.2 | 16.9 | 16.4 | 9.5 | 8.5 |
| Cinema | 26.9 | 25.6 | 26.9 | 28.1 | 15.6 | 5.0 | 8.8 | 14.4 | 16.9 | 11.2 | 24.4 | 15.0 | 7.5 | 12.5 |
| Technology | 43.4 | 38.8 | 41.4 | 48.7 | 28.3 | 13.2 | 17.8 | 22.4 | 36.2 | 25.7 | 31.6 | 25.7 | 25.7 | 17.8 |
| Ancient World | 40.0 | 39.0 | 46.0 | 46.0 | 32.0 | 7.0 | 17.0 | 22.0 | 37.0 | 20.0 | 21.0 | 20.0 | 21.0 | 14.0 |
| Games | 23.4 | 14.9 | 20.2 | 26.6 | 17.0 | 5.3 | 4.3 | 8.5 | 10.6 | 9.6 | 19.1 | 11.7 | 8.5 | 9.6 |
| Sports | 26.1 | 23.9 | 23.9 | 27.3 | 14.8 | 5.7 | 4.5 | 12.5 | 17.0 | 13.6 | 11.4 | 5.7 | 8.0 | 12.5 |
| Geography | 29.0 | 30.4 | 29.0 | 36.2 | 21.7 | 5.8 | 13.0 | 14.5 | 26.1 | 21.7 | 23.2 | 15.9 | 18.8 | 21.7 |
| Psychology | 40.0 | 40.0 | 40.0 | 42.0 | 26.0 | 8.0 | 18.0 | 18.0 | 38.0 | 22.0 | 30.0 | 16.0 | 16.0 | 20.0 |
| Design | 25.0 | 18.8 | 20.8 | 31.2 | 10.4 | 2.1 | 0.0 | 2.1 | 12.5 | 8.3 | 12.5 | 8.3 | 2.1 | 4.2 |
| Numismatics | 45.5 | 36.4 | 63.6 | 36.4 | 45.5 | 9.1 | 9.1 | 36.4 | 18.2 | 27.3 | 45.5 | 18.2 | 27.3 | 27.3 |
| Other | 28.1 | 25.5 | 24.2 | 29.4 | 16.3 | 3.3 | 7.2 | 8.5 | 22.2 | 11.8 | 15.7 | 8.5 | 10.5 | 10.5 |

**Figure 3.** Accuracy of LLMs across question topics

Questions from literature and art proved moderately challenging, with accuracies not exceeding 30%. The best results in these categories were achieved by `DeepSeek-V3-0324` (24.8% in literature; 29.5% in art), followed by `Qwen3-235B-A22B-Thinking` (23.9% / 29.1%), while smaller dense models often remained below 15%. In history and the ancient world, reasoning-capable and large-scale models performed much better: `DeepSeek-V3.1` and `Qwen3-235B-A22B-Thinking` reached 33–46%, and even mid-size models such as `Llama-4-Maverick` achieved moderate accuracy (2437%), suggesting that historical knowledge is relatively well represented in training corpora.

Nature and science questions achieved relatively high accuracies. `Qwen3-235B-A22B-Thinking` scored 41.2% in nature and 38.7% in science, while `DeepSeek-R1` also performed strongly (38.5% and 32.0%). Smaller dense models again fell below 20%, indicating that large MoE architectures with reasoning support are especially effective for factual and explanatory domains. In technology, results were even stronger: `Qwen3-235B-A22B-Thinking` reached 48.7%, the best score across all categories, with `DeepSeek-R1` and `DeepSeek-V3.1`

exceeding 40%. This likely reflects both rich representation of contemporary technological concepts in pretraining data and the reasoning-friendly structure of such questions.

Performance in specialized domains varied widely. In numismatics, `DeepSeek-V3.1` achieved 63.6% accuracy (the single highest category-level result), but this figure is based on only 11 questions, so it should be interpreted with caution. Design questions proved difficult for all models, with a maximum of 31.2% by `Qwen3-235B-A22B-Thinking`. Etymology also challenged most systems, with top results below 26%.

In more culturally grounded or popular categories such as cinema, games, and sports, even the strongest models rarely exceeded 27%. Here, reasoning-enabled models (`Qwen3-235B-A22B-Thinking` and `DeepSeek-R1`) maintained a relative edge but still lagged behind their performance in science and technology. Geography and psychology showed stronger outcomes: `Qwen3-235B-A22B-Thinking` reached 36.2% and 42.0% respectively, while smaller dense models rarely surpassed 20%.

Finally, in the heterogeneous "Other" category, large reasoning-enabled models again outperformed their non-reasoning counterparts (`Qwen3-235B-A22B-Thinking` at 29.4% vs. 16.3% for `Qwen3-235B-A22B`), while mid-scale models typically stayed around 1522%.

Overall, the thematic breakdown confirms that reasoning-enabled large-scale MoE models consistently lead across domains, but their relative advantage varies depending on the knowledge area, with particularly strong gains in technology, ancient world, psychology, and nature.

## 4.4. Comparison with Human Teams

Since the dataset includes empirical measures of human team performance for each question, we first computed the *average team success rate*, defined as the proportion of correct answers across all teams and all questions. This quantity reflects the expected probability that a randomly selected team would answer a randomly selected question correctly.

A direct comparison of this measure with model accuracy, however, is not fully appropriate. The *average human success rate* aggregates performance across a population of teams and captures the distribution of abilities in the sample, whereas *model accuracy* describes the performance of a single agent answering each question once. Thus, while both metrics are probabilities of success on a random question, they represent different types of averages: one collective, the other individual. For this reason, numerical values cannot be interpreted as strictly equivalent.

Nevertheless, with these limitations in mind, the overall level of performance of the models can be interpreted against the human benchmark. In our case, the mean success rate of human teams was **45.8%**, whereas the best-performing model reached **32.4%** accuracy. This indicates that the model underperforms relative to the average human team, although the comparison should be interpreted with caution.

To assess whether models and humans perceive question difficulty in a comparable way, we examined the relationship between human success rates per question and model outcomes. We computed *Pearsons correlation* (sensitive to linear relationships) and *Spearmans rank correlation* (robust to monotonic but non-linear dependencies). Both coefficients converged to the same result: a weak but statistically significant positive correlation ($r \approx 0.19$, $p < 10^{-22}$). This indicates that questions which are easier for humans tend to be somewhat easier for the model as well, though the strength of the relationship is limited.

We further stratified questions into ten groups according to human success rate (from 0–10% up to 90–100%) and plotted model accuracy in each bin (Fig. 4). The barplot shows a clear

upward trend: model accuracy rises from about 13% on the hardest questions (0–10% human success) to about 62% on the easiest questions (90–100% human success). However, the model consistently lags behind human teams, most strikingly on easy questions, where human success approaches 100% while the model remains far below this ceiling.
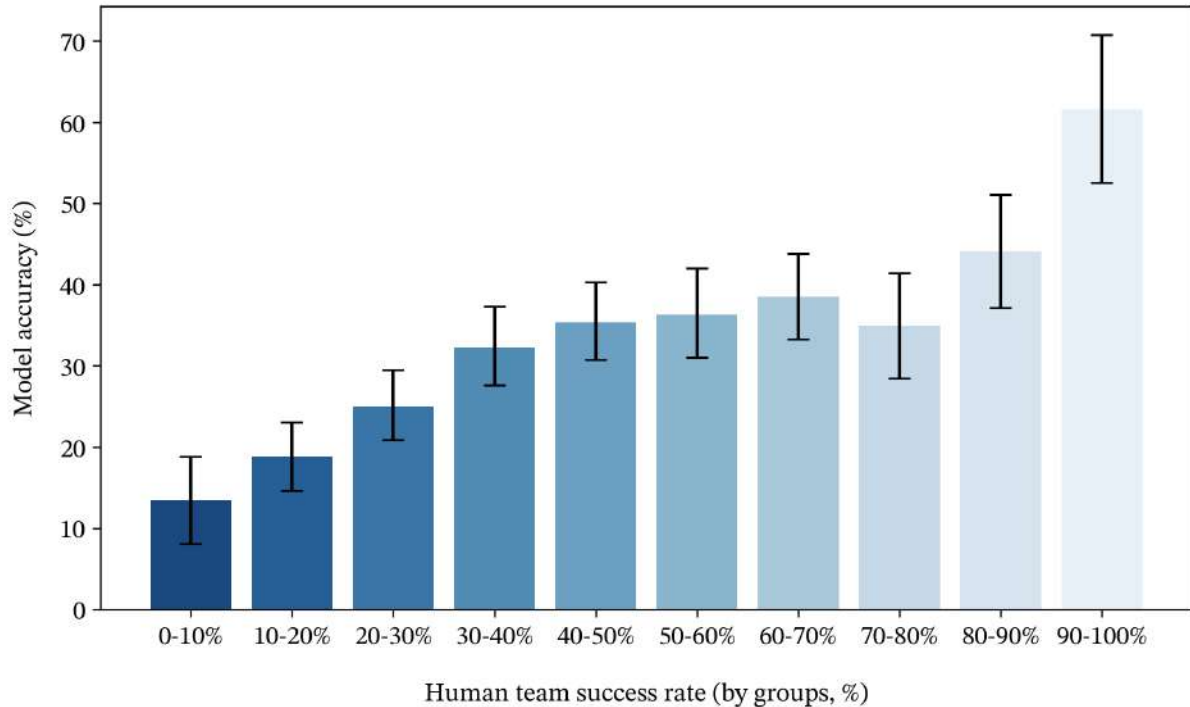


**Figure 4.** Accuracy of the model on questions grouped by human team success rate (%). Darker colors indicate more difficult questions, lighter colors indicate easier ones. Error bars show 95% bootstrap confidence intervals

In summary, models exhibit systematically lower average performance compared to human teams and only partially align with human judgments of question difficulty. While there is evidence of a shared gradient of difficulty (harder questions for humans also tend to be harder for models), the relatively low correlations and persistent performance gap indicate substantial differences in underlying problem-solving strategies.

## 4.5. Illustrative Case Studies

To complement the aggregate statistics, we present qualitative examples that illustrate where human teams and models converge or diverge. We organize the cases into four categories (Tab. 6):
- difficult for humans, answered correctly by the model;
- easy for humans, answered incorrectly by the model;
- difficult for both humans and the model;
- easy for both.

These examples suggest that models tend to perform better on questions that hinge on recalling well-known facts or stable cultural associations, while they are less reliable when success depends on reconstructing hidden logic, metaphors, or wordplay.

**Table 6.** Examples of questions illustrating differences between human and model performance. Answers produced by `Qwen3-235B-A22B-Thinking` model are shown. Correct model answers are highlighted in green, incorrect ones in red

| Question | Human success | Model answer | Correct answer |
|---|---|---|---|
| *Difficult for humans, answered correctly by the model* | | | |
| "The achievement a player receives for the first queen checkmate on chess.com has the same name as a 1974 work. Name this work." | 0.0% | "Killer Queen" | Killer Queen |
| "The slogan Let your fingers do the walking belonged to THEM. In a note about THEM, Dennis Yu mentioned several tens of millions of trees. Give THEIR name or say what THEY are." | 2.5% | "Yellow Pages" | Yellow Pages |
| *Easy for humans, answered incorrectly by the model* | | | |
| "Recently, a prisoner serving a life sentence demanded release after HER. Name HER in two words." | 96.8% | "Constitutional Court" | clinical death |
| "On the tram line between Dsseldorf and Duisburg, at certain times of the day, THEY used to attach HIM to the trains. Name HIM." | 95.1% | "trailer car" | dining car |
| *Difficult for both humans and the model* | | | |
| "Above the Goshin grove rise several bare tree trunks, as if struck by lightning, though that could not have happened. Goshin is one of the most famous examples of group of THESE. Name THIS with an indeclinable word." | 0.0% | ¡No answer¿ | bonsai |
| "Folk etymology traces THIS WORD to a numeral. Researchers, however, believe it goes back to a Proto-Indo-European root meaning to cause pain. Name THIS WORD." | 0.5% | "trizna" (Old Russian funeral feast) | duel |
| *Easy for both humans and the model* | | | |
| "The national motto of Switzerland emerged in the second half of the 19th century. The motto emphasizes that, despite the division of the country into many cantons, if something happens, each canton will support the others. Name this motto." | 100.0% | "Unus pro omnibus, omnes pro uno" | One for all and all for one |
| "The airtight armor made for Henry VIII in the 16th century allowed him to move quite freely. Dan Snow writes that several centuries later, Henry VIIIs armor attracted the interest of specialists working on a commission for which organization?" | 99.0% | "NASA" | NASA |

## Conclusion

In this study, we introduced a new dataset of 2600 *What? Where? When?* questions collected from 2018–2025 and enriched with empirical human success rates. Using structural and thematic clustering, we provided a fine-grained view of question types and knowledge domains, and evaluated 14 recent open-source LLMs with both automatic metrics and an LLM-as-a-Judge approach.

Our results show that the strongest open models, such as `Qwen3-235B-A22B-Thinking` and `DeepSeek-R1`, approach but do not surpass the average human team performance. Large-scale reasoning-enabled architectures consistently outperformed their non-reasoning counterparts, particularly in domains like technology, ancient world, psychology, and nature, while smaller dense models lagged behind across categories. At the same time, omission and wordplay-based questions remained difficult for all systems, underscoring persistent weaknesses in handling associative reasoning and linguistic creativity.

The inclusion of human answer rates allowed us to directly compare model accuracy with human performance. Although correlations between model and human difficulty patterns were statistically significant ($r \approx 0.19$, $p < 10^{-22}$), they were weak, suggesting that humans and models rely on different problem-solving strategies. Qualitative examples further confirmed that models excel more often at fact recall than at reconstructing hidden logic.

Our *What? Where? When?* benchmark is substantially harder than prior Russian quiz datasets. Under the same EM metric, the best result on our data is EM = 0.255 for `Qwen3-235B-A22B-Thinking` (EM = 0.222 for `DeepSeek-V3-0324`), whereas on *CheGeKa* (MERA) `DeepSeek-V3-0324` reaches EM = 0.442; proprietary `Gemini 1.5 Pro` and `Claude 3.7 Sonnet` achieve EM = 0.534 and 0.526, and the human benchmark stands at EM = 0.645 [1]. For metric alignment, we compare EM to EM (MERA reports token-wise F1 and EM), and we avoid contrasting F1 with our judge-based Accuracy. For context, the strongest models judge-based Accuracy on our benchmark is 32.4%.

These findings highlight both the progress of modern open LLMs and their current limitations in intellectual quiz-style reasoning. Future work may expand the dataset, explore interactive multi-agent approaches, and integrate richer evaluation of reasoning traces, bringing automated systems closer to the cognitive style of human quiz players.

## References

1. MERA Leaderboard. `https://mera.a-ai.ru/en/text/leaderboard`, accessed: 2025-09-08

2. What? Where? When? `https://en.wikipedia.org/wiki/What%3F_Where%3F_When%3F`, accessed: 2025-09-08

3. Aßenmacher, M., Karrlein, L., Schiele, P., *et al.*: Introducing wwm-german-18k - can LLMs crack the million? (or win at least 500 euros?). In: Proceedings of the 7th International Conference on Natural Language and Speech Processing (ICNLSP 2024). pp. 287–296 (2024), `https://aclanthology.org/2024.icnlsp-1.31/`

4. Campello, R.J.G.B., Moulavi, D., Sander, J.: Density-based clustering based on hierarchical density estimates. In: Advances in Knowledge Discovery and Data Mining. Lecture Notes

in Computer Science, vol. 7819, pp. 160–172. Springer (2013). `https://doi.org/10.1007/978-3-642-37456-2_14`

5. Chen, A., Stanovsky, G., Singh, S., *et al.*: Evaluating question answering evaluation. In: Proceedings of the 2nd Workshop on Machine Reading for Question Answering. pp. 119–124 (2019). `https://doi.org/10.18653/v1/D19-5817`

6. Chi, N., Malchev, T., Kong, R., *et al.*: ModeLing: A Novel Dataset for Testing Linguistic Reasoning in Language Models. In: Proceedings of the 6th Workshop on Research in Computational Linguistic Typology and Multilingual NLP. pp. 113–119 (2024), `https://aclanthology.org/2024.sigtyp-1.14/`

7. Cobbe, K., Kosaraju, V., Bavarian, M., *et al.*: Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168 (2021), `https://arxiv.org/abs/2110.14168`

8. Foster, E.J., Friedlander, K.J., Fine, P.A.: Mastermind and expert mind: A qualitative study of elite quizzers. Journal of Expertise 8(1), 38–71 (2025), `https://www.journalofexpertise.org/articles/volume8_issue1/JoE_8_1_Foster_etal.html`

9. Grootendorst, M.: BERTopic: Neural topic modeling with a class-based TF-IDF procedure. arXiv preprint arXiv:2203.05794 (2022), `https://arxiv.org/abs/2203.05794`

10. Hendrycks, D., Burns, C., Basart, S., *et al.*: Measuring massive multitask language understanding. In: Proceedings of the International Conference on Learning Representations (ICLR) (2021), `https://openreview.net/forum?id=d7KBjmI3GmQ`

11. Hu, L., Li, Q., Xie, A., *et al.*: GameArena: Evaluating LLM reasoning through live computer games. In: The Thirteenth International Conference on Learning Representations (ICLR) (2025), `https://openreview.net/forum?id=SeQ8l8xo1r`

12. Joshi, M., Choi, E., Weld, D.S., *et al.*: TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics. pp. 1601–1611 (2017). `https://doi.org/10.18653/v1/P17-1147`

13. Khan, M.A., Yadav, N., Masud, S., *et al.*: QUENCH: Measuring the gap between Indic and non-Indic contextual general reasoning in LLMs. In: Proceedings of the 31st International Conference on Computational Linguistics. pp. 4493–4509 (2025), `https://aclanthology.org/2025.coling-main.303/`

14. Lifar, M., Protsenko, B., Kupriianenko, D., *et al.*: LlaMa meets Cheburashka: impact of cultural background for LLM quiz reasoning. In: Language Gamification - NeurIPS 2024 Workshop (2024), `https://openreview.net/forum?id=xCAzTXumhh`

15. 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/`

16. McInnes, L., Healy, J., Saul, N., *et al.*: Umap: Uniform manifold approximation and projection for dimension reduction. The Journal of Open Source Software 3(29), 861 (2018), `https://joss.theoj.org/papers/10.21105/joss.00861`

17. Mikhalkova, E., Khlyupin, A.A.: Russian Jeopardy! Data Set for Question-Answering Systems. In: Proceedings of the Thirteenth Language Resources and Evaluation Conference. pp. 508–514 (2022), `https://aclanthology.org/2022.lrec-1.53/`

18. Papineni, K., Roukos, S., Ward, T., *et al.*: BLEU: a Method for Automatic Evaluation of Machine Translation. In: Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics. pp. 311–318 (2002). `https://doi.org/10.3115/1073083.1073135`

19. Rodriguez, P., Feng, S., Iyyer, M., *et al.*: Quizbowl: The case for incremental question answering (2021), `https://arxiv.org/abs/1904.04792`

20. Srivastava, A., Rastogi, A., Rao, A., *et al.*: Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. Transactions on Machine Learning Research (2023), `https://openreview.net/forum?id=uyTL5Bvosj`

21. Taktasheva, E., Shavrina, T., Fenogenova, A., *et al.*: TAPE: Assessing few-shot Russian language understanding. In: Findings of the Association for Computational Linguistics: EMNLP 2022. pp. 2472–2497 (2022). `https://doi.org/10.18653/v1/2022.findings-emnlp.183`

22. Xian, N., Fan, Y., Zhang, R., *et al.*: An empirical study of evaluating long-form question answering. In: Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 1141–1151. SIGIR '25 (2025). `https://doi.org/10.1145/3726302.3729895`

23. Yang, Z., Qi, P., Zhang, S., *et al.*: HotpotQA: A dataset for diverse, explainable multi-hop question answering. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. pp. 2369–2380 (2018). `https://doi.org/10.18653/v1/D18-1259`

24. Zhang, Y., Wang, M., Li, X., *et al.*: TurnBench-MS: A benchmark for evaluating multi-turn, multi-step reasoning in large language models. In: Christodoulopoulos, C., Chakraborty, T., Rose, C., Peng, V. (eds.) Findings of the Association for Computational Linguistics: EMNLP 2025. pp. 19892–19924. Association for Computational Linguistics, Suzhou, China (Nov 2025). `https://doi.org/10.18653/v1/2025.findings-emnlp.1084`

25. Zheng, L., Chiang, W.L., Sheng, Y., *et al.*: Judging LLM-as-a-judge with MT-bench and Chatbot Arena. In: Proceedings of the 37th International Conference on Neural Information Processing Systems. NIPS '23, Curran Associates Inc. (2023). `https://doi.org/10.5555/3666122.3668142`

# One-Shot Prompting for Russian Dependency Parsing

*Elena D. Shamaeva*[1] iD *, Mikhail M. Tikhomirov*[1] iD *,*
*Natalia V. Loukachevitch*[1] iD

This study investigates the application of Large Language Models (LLMs) for dependency parsing of Russian sentences. We evaluated several models (including Qwen, RuAdapt, YandexGPT, T-pro, T-lite, and Llama) in a one-shot mode across multiple Russian treebanks: SynTagRus, GSD, PUD, Poetry, and Taiga. Among the models tested, Llama70 achieved the highest scores in both UAS and LAS. Furthermore, we observed a general trend where larger models tended to perform better. Our analysis also revealed that parsing quality for Qwen4 and RuAdapt4 on the Taiga treebank was notably sensitive to prompt design. However, the results from all LLMs remained lower than those obtained from classical neural parsers. A key challenge encountered by many models was a difference between generated token sets and gold token sets, which was observed in a considerable portion of each treebank. Additionally, the T-pro and T-lite models produced a significant number of extra lines. The implementation for this study is publicly available at `https://github.com/Derinhelm/llm_parsing/tree/main`.

*Keywords: LLM, parsing, dependency tree, one-shot, prompt-tuning, Russian language.*

## Introduction

In the era of Large Language Models (LLMs), syntax parsing remains an important task in Natural Language Processing (NLP), because syntax parsers allow for obtaining more interpretable results. These parsers are used as an auxiliary tool in tasks such as assessing text complexity [11], paraphrasing [10], named entity recognition [1, 9, 12, 18], and plagiarism detection [15]. Additionally, parsers are used for linguistic text analysis [3].

For syntax parsing, both classical neural networks and LLMs can be applied. While classical neural parsers have achieved a high level of performance[2], the application of LLMs to this task represents an emerging field of research. In this field, the widely adopted technique is prompt-based tuning of LLMs [5]. An important aspect of this prompt-based approach is the optimal design of prompts, particularly the selection of prompt examples. This research direction, however, remains notably underexplored for syntax parsing of the Russian language.

The article describes the experiment on applying LLMs in one-shot mode for syntax parsing of Russian sentences. The evaluation was conducted on five test samples from Russian corpora, which contain sentences with syntax annotation. This research evaluates models such as Qwen, RuAdapt, Llama, T-pro, T-lite, and YandexGPT. A significant feature of this study is the specification of the gold token set in the prompt.

The article is organized as follows. Section 1 discusses related works. Section 2 provides information about syntax parsing. Section 3 outlines the experimental setup. Section 4 describes the results. The Conclusion summarizes the study and suggests directions for future work.

---

[1]Lomonosov Moscow State University, Moscow, Russian Federation
[2]For the Russian language [14], parsers DeepPavlov and Stanza exceed 0.9 by the metric UAS (Unlabeled Attachment Score) on the PUD and SynTagRus treebanks, also, the parsers exceed 0.75 on the treebanks Taiga, Poetry and GSD.

# 1. Related Works

In the area of syntax parsing, an emerging field is the application of LLMs with an autoregressive decoder architecture. The main approaches for this purpose are both fine-tuning [5] and prompt-tuning [6] methods. Furthermore, these parsers are distinguished by the employed syntactic representation: dependency trees or constituency structures. These key distinctions are summarized in Tab. 1.

**Table 1.** Related works

| Article | Research type | Russian language | Syntax structure | Prompting mode |
|---|---|---|---|---|
| [6] | Fine-tuning | yes | Dependency | — |
| [19] | Fine-tuning | no | Dependency | — |
| **[2]** | Both | no | Constituency | Zero-shot, few-shot |
| **[5]** | Prompt-tuning | no | Dependency | **One-shot** |
| **[16]** | Prompt-tuning | no | Constituency | Zero-shot, others* |
| [7] | Other | no | Constituency | — |
| **This article** | **Prompt-tuning** | **yes** | **Dependency** | One-shot |

\* Five-shot prompt-tuning, and zero-shot prompt-tuning in Chain-of-Thought mode.

The articles about prompt-tuning do not consider the Russian language. Moreover, among them, only article [5] is devoted to prompt-tuning for dependency trees[3], and therefore its prompt design is the basis for our work. However, since that article did not consider Russian sentences, the results of these studies are not directly comparable.

# 2. Syntax

## 2.1. Syntax Representation

The most popular ways to represent the syntax structure of a sentence are constituency structure and dependency structure (also called dependency tree). Figure 1 shows examples of the structures. For the Russian language most syntax datasets are represented by dependency trees. A dependency tree is a directed graph, the nodes of which correspond to sentence tokens (words, punctuation marks, etc.) and edges correspond to relations between tokens. Each token is connected to one main token, which is called parent token. The root token of the tree is connected to the auxiliary token ROOT.
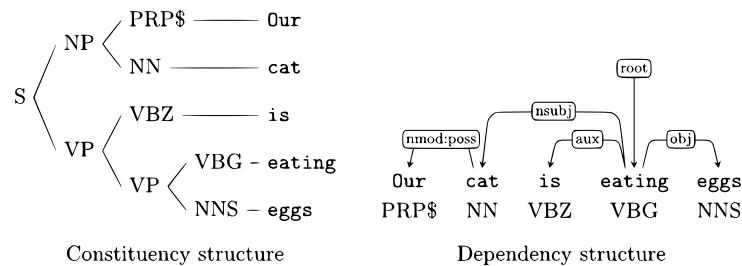


**Figure 1.** An example of constituency and dependency structures [8]

---

[3]An example of a prompt from [5] is provided in Appendix A.

The parsing by LLM is the task of generating a text sequence which describes the syntax structure of a sentence. So, a dependency tree should be represented in text format. There are two widely used formats: the CoNLL-U format and the bracket sequence format. The CoNLL-U format is also used in datasets of sentences with syntax markup (also called treebanks). In the CoNLL-U format, each line (except comment) corresponds to a token and consists of ten values, splitted by the tab character. The first value is the token identifier, the second is the token text, the seventh is the identifier of the parent token, the eighth is the relation tag. While the remaining values are morphological and other features of the token. Figure 2 shows an example of a CoNLL-U sentence.
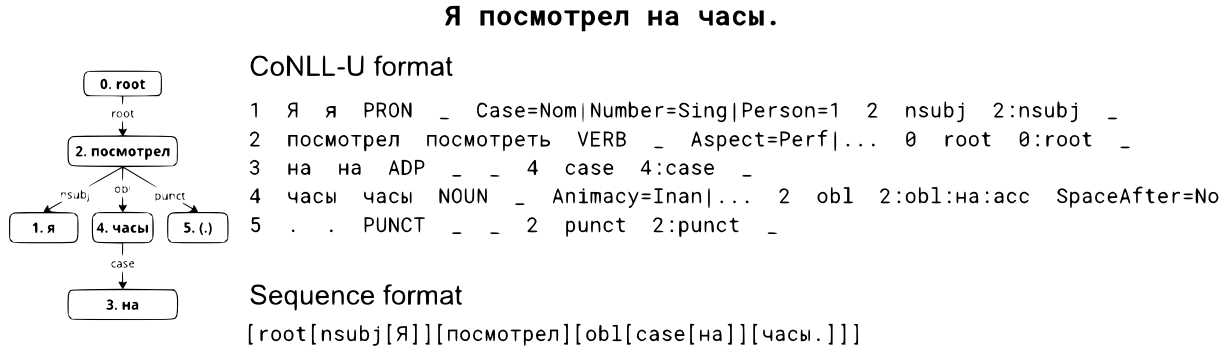


**Я посмотрел на часы.**

CoNLL-U format

```
1  Я  я  PRON  _  Case=Nom|Number=Sing|Person=1  2  nsubj  2:nsubj  _
2  посмотрел  посмотреть  VERB  _  Aspect=Perf|...  0  root  0:root  _
3  на  на  ADP  _  _  4  case  4:case  _
4  часы  часы  NOUN  _  Animacy=Inan|...  2  obl  2:obl:на:acc  SpaceAfter=No
5  .  .  PUNCT  _  _  2  punct  2:punct  _
```

Sequence format

```
[root[nsubj[Я]]][посмотрел][obl[case[на]][часы.]]]
```

**Figure 2.** Examples of a sentence in the CoNLL-U and sequence formats

Datasets of dependency trees (treebanks) are stored primarily in the CoNLL-U format. So, the CoNLL-U format can be used by LLMs without fine-tuning [5], while fine-tuning is required for generating a dependency tree in the bracket sequence format [6]. Often, only four syntax columns of the CoNLL-U format are generated: token ID, form, parent ID, and relation type, while columns lexeme, part of speech and morphological features are replaced with the underscore symbol.

## 2.2. Evaluation of Syntax Parser

The standard way to evaluate parser results is to calculate metrics UAS (Unlabeled Attachment Score) and LAS (Labeled Attachment Score). The preliminary stage before evaluation is the aligning. At this stage, a correspondence is established between tokens from the dataset and tokens from the dependency tree, created by parser. After that, F1-score of UAS and LAS is calculated. Equations 1-6 show formulas for the calculation. G is a set of gold token, P is a set of dependency tree tokens, while $\mathbf{p(t)}$ is a function that maps a token to the parent token, $\mathbf{d(t)}$ is a function that maps a token to the relation between the token and the parent token.

$$UAS\_precision = \frac{\| \, gt|gt = pt, gt \in G, pt \in P, \mathbf{p}(gt) = \mathbf{p}(pt) \, \|}{\| \, pt|pt \in P \, \|}, \tag{1}$$

$$UAS\_recall = \frac{\| \, gt|gt = pt, gt \in G, pt \in P, \mathbf{p}(gt) = \mathbf{p}(pt) \, \|}{\| \, gt|gt \in G \, \|}, \tag{2}$$

$$LAS\_precision = \frac{\| \, gt|gt = pt, gt \in G, pt \in P, \mathbf{p}(gt) = \mathbf{p}(pt), \mathbf{d}(gt) = \mathbf{d}(pt) \, \|}{\| \, pt|pt \in P \, \|}, \tag{3}$$

$$LAS\_recall = \frac{\| gt|gt = pt, gt \in G, pt \in P, \mathbf{p}(gt) = \mathbf{p}(pt), \mathbf{d}(gt) = \mathbf{d}(pt) \|}{\| gt|gt \in G \|}, \qquad (4)$$

$$UAS\_F1 = \frac{2 * UAS\_precision * UAS\_recall}{UAS\_precision + UAS\_recall}, \qquad (5)$$

$$LAS\_F1 = \frac{2 * LAS\_precision * LAS\_recall}{LAS\_precision + LAS\_recall} \qquad (6)$$

## 3. Experimental Setup

### 3.1. Prompts

In our article, as in [5], the prompt contains one sentence with the gold dependency tree in the CoNLL-U format. In both articles, the length of gold sentence is from 4 to 7. The sentence is used more as an example of the format than as a source of linguistic information. In [5] test sentences are randomly selected from the train set of a treebank, while in our work each prompt contains one gold dependency tree. Another difference from [5] lies in the format of token set representation. While in the original article, tokens from test sentences are simply separated by spaces, in our article the token set is represented as a python list. This approach is based on the assumption that LLMs work well with program code. In future, comparison between the token set representations will be planned.

Figure 3 shows an example of our prompt. In addition to prompts from [5], the article prompts also include gold token set and some restrictions. In our study, we experiment with 10 prompts. In our article, as in [5], one-shot prompting is considered. Each test prompt is passed to a tested LLM once. Figure 4 shows the gold sentences for the article prompts.

```
Пример: Предложение <Рядом проходит автомобильная дорога .> в формате
CONLL:
1 Рядом _ _ _ _ 2 advmod _ _
2 проходит _ _ _ _ 0 root _ _
3 автомобильная _ _ _ _ 4 amod _ _
4 дорога _ _ _ _ 2 nsubj _ _
5 . _ _ _ _ 2 punct _ _
Задание: Верни в формате CONLL предложение <С 2012 года центр
занимается также вопросом об освещении изменения климата .>:
Результат должен состоять из 12 строк в формате CONLL. Во втором
столбце должны быть токены ['С', '2012', 'года', 'центр', 'занимается',
'также', 'вопросом', 'об', 'освещении', 'изменения', 'климата', '.'].
Нельзя нарушать порядок токенов. Нельзя добавлять токены. Нельзя
удалять токены.
```

**Figure 3.** The prompt example of the article

1. Рядом проходит автомобильная дорога.
2. Соглашение рассчитано на два года.
3. Доцент Саратовского государственного университета.
4. Девушка ждала его 4 года.
5. Началу работ препятствовал недостаток финансирования.
6. Сила трения незначительная.
7. В России встречаются 2 вида.
8. Николай Резанов родился в Ленинграде.
9. Жена: Ольга Александровна Михайлова.
10. Женат, имеет трех сыновей.

**Figure 4.** Gold sentences

## 3.2. Models

In the study, we considered ten LLMs, working with the Russian language. The models are selected from open-sources projects Qwen[4], RuAdapt[5] [17], T-Tech[6], Yandex[7], LLaMA[8]. RuAdapt models are Russian adapted versions of Qwen models. Additionally, LLMs from T-Tech are based on Qwen models.

The study considers LLMs with different amount of parameters. Table 2 shows the values.

**Table 2.** Statistics on the amount of LLM parameters

| Model | Parameters (billions) |
| --- | --- |
| Qwen/Qwen3-32B | 32.8 |
| Qwen/Qwen3-8B | 8.19 |
| Qwen/Qwen3-4B | 4.02 |
| RefalMachine/RuadaptQwen3-32B-Instruct | 32.7 |
| RefalMachine/RuadaptQwen3-8B-Hybrid | 8.14 |
| RefalMachine/RuadaptQwen3-4B-Instruct | 4.01 |
| t-tech/T-pro-it-2.0 | 32.8 |
| t-tech/T-lite-it-1.0 | 7.61 |
| yandex/YandexGPT-5-Lite-8B-instruct | 8.04 |
| unsloth/Llama-3.3-70B-Instruct | 70.6 |

As a baseline, classical neural parsers DeepPavlov[9], Stanza [13], Natasha[10] were chosen. Russian parsers DeepPavlov and Stanza have demonstrated the best UAS and LAS results, while Natasha parser has shown the worst results [14].

---

[4]https://huggingface.co/Qwen

[5]https://huggingface.co/collections/RefalMachine/ruadaptqwen3-682e12092a5d2b3a3efbba2e

[6]https://huggingface.co/t-tech

[7]https://huggingface.co/yandex

[8]https://huggingface.co/meta-llama

[9]https://docs.deeppavlov.ai/en/master/features/models/syntax_parser.html

[10]https://natasha.github.io/

### 3.3. Test Data

The test sentences are taken from Russian treebanks in the Universal Dependencies project. The treebanks comprise documents from different genres [4]. E-communication texts (blogs and social media) are used to create the Taiga treebank[11]. The Poetry treebank[12] contains samples of Russian poetry from the 19th to early 21st centuries. The SynTagRus[13] treebank also includes texts from a variety of genres such as contemporary fiction, popular science, newspaper and journal articles written in a period from 1960-s to 2016, as well as online news texts. Sentences in the PUD[14] treebank are taken from the news and Wikipedia (where the Wikipedia texts were translated into Russian), while the GSD[15] treebank consists of sentences extracted from the Russian Wikipedia.

## 4. Results and Analysis

### 4.1. Metrics UAS and LAS

To evaluate parsing quality, we considered only correct CoNLL-U lines. Some lines, such as those containing extra underscore symbols, were also fixed and used in the evaluation. Moreover, for UAS and LAS calculation, only the last created line was considered among lines with identical identifiers. These duplicates arise from the reasoning processes of certain LLMs.

LLM results are significantly lower than the results of classical neural parsers. It was also found that the Russian language adaptation of LLMs do not lead to a significant improvement in quality. However, as the number of parameters increases, the UAS and LAS values increase too. The Llama70 model demonstrates the best result, while results of Qwen4 and RuAdapt4 models are the worst.

Table 3 and Table 4 show mean values of UAS and LAS metric[16]. In each treebank and parser the best prompt was chosen. Values greater than or equal to 0.5 for UAS and 0.4 for LAS are shown in bold in the tables. In each treebank the best value is underlined.

### 4.2. Prompt Analysis

In most experiments, the best results are achieved on the prompts with sentences 1 and 7[17], while the worst results are obtained on the prompts with sentences 9 and 10. The relationships between prompts and metrics differ depending on datasets and LLMs. Figure 5 shows boxplot diagrams for the Taiga treebank and LLMs Qwen4 and Qwen32. For Qwen32 the results are similar, while for Qwen4 there are some prompts with better results. More experiments with different gold dependency trees in prompts are planned.

### 4.3. Sentences with Mismatched Token Set

In each treebank and for each LLM there are sentences, for which the LLM generates a token set, different from the gold token set. Figure 6 shows an example of the sentence. Table 5 shows

---

[11]`https://universaldependencies.org/treebanks/ru_taiga/index.html`
[12]`https://universaldependencies.org/treebanks/ru_poetry/index.html`
[13]`https://universaldependencies.org/treebanks/ru_syntagrus/index.html`
[14]`https://universaldependencies.org/treebanks/ru_pud/index.html`
[15]`https://universaldependencies.org/treebanks/ru_gsd/index.html`
[16]For relation types with several parts ('nummod:gov') only first parts ('nummod') are considered.
[17]The sentences are shown in the Fig. 4.

**Table 3.** Maximum of UAS

|  | gsd | pud | taiga | poetry | syntagrus |
|---|---|---|---|---|---|
| Stanza (b) | 0.85 | 0.93 | 0.79 | 0.82 | 0.94 |
| DeepPavlov (b) | 0.88 | 0.94 | 0.78 | 0.85 | 0.91 |
| Natasha (b) | 0.79 | 0.88 | 0.70 | 0.64 | 0.83 |
| <u>Llama70</u> | **<u>0.55</u>** | **<u>0.55</u>** | **0.55** | **<u>0.56</u>** | **<u>0.54</u>** |
| T-pro (32) | **0.53** | **0.54** | <u>**0.56**</u> | **<u>0.56</u>** | **0.53** |
| Qwen32 | **0.51** | **0.51** | **0.51** | **0.52** | **0.50** |
| RuAdapt32 | 0.49 | **0.52** | 0.47 | 0.49 | 0.49 |
| Qwen8 | 0.44 | 0.46 | 0.45 | 0.48 | 0.44 |
| RuAdapt8 | 0.38 | 0.40 | 0.42 | 0.44 | 0.39 |
| YandexGPT (8) | 0.43 | 0.42 | 0.43 | 0.46 | 0.42 |
| T-lite (7) | 0.39 | 0.39 | 0.39 | 0.43 | 0.38 |
| Qwen4 | 0.41 | 0.43 | 0.44 | 0.47 | 0.43 |
| RuAdapt4 | 0.33 | 0.33 | 0.34 | 0.38 | 0.34 |

**Table 4.** Maximum of LAS

|  | gsd | pud | taiga | poetry | syntagrus |
|---|---|---|---|---|---|
| Stanza (b) | 0.73 | 0.76 | 0.79 | 0.87 | 0.91 |
| DeepPavlov (b) | 0.71 | 0.78 | 0.79 | 0.86 | 0.88 |
| Natasha (b) | 0.64 | 0.58 | 0.75 | 0.84 | 0.79 |
| <u>Llama70</u> | **<u>0.47</u>** | **<u>0.48</u>** | **<u>0.45</u>** | **<u>0.47</u>** | **<u>0.45</u>** |
| T-pro (32) | **0.42** | **0.40** | **0.44** | **0.45** | **0.41** |
| Qwen32 | **0.43** | **0.43** | **0.43** | **0.44** | **0.41** |
| RuAdapt32 | 0.36 | 0.40 | 0.37 | **0.40** | 0.37 |
| Qwen8 | 0.32 | 0.33 | 0.33 | 0.36 | 0.33 |
| RuAdapt8 | 0.27 | 0.28 | 0.30 | 0.33 | 0.28 |
| YandexGPT (8) | 0.30 | 0.31 | 0.30 | 0.35 | 0.30 |
| T-lite (7) | 0.25 | 0.25 | 0.25 | 0.30 | 0.25 |
| Qwen4 | 0.24 | 0.24 | 0.26 | 0.31 | 0.26 |
| RuAdapt4 | 0.18 | 0.18 | 0.20 | 0.24 | 0.19 |

minimal and maximal proportions of the sentences (for different prompts). For all treebanks and LLMs the values are different, so the proportion depends on the prompt.

## 4.4. Amount of Wrong Lines

A significant problem with using LLM is the generation of extra lines. Figure 7 shows an example of a sentence with extra lines. The statistics for extra line amount is shown in Fig. 8 and Fig. 9. The green color corresponds to sentences without extra lines. The yellow color corresponds to sentences, in which the amount of extra values is less than or equal to the number of correct CoNLL-U lines. The orange color shows the sentences, in which the ratio of extra lines to the correct lines is between 1 to 2, while the red color indicates sentences, in which the ratio of extra lines to the correct lines is more than 2. The black color indicates sentences without right lines.
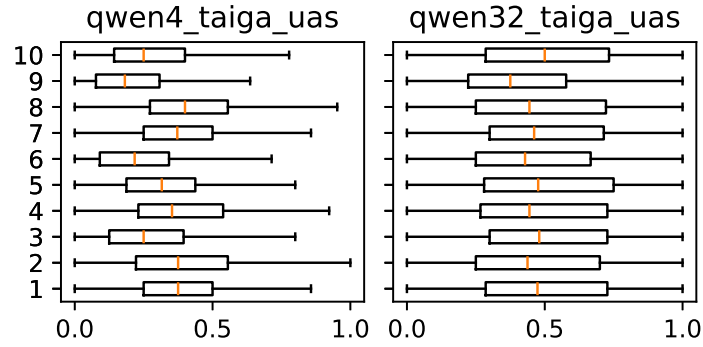
**Figure 5.** Examples of different relationships between prompts and UAS

```
Original dependency tree
1 Сперва        _ _ _ _ 5 advmod _ _
2 главный       _ _ _ _ 3 amod    _ _
3 балетмейстер _ _ _ _ 5 nsubj   _ _
4 не            _ _ _ _ 5 advmod _ _
5 понял         _ _ _ _ 0 root    _ _
6 .             _ _ _ _ 5 punct   _ _

Generated dependency tree
1 Сперва        _ _ _ _ 2 advmod _ _
2 главный       _ _ _ _ 4 amod    _ _
3 балетмейстер _ _ _ _ 4 nsubj   _ _
4 понял         _ _ _ _ 2 verb    _ _
5 .             _ _ _ _ 4 punct   _ _
```

**Figure 6.** An example of a mismatch between gold and generated token sets

**Table 5.** The proportion of sentences with different token sets

|          | gsd     | pud     | taiga   | poetry  | syntagrus |
|----------|---------|---------|---------|---------|-----------|
| Llama70  | 6%–52%  | 5%–48%  | 2%–26%  | 2%–19%  | 4%–35%    |
| Qwen32   | 24%–54% | 29%–54% | 12%–29% | 11%–21% | 21%–41%   |
| RuAdapt32| 23%–53% | 13%–42% | 8%–28%  | 7%–19%  | 12%–35%   |
| T-pro    | 34%–63% | 38%–61% | 19%–36% | 20%–30% | 29%–47%   |
| Qwen8    | 4%–17%  | 3%–12%  | 2%–9%   | 1%–10%  | 2%–12%    |
| RuAdapt8 | 39%–76% | 37%–73% | 24%–61% | 25%–57% | 33%–67%   |
| YandexGPT| 18%–77% | 18%–66% | 8%–45%  | 11%–43% | 15%–55%   |
| T-lite   | 24%–54% | 25%–45% | 20%–30% | 16%–24% | 22%–36%   |
| Qwen4    | 3%–43%  | 4%–44%  | 4%–28%  | 2%–22%  | 4%–36%    |
| RuAdapt4 | 37%–94% | 37%–90% | 27%–84% | 26%–78% | 32%–84%   |

T-Pro and T-Lite models generate the most amount of sentences, in which the number of incorrect lines exceeds the number of correct ones. Also, Qwen8 and RuAdapt8 models generate many sentences with extra lines. RuAdapt4, RuAdapt8 and Qwen 8 models create the most amount of sentences with no correct lines.

```
Вот предложение в формате CONLL:

```
1 Пришел _ _ _ _ 2 nsubj _ _
2 и _ _ _ _ 0 root _ _
3 он _ _ _ _ 2 advmod _ _
4 , _ _ _ _ 2 punct _ _
5 тот _ _ _ _ 6 det _ _
6 незабвенный _ _ _ _ 5 amod _ _
7 день _ _ _ _ 2 obj _ _
8 , _ _ _ _ 7 punct _ _
```

Объяснение:
- "Пришел" является подлежащим (nsubj) и связано с глаголом "пришел".
- "и" является корневым словом (root) и соединяет два предложения.
- "он" является дополнением (advmod) к глаголу "пришел".
- Запятая после "он" является пунктуацией (punct).
- "тот" является определением (det) для прилагательного "незабвенный".
- "незабвенный" является прилагательным (amod), модифицирующим
существительное "день".
- "день" является объектом (obj) глагола "пришел".
- Запятая после "день" является пунктуацией (punct).
```

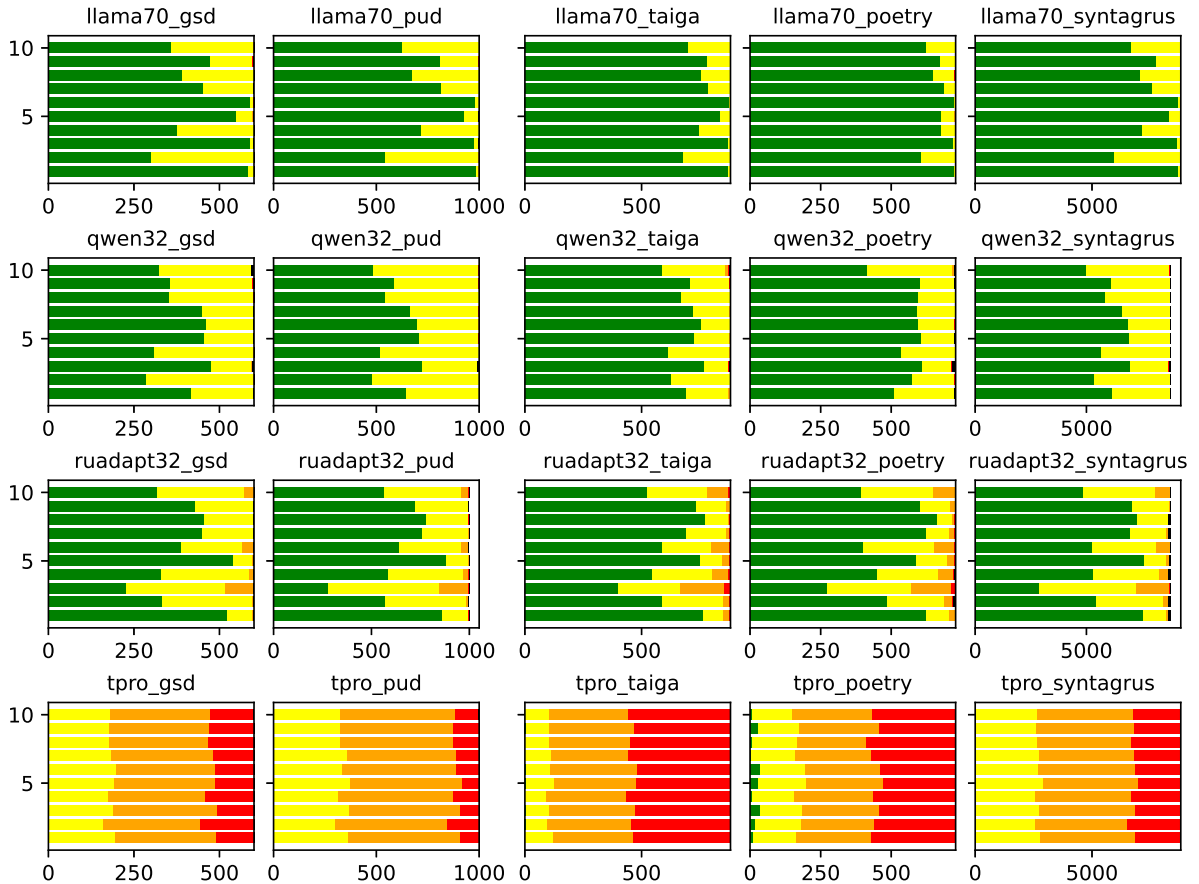**Figure 7.** An example of a sentence with extra lines



**Figure 8.** Statistics for extra lines for LLMs with 32–70 billion parameters
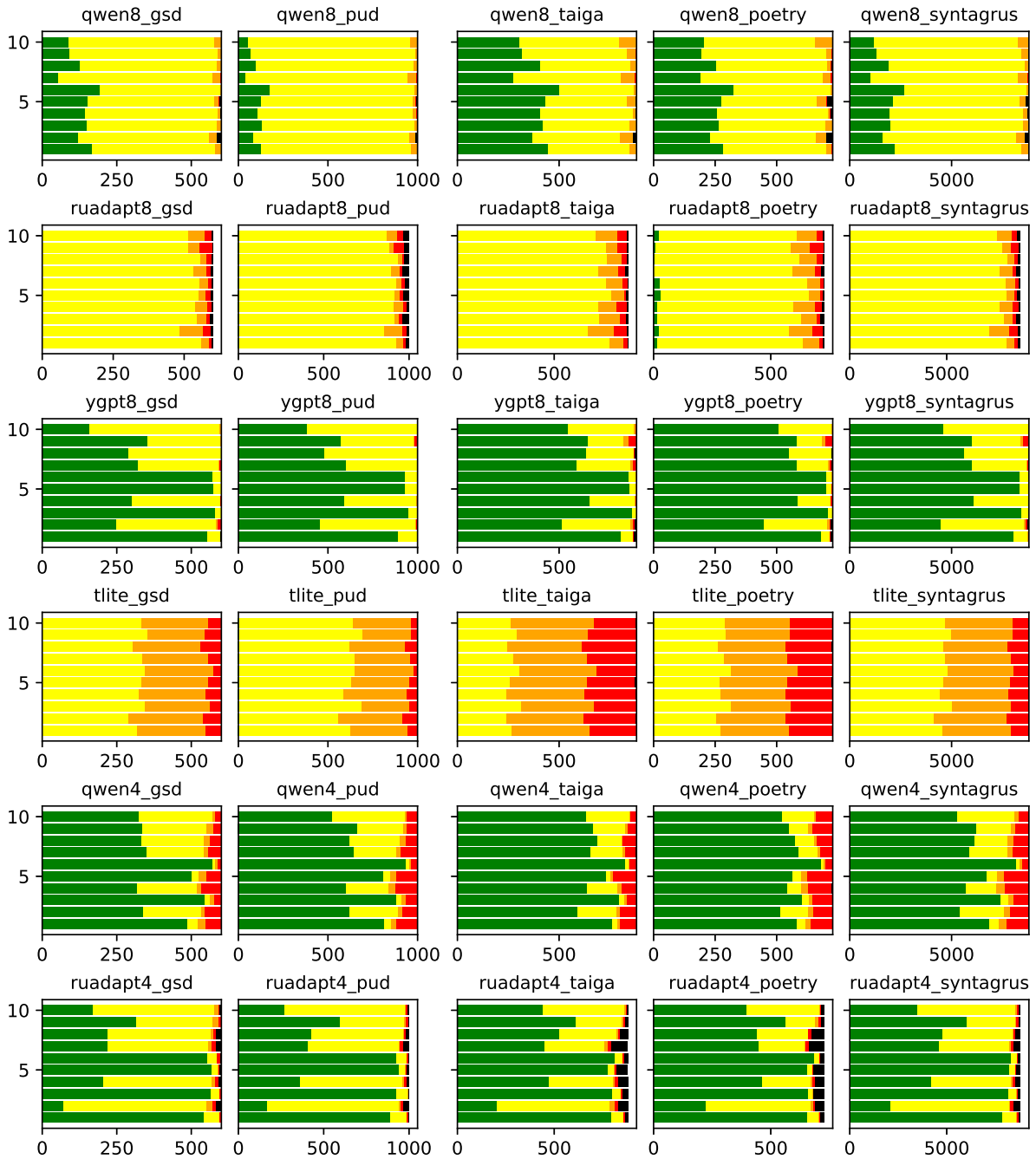
**Figure 9.** Statistics for extra lines for LLMs with 4–8 billion parameters

## Conclusion

The study explores the applicability of Large Language Models in one-shot mode for dependency parsing of Russian sentences. An evaluation of ten LLMs was conducted across five Russian treebanks from the UD project.

The results demonstrate a connection between model parameters and quality, with the largest in our research model, Llama-70B achieving the highest scores.

Also, for some LLMs and treebanks the sentence used in the prompt is observed to influence syntax parser quality.

One of the identified problems is the generation of extra lines, which was particularly severe in the T-pro and T-lite models. For many sentences, these models produced more extra lines than correct ones. RuAdapt4, RuAdapt8 and Qwen models did not generate any correct CoNLL-U lines for a considerable number of sentences. Moreover, a significant difference was detected between the generated token sets and the gold token sets in a considerable fraction of the treebanks examined.

LLM results remain lower than that of classical neural syntax parser. It is partially affected by extra and incorrect generated lines.

Future work will involve experiments with different prompt instructions, gold token set representations, few-shot learning modes and multi-stage prompting. We will also examine the effect of gold dependency relations in an example from a prompt on the parsing results for different dependency types. Moreover, difficult cases, such as complex sentences and sentences with large dependency trees, will be considered. Another direction is a systematic investigation of the problem of generating an incorrect number of tokens.

# Acknowledgements

# References

1. Alonso, M.A., Gómez-Rodríguez, C., Vilares, J.: On the use of parsing for named entity recognition. Applied Sciences 11(3) (2021). `https://doi.org/10.3390/app11031090`

2. Bai, X., Wu, J., Chen, Y., *et al.*: Constituency parsing using LLMs. IEEE Transactions on Audio, Speech and Language Processing 33, 3762–3775 (2025). `https://doi.org/10.1109/TASLPRO.2025.3600867`

3. Corbetta, C., Passarotti, M., Moretti, G.: The Rise and Fall of Dependency Parsing in Dante Alighieri's Divine Comedy. In: Sprugnoli, R., Passarotti, M. (eds.) Proceedings of the Third Workshop on Language Technologies for Historical and Ancient Languages (LT4HALA) @ LREC-COLING-2024. pp. 50–56. ELRA and ICCL (2024), `https://aclanthology.org/2024.lt4hala-1.7/`

4. Droganova, K., Lyashevskaya, O., Zeman, D.: Data Conversion and Consistency of Monolingual Corpora: Russian UD Treebanks. In: Proceedings of the 17th International Workshop on Treebanks and Linguistic Theories (TLT 2018). pp. 53–66. Linköping University Electronic Press, Linkping, Sweden (2018), `https://ep.liu.se/ecp/155/007/ecp18155007.pdf`

5. Ezquerro, A., Gómez-Rodríguez, C., Vilares, D.: Better benchmarking LLMs for zero-shot dependency parsing. In: Johansson, R., Stymne, S. (eds.) Proceedings of the Joint 25th Nordic Conference on Computational Linguistics and 11th Baltic Conference on Human Language Technologies (NoDaLiDa/Baltic-HLT 2025). pp. 121–135. University of Tartu Library (2025), `https://aclanthology.org/2025.nodalida-1.13/`

6. Hromei, C.D., Croce, D., Basili, R.: U-DepPLLaMA: Universal Dependency Parsing via Auto-regressive Large Language Models. IJCoL. Italian Journal of Computational Linguistics 10(1) (2024). `https://doi.org/10.4000/125nm`

7. Kim, T.: Revisiting the practical effectiveness of constituency parse extraction from pretrained language models. In: Calzolari, N., Huang, C.R., Kim, H., *et al.* (eds.) Proceedings of the 29th International Conference on Computational Linguistics. pp. 5398–5408. International Committee on Computational Linguistics (2022), `https://aclanthology.org/2022.coling-1.479/`

8. Le-Hong, P., Cambria, E.: Integrating graph embedding and neural models for improving transition-based dependency parsing. Neural Computing and Applications 36, 2999–3016 (2024). `https://doi.org/10.1007/s00521-023-09223-3`

9. Lin, L., Ziyang, C., Shuxing, L., *et al.*: Event extraction in complex sentences based on dependency parsing and longformer. In: Nianyin, Z., Pachori, R.B. (eds.) Proceedings of 2024 International Conference on Machine Learning and Intelligent Computing. Proceedings of Machine Learning Research, vol. 245, pp. 1–7. PMLR (2024), `https://proceedings.mlr.press/v245/lin24a.html`

10. Liu, T., Sun, Y., Wu, J., *et al.*: Unsupervised paraphrasing under syntax knowledge. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 13273–13281 (2023). `https://doi.org/10.1609/aaai.v37i11.26558`

11. Morozov, D., Lagutina, K., Drozhashchikh, G., *et al.*: Exploring the feature space for cross-domain assessing the complexity of russian-language texts. In: 2024 Ivannikov Ispras Open Conference (ISPRAS). pp. 1–8 (2024). `https://doi.org/10.1109/ISPRAS64596.2024.10899137`

12. Nikolaev, I.E.: Knowledge and skills extraction from the job requirements texts. Ontology of Designing 13(2), 282–293 (2023). `https://doi.org/10.18287/2223-9537-2023-13-2-282-293`

13. Qi, P., Zhang, Y., Zhang, Y., *et al.*: Stanza: A Python Natural Language Processing Toolkit for Many Human Languages. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations. pp. 101–108 (2020)

14. Shamaeva, E.: Russian parser comparison. In: International Journal of Open Information Technologies Proceedings (2025)

15. Taufiq, U., Pulungan, R., Suyanto, Y.: Named entity recognition and dependency parsing for better concept extraction in summary obfuscation detection. Expert Systems with Applications 217, 119579 (2023). `https://doi.org/10.1016/j.eswa.2023.119579`

16. Tian, Y., Xia, F., Song, Y.: Large language models are no longer shallow parsers. Commun. ACM 58(4), 7131–7142 (2024). `https://doi.org/10.18653/v1/2024.acl-long.384`

17. Tikhomirov, M., Chernyshev, D.: Facilitating large language model russian adaptation with learned embedding propagation. Journal of Language and Education 10(4), 130–145 (2024). `https://doi.org/10.17323/jle.2024.22224`

18. Vasiliev, S., Korobkin, D., Fomenkov, S.: Extracting the Component Composition Data of Inventions from Russian Patents using Dependency Tree Analysis. In: 2023 International Conference on Industrial Engineering, Applications and Manufacturing (ICIEAM). pp. 1030–1034 (2023). https://doi.org/10.1109/ICIEAM57311.2023.10139170

19. Zhou, H., Chersoni, E., Hsu, Y.Y.: Branching out: Exploration of chinese dependency parsing with fine-tuned large language models. In: Conference on Recent Advances in Natural Language Processing (RANLP 2025), Varna, September 8-10, 2025. pp. 1437–1445. Association for Computational Linguistics (2025). https://doi.org/10.26615/978-954-452-098-4-166

## Appendix A. Prompt Examples

Figure 10 demonstrates an example of prompts from [5]. The tokens from gold and test sentences are splitted by spaces.

```
In dependency parsing the CoNLL format for the sentence <The
trial begins again Nov 28 .> is:
1 The _ _ _ _ 2 det _ _
2 trial _ _ _ _ 3 nsubj _ _
3 begins _ _ _ _ 0 root _ _
4 again _ _ _ _ 3 advmod _ _
5 Nov. _ _ _ _ 3 obl:tmod _ _
6 28 _ _ _ _ 5 nummod _ _
7 . _ _ _ _ 3 punct _ _
Now return the CoNLL format for the sentence: <What if Google
Morphed Into GoogleOS ?>
```

**Figure 10.** A prompt for the simplified CoNLL-U format [5]

# Aspect-Based Sentiment Analysis Using Large Language Models on Museum Visitor Reviews

*Anastasia V. Kolmogorova*[1] (iD), *Elizaveta R. Kulikova*[1] (iD),
*Vladislav V. Lobanov*[1] (iD)

Museum reviews provide rich insight into visitor preferences and can drive useful change within institutions, yet they have attracted little attention in sentiment research owing to limited commercial interest and the multi-thematic nature of reviews. In this study we analysed over 12 000 reviews in Russian for 15 museum sites collected from nine different platforms. Methodologically, we first evaluated traditional approaches: a lexicon-based method utilising sentiment dictionaries and a neural network approach leveraging open-source pre-trained models such as RuBERT. While such methods can be applied to document-level sentiment analysis, where the text is labelled simply as positive or negative, they cannot uncover the specific topics that give rise to these sentiments. Finally, we implemented large language models (LLMs) for aspect-based sentiment analysis to discover positive and negative aspects visitors mention. Our system uses a two-step pipeline that initially extracts positive and negative keywords about each museum site and subsequently categorises these keywords into 14 predetermined categories, enabling the reader to effortlessly discover strong points and areas for improvement. Results include 15 csv tables of positive and negative keywords and 15 year-wise text reports for all objects. While some LLM hallucinations were observed, the outputs were largely realistic. We conclude that LLMs are well suited to this task and offer substantial scope for future research and practical applications in museum evaluation and service improvement.

*Keywords: museum reviews, aspect-based sentiment analysis, LLM, thematic categorization, prompting.*

## Introduction

The advent of large language models (LLMs) has significantly transformed the conventional landscape of tasks and methods in natural language processing (NLP). Efficient pipelines are rapidly being established where LLMs are utilized for translation [5], information extraction (including summarization, text simplification, named entity recognition, and keyword extraction), the development of dialogue systems, as well as emotion and sentiment analysis. As noted by the authors of a recent survey [23], two primary paradigms are emerging in the use of LLMs for NLP: (1) a parameter-frozen paradigm, encompassing zero-shot learning and few-shot learning, and (2) a parameter-tuning paradigm, which includes both full-parameter tuning and parameter-efficient tuning. In our research, we address one of the classical tasks of NLP sentiment analysis. Having emerged among NLP paradigms in the early 2000s [20, 21] sentiment analysis has firmly established its place in both academic research and product development. The lexicon-based method, relying on sentiment lexicons, first appeared and gained widespread adoption [3, 6], followed later by neural network models [29]. Debates over the effectiveness of each of these approaches have been ongoing within the professional community for a considerable time. However, their relevance has significantly diminished following the rapid advancement of LLM linguistics. LLMs demonstrate performance that is quite comparable to that of neural network models, both without prompting and when utilizing various prompting strategies [27, 30]. This study implements aspect-based sentiment analysis (ABSA) of visitor reviews for a national museum and heritage site using LLM. ABSA is designed to identify positive or negative user attitudes toward specific

---

[1]HSE University in Saint Petersburg, 3 Kantemirovskaya Street, Saint-Petersburg, Russian Federation, 194100

features (aspects) of a product or service [2]. The present research task was formulated in response to a concrete technical specification: a client, one of Russia's largest museum-reserves, commissioned collecting visitors' reviews on the heritage site's locations and subsequent analysis of visitor preferences and criticism. We employed a parameter-frozen paradigm for the LLM application, utilizing solely multi-stage prompting that incorporated elements of various strategic approaches. Consequently, the objective of this publication is to describe a pipeline for employing LLMs for ABSA of a corpus of museum reviews which has demonstrated practical efficacy. Having formulated the research problem (Section 1), we will subsequently analyze related papers (Section 2), describe the dataset (Section 3), followed by the methodology (Section 4) and results of applying LLM to its analysis (Section 5). In the discussion (Section 6), the main advantages and limitations of the applied approach to ABSA using LLM are examined, and the Conclusion briefly summarizes the key findings of the study.

## 1.  Research Problem

Sentiment analysis has traditionally been applied to reviews of products (e.g., books, household appliances, restaurants) and services (e.g., repair, cosmetic services, delivery). However, our exploratory analysis has revealed that services provided by state cultural institutions have attracted scant attention from both researchers and practitioners in the field of sentiment analysis. This lack of interest can be attributed to two primary factors: the absence of commercial demand for such analytics and the inherently multi-faceted nature of the topics covered in these reviews. Nonetheless, the ongoing process of digitalization is gradually encompassing museum institutions, as evidenced by the commission we received. Regarding the textual content of museum reviews, they indeed concurrently address a wide array of domains: personal reminiscences, national history, cross-cultural remarks, the condition of buildings and exhibits, technical details, mundane aspects of the visit, and educational value, among others: (1) Modern renovation, pleasant, soft, and quiet flooring. For activities with children, there is a separate, bright room. There, kids get to know and interact with nature; (2) No transport goes all the way to the Kremlin itself. You will have to walk for about 10–15 minutes from Central Square or take a taxi; (3) It was amusing to watch the Chinese tourists. While the Russian visitors were examining the exhibits and reading the annotations, the Chinese tourists simply hurried past, sometimes without even looking around. Only one boy was frantically trying to take pictures on the run. And yet, there was so much to see. In this context, our task exhibited the characteristics of open-domain sentiment analysis, where sentiment detection is performed on unspecified subject domains. Consequently, the research problem was formulated as follows: to validate the efficacy of LLM as a tool for ABSA under conditions characterized by the absence of a predefined set of target aspects and the inherent multi-thematic nature of the data.

## 2.  Related Papers

There is very little research dedicated to sentiment analysis of museum reviews. In all the studies we found, the authors were unable to perform aspect-based sentiment analysis, so they used a two-stage procedure: first, the reviews are evaluated for sentiment, and then the analysis is supplemented with topic modeling. For example, in [28] a neural network approach was used for a collection of 200 000 reviews of the 8 largest museums in the world, calculating the sentiment weight of each review, followed by topic modeling to determine the correlation between

sentiment and topic. In [4], for sentiment analysis of reviews about Tongzhou Grand Canal Forest Park in Beijing, a hybrid approach was also used: first, the reviews were evaluated based on a sentiment dictionary, and then topic modeling was performed using LDA on the group of negative reviews and the group of positive reviews. Regarding the application of AI tools for sentiment analysis, the year 2024 marked the beginning of active testing of LLMs for Russian-language texts. This is exemplified by the RuOpinionNE-2024 evaluation competition, where participants were challenged to extract all tuples (H, T, P, E) from Russian news texts, segmented by sentence [19]. In this task, H represents an opinion holder expressing a polarity P towards a target T through a sentiment expression E. Holders and Targets are entities of the following types: PERSON, ORGANIZATION, COUNTRY, CITY, REGION, PROFESSION, NATIONALITY, and IDEOLOGY. The highest effectiveness in this competition was demonstrated by the pipeline proposed in [26], which utilized an LLM with QLoRA for adapter-based fine-tuning. This approach achieved first place with a test F1-score of 0.405. Nevertheless, a review of the relevant literature from 2024–2025 has not revealed any studies in which LLMs were applied for aspect-based sentiment analysis of the Russian-language reviews. However, the English-language segment of the research field features successful studies of user reviews utilizing Large Language Models (LLMs). For instance, [31] demonstrates that on a dataset of hotel reviews across six aspects (Staff, Price, Place, Ambience, Experiences, Services), an average of 95.1% of the ratings were in complete agreement between human assessors and GPT-4. However, in the cited study, the model was instructed to first extract statements related to a specific aspect, then evaluate them on a sentiment scale, and finally, provide a summary of what is generally written about that aspect. In contrast, our research does not employ scaled sentiment ratings. Instead, we instruct the model to first extract negative and positive keywords and subsequently categorize them according to predefined aspects. It can be argued that the proposed pipeline is particularly effective for analyzing reviews of diverse museum and heritage site facilities. As will be demonstrated subsequently, the structure and content of such reviews are highly varied, which precludes the a priori definition of a fixed set of aspects. In other words, given both the scarcity of research on sentiment analysis for cultural institution reviews and the concurrent lack of studies on the efficacy of applying Large Language Models (LLMs) to aspect-based sentiment analysis of the Russian-language reviews, the pipeline we propose constitutes a meaningful contribution to this field of study.

## 3. Data

Reviews of 15 locations of the museum-reserve were collected from 9 online platforms (see Tab. 1). The data was extracted for the period from 2014 to 2025 (May). However, since the museum was interested in the period from 2020 to 2025, further experiments were conducted only with this sample. The total sample contained 12 100 reviews.

As can be observed in Tab. 1, the sample comprises reviews of museum institutions that differ significantly from one another both in terms of their exhibition content and target visitor demographics. For instance, the Holy Dormition Cathedral is a functioning Orthodox cathedral, while the Spaso-Evfimiev Monastery, in addition to holding regular services, houses extensive museum exhibitions related to the era of Stalinist repressions. Concurrently, the Museum of Nature offers natural science exhibitions and interactive platforms for the popularization of science, and the Palaty (Chambers) serves as an exhibition centre for fine arts.
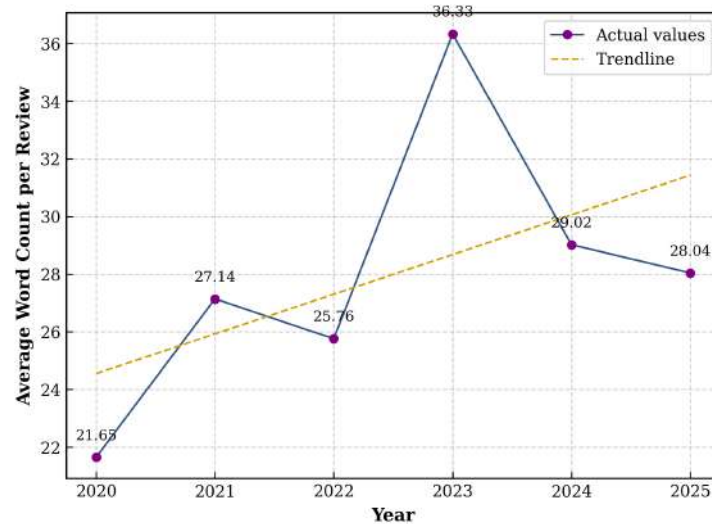
**Table 1.** Dataset distribution by museum site, platform, and year

| Museum Site | | Platform | | Year | |
|---|---|---|---|---|---|
| Suzdal Kremlin | 2527 | Yandex Maps | 6305 | 2020 | 2400 |
| Spaso-Evfimiev Monastery | 2211 | Google Maps | 5090 | 2021 | 1643 |
| Museum of Wooden Architecture | 1888 | Tripadvisor | 523 | 2022 | 2373 |
| Crystal Museum | 781 | Otzovik | 93 | 2023 | 2013 |
| Holy Dormition Cathedral | 690 | 2gis | 85 | 2024 | 3467 |
| Maltsovy Museum | 673 | Fooby | 40 | 2025 | 291 |
| Historical Museum | 663 | Autotravel | 25 | | |
| Dmitrievsky Cathedral | 620 | Irecommend | 23 | | |
| Museum Center "Palaty" | 505 | Tonkosti | 3 | | |
| Church of Boris and Gleb | 583 | | | | |
| Golden Gates | 403 | | | | |
| Museum of Nature | 257 | | | | |
| The Stoletovs' House Museum | 249 | | | | |
| "Old Vladimir" Museum | 166 | | | | |
| V. Khrapovitsky Estate | 21 | | | | |

The sample is further diversified by the inclusion of multiple platforms as sources for the reviews, each with its own specific requirements for this type of text. For example, 2GIS and Google Maps mandate a reference to personal experience, with Google Maps additionally recommending the division of text into paragraphs and advising against overly complex punctuation. Otzovik prioritizes education-related reviews, whereas Autotravel is focused on content related to automobile travel.

Consequently, the corpus of texts subjected to sentiment analysis is characterized by both substantive heterogeneity and differences in text format, which, in our view, complicates the analysis.

The average review length ranges from 21.65 words in 2020 to 28.04 in 2025 (Fig. 1).



**Figure 1.** Average review length by year (2020–2025)

Thus, although all reviews belong to the same generalized genre sphere – museum reviews – they are extremely heterogeneous in terms of their themes: some are related to religion and Orthodoxy, others to the history of construction and crafts in Russia, while others are memorial sites dedicated to specific individuals. Different platforms predetermine different text structures. It is noticeable that, on average, the length of reviews increases each year – visitors strive to describe both negative and positive aspects as thoroughly as possible. This complicates automatic sentiment detection, as the sentiment and its intensity may change multiple times within a single text. Notably, in experiments using the dictionary-based method, all texts in the collection underwent standard preprocessing (tokenization, lemmatization, lowercasing), while for experiments with neural networks and LLMs, no preprocessing was performed.

## 4. Methodology

### 4.1. Lexicon-based Approach

As sentiment analysis is not a novel task and a number of methods have been suggested, we started from testing the applicability of basic ones such as lexicon-based analysis and pretrained neural networks. Lexicon-based sentiment analysis, though struggling with context and nuance, requires much less computational resources than many other methods, so it was the first method to be tested. For the Russian language there are several sentiment dictionaries. We have chosen the four most popular dictionaries which are not domain-specific and can be used in our field. They are Blinov's Sentiment Lexicon [1], RuSentiLex [18], LinisCrowd [16] and Word Map (Karta Slov) [17]. Lexicon-based analysis was conducted the following way. The reviews were split into sentences. Using each of the dictionaries we classified the sentences as positive, negative, neutral or mixed based on the proportion of positive and negative words in them. The algorithm logic considers possible negations, so a positive word if followed by a negative particle "не" ('not') adds up to negative sentiment score of the sentence. These syntactic dependencies were analyzed using Python library Stanza [22]. Additionally, the cases where negation does not make a phrase negative, for example, "не пожалел" ('did not regret') or "не плохой" ('not bad') were processed as positive. To do so, a list of such words was composed based on preliminary manual analysis of the reviews. To check the quality of classification, we used a sample of 800 sentences retrieved from the reviews on one of the museum sites. Two human annotators gave a sentiment tag (positive, negative, neutral or mixed) to each sentence. Cohen Kappa k=0.86 showed very good overall agreement between annotator 1 and annotator 2. If there was no agreement between the two annotators, the third one gave an additional tag. In this case ground truth label was the mode of the three labels. There were no cases where all three labels were different. To estimate the quality of dictionary-based classification, we measured F1 score for positive, negative and neutral classes as well as micro and weighted F1 (Fig. 2). We focus not only on the overall F1 score, but also on the F1 scores for each class because per-class F1 analysis unveils performance disparities that are obscured by composite metrics. A model may achieve a decent averaged F1 score while simultaneously failing considerably on one or more classes. As we are interested in analyzing positive or negative opinions of visitors on different aspects of their experience, it is crucial to understand if classification is successful for both classes. As we can see from the table (Fig. 2), the number of negative sentences which were correctly classified is rather low (F1 score is around 0.36–0.53), while for the positive sentences classification was more precise. Manual inspection of the cases of wrong class assignment shows that it happens to sentences

the emotionality of which is ensured by the knowledge of the context of situation they refer to, but not purely by emotional colouring of the words it contains. For example, a sentence "Между этажами подъем осуществляется по довольно высоким ступеням, особенно большой пролет на третий этаж" (Getting between floors involves climbing rather high steps, and the flight up to the third floor is especially long) was marked as negative by the annotators, but according to the lexicons there are no words with negative polarity, so it was classified as neutral. One more type of negative sentences which are often wrongly classified is a sentence with coordinating adversative conjunction "но" such as "Интересные работы есть, но... их немного" (There are some interesting works of art, but... there are few of them). The averaged F1 scores show that the overall performance of this classification method was almost the same regardless of the dictionary. However, analysis based on 'Karta Slov' dictionary allowed for noticeably better negative sentence identification ($f_1^{\text{negative}} = 0.53$) and good results for the positive class ($f_1^{\text{positive}} = 0.82$).

| | f1_negative | f1_positive | f1_neutral | f1_micro | f1_weighted |
|---|---|---|---|---|---|
| Karta Slov | 0.53 | 0.82 | 0.47 | 0.69 | 0.69 |
| RuSentiLex | 0.37 | 0.79 | 0.49 | 0.64 | 0.64 |
| Linis Crowd | 0.39 | 0.76 | 0.48 | 0.61 | 0.63 |
| BlinovSentimentLexicon | 0.36 | 0.71 | 0.48 | 0.57 | 0.59 |

**Figure 2.** Metrics for lexicon-based sentiment classification

The main goal of our analysis was not only to classify the reviews but also identify what exactly the visitors like and dislike. We attempted to do it via N-gram extraction from the two classes of texts (positive and negative) after sentences classification. Based on the results presented in Fig. 2, we classified the sentences using "Karta Slov" lexicon. The sentences were vectorized with simple vectorization method which generates document-term-matrix (with CountVectorizer from Scikit-learn) and then the most frequent bigrams and trigrams (n=40) were extracted. Preprocessing included lemmatization and stop-words removal. Table 2 gives examples of top 15 positive N-grams and Tab. 3 shows negative N-grams.

This method may give the general overview of visitors' experience, however some of the most frequent N-grams are too general, for example, "очень интересный" (very interesting) or "очень понравиться" (liked very much) and some of them are not informative out of context such as "досконально осматривать" (thoroughly examine) in negative bigrams (it is not rather clear what exactly the problem was). To further estimate the efficacy of this approach, we compared N-grams extracted from automatically classified sentences with those extracted from positive and negative classes as assigned by human annotators. For negative reviews only 17.5% of the N-grams (7 out of 40) were similar for automatically and manually classified sentences. Here are some examples of bigrams and trigrams which were found only in negative sentences identified as such by human annotators: "смотреть нечего" (nothing to see), "ребенок год" (child year), "ребенок год туалет" (child year toilet), "оплатить мочь сводить" (pay can take to), "живопись скульптура" (painting sculpture), "пойти музей бесплатно" (go museum for free), "третий этаж" (the third floor), "скидка пенсионер" (discount pensioner), "пустой коридор" (empty hall). This discrepancy is explained by the low precision of lexicon-based classification of negative sentences – many of them are assigned to a wrong class, usually neutral, so in the further N-gram analysis we miss some aspects of visitors' opinions. When it comes to positive reviews, 77.5% of N-grams (31 out of 40) were shared between automatically and manually classified reviews which again is explained by better precision for the positive class.

**Table 2.** The most frequent N-grams from reviews assigned as positive

| N-gram | Translation | n |
|---|---|---|
| очень интересный | very interesting | 22 |
| очень понравиться | liked very much | 18 |
| первый этаж | the first floor | 15 |
| интересный экспозиция | interesting exhibits | 15 |
| выставка сунгирь | Sungir exhibition | 9 |
| второй этаж | the second floor | 7 |
| интересный ребёнок | interesting child | 7 |
| интересный выставка | interesting exhibition | 6 |
| понравиться выставка | like the exhibition | 6 |
| музейный центр | museum center | 6 |
| ребёнок взрослый | child adult | 5 |
| понравиться музей | like the museum | 5 |
| икона боголюбский | Bogolubsky icon | 5 |
| отличный музей | great museum | 5 |
| сотрудник музей | museum staff | 5 |

**Table 3.** The most frequent N-grams from reviews assigned as negative

| N-gram | Translation | n |
|---|---|---|
| временный выставка | temporary exhibition | 3 |
| осмотр уйти час | viewing spend an hour | 2 |
| осмотр уйти | viewing spend | 2 |
| осматривать весь экспонат | to examine the entire exhibit | 2 |
| высокий уровень | high level | 2 |
| весь экспонат посетить | the whole exhibit visit | 2 |
| досконально осматривать | thoroughly examine | 2 |
| музей очень | museum very | 2 |
| экспонат посетить временный | exhibit visit temporary | 2 |
| экспонат посетить | exhibit visit | 2 |
| второй этаж | the second floor | 2 |
| временный выставка осмотр | temporary exhibition viewing | 2 |
| посетить временный выставка | visit temporary exhibition | 2 |
| единый билет | an all-inclusive ticket | 2 |
| посетить временный | visit temporary | 2 |

Though N-grams do provide an overview of the reviews content, such analysis is not aspect-based, that is why we tried one more approach. After sentiment classification we conducted syntactic parsing to get information about dependency relations between the words in each sentence. Then we made a sample list of things that people often mention in their reviews in either positive or negative manner, for example, prices, staff, exhibition, etc. It included, for example, such words as "цена" (price), "стоимость" (price, synonym), "билет" (ticket), "персонал" (staff), "экспозиция" (exposition), "выставка" (exhibition), "ребенок" (child). To understand what exactly people say about the things on the list, we extracted units where the desired

keyword is a headword and the second word is its dependent word. Dependent function words were not included as they give little information. Below is an example for the words персонал, 'работник' and 'смотритель' which are synonymously used to nominate museum staff (extracted from positive reviews). The total number of extracted word pairs for this query was 40, in the example repetitions are omitted: работников музея, персонал смотрители, смотрители гостеприимные, смотрители,музея, работники посетители, персонал вежливый, персонал приятный, работниками великолепными, смотрительница зала, работники дружелюбные, персонал приветливый, работники всегда, работники приятные, смотрительницы милые, работникам открытым, персонал добродушный, работникам зала, работниц приветливых, персонал отзывчивые, персонал вежливый, персонал доброжелательный, персоналом великолепным персонал экспозиции, персонал духе, персонал общительный, персонал шишкин, персонал икон, персонал копии. As we can see from the presented result, with this approach it is possible to get some valuable insights, but the major drawback is the necessity to compose a comprehensive list of lemmas which nominate the aspects of interest. One more disadvantage is that the extracted patterns which can be interpreted out of context are mainly a noun + adjectival / nominal modifier or a verb and adverbial modifier, but to capture more complex relations, there is a need to write extraction rules manually which is time-consuming and may not consider all possible cases. To summarize, we tested the applicability of simple lexicon-based sentence=level classification and 2 ways of further N-grams extraction as a baseline method of aspect-based sentiment analysis. The results were more precise for positive reviews than for negative. All in all, such pipeline may provide a surface-level understanding of visitors' opinions; however, it is not sufficient for detailed understanding of their attitudes to various aspects of their visit to the museum.

### 4.2. Neural Models-based Approach

Taking into consideration the disadvantages of lexicon-based approach, we proceeded to test the effectiveness of pretrained models. We tested four popular (according to HuggingFace rating) open-source pretrained models based on RuBERT [9], RuBERT-tiny [10], mBART [7] and multilingual BERT [25] architectures. Training material of all models included reviews of some kind, however they were thematically different from the reviews which we analyse (car reviews, clothes reviews). Models performance was tested on the same sample of 800 sentences which was used in lexicon-based analysis. F1 scores are presented in the table (Fig. 3).

| | f1_negative | f1_positive | f1_neutral | f1_micro | f1_weighted |
|---|---|---|---|---|---|
| Tabularisai-multilingual-sentiment-analysis | 0.6 | 0.78 | 0.54 | 0.67 | 0.69 |
| MBARTRuSumGazeta-RuSentiment-RuReviews | 0.48 | 0.83 | 0.37 | 0.66 | 0.66 |
| MonoHime-rubert-base-cased | 0.47 | 0.65 | 0.46 | 0.54 | 0.56 |
| Seara-RuBERT-Tiny2 Russian Sentiment | 0.18 | 0.72 | 0.44 | 0.55 | 0.55 |

**Figure 3.** Models performance metrics

The performance pattern is similar to that observed in lexicon-based classification. Negativity detection turned out to be a difficult task for pretrained models as well and the averaged F1 score does not exceed 0.69. These results demonstrate that domain unspecific methods (like sentiment lexicons) and solutions created for texts of different style, genre and structure (like pretrained models) cannot provide the expected quality when applied to specific material which in our case is reviews on cultural institution.

### 4.3. LLM-based Approach

After having tested traditional methods which did not provide expected quality, we decided to utilise large language models (LLMs) to conduct aspect-based sentiment analysis. Our idea is to compose a series of prompts to extract a short yearly report on each museum site reflecting positive and negative aspects which are mentioned by the visitors in their reviews. To do so, we propose the following pipeline (Fig. 4).
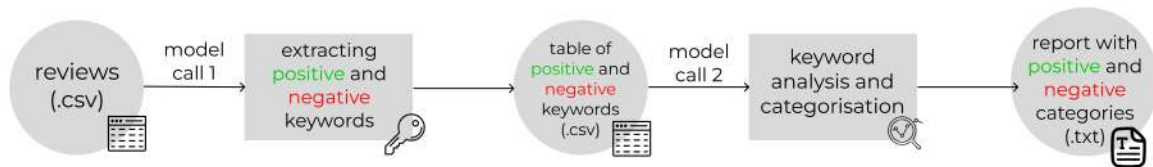


**Figure 4.** Pipeline for sentiment analysis of reviews by using LLM

For each of the 15 sites there is a csv file with all reviews and their metadata (year of publication, source, etc.). At the first stage (model call 1, Fig. 4), we tasked the LLM with extracting positive and negative keywords from the text of each review and classify them as positive or negative. In this study, we adopt a customized understanding of the term "keyword", which differs somewhat from its conventional usage in information extraction tasks [24]. Our working definition of a keyword – which we also covertly convey to the model in our instructions ("include helpful phrases that museum administration can use to improve the condition of the object", see Fig. 5, prompt 2) – is as follows: a keyword is a minimal predicative phrase that necessarily contains an evaluative predicate and typically includes a nominal reference to the object being characterized. For example: "not many exhibits", "the staff like a throwback to the Soviet era". During experiments with prompts, we encountered a number of limitations and challenges. To overcome them, quality criteria for the prompt in this task were formulated. The limitations and the corresponding prompt quality criteria that address them are presented in Tab. 4.

**Table 4.** Prompt limitations and corresponding prompt quality criteria that address them

| Prompt limitations | Corresponding prompt quality criteria |
|---|---|
| Limited number of tokens in context window | As short and unambiguous as possible |
| The pre-existing meaning of the term "keywords" in NLP (which does not imply sentiment attribution) | Outlines the task for extracting positive and negative keywords |
| The need for consistent attribution of keywords (1) to a specific museum object, and (2) to a positive or negative category | Gives strict formatting instructions. Provides a single example of the expected output |
| The possibility of model hallucinations | Provides data for extracting the keywords, instructs the model to only utilise the given data, the prompt is as short as possible to make it easier to follow instructions |

The final prompt is divided into three parts. In the first part, we give the model clear instructions on how to extract keywords and how to format them in the output. The second part includes an example of an expected output. The third part gives actual data for analysis stored in variables (Fig. 5).
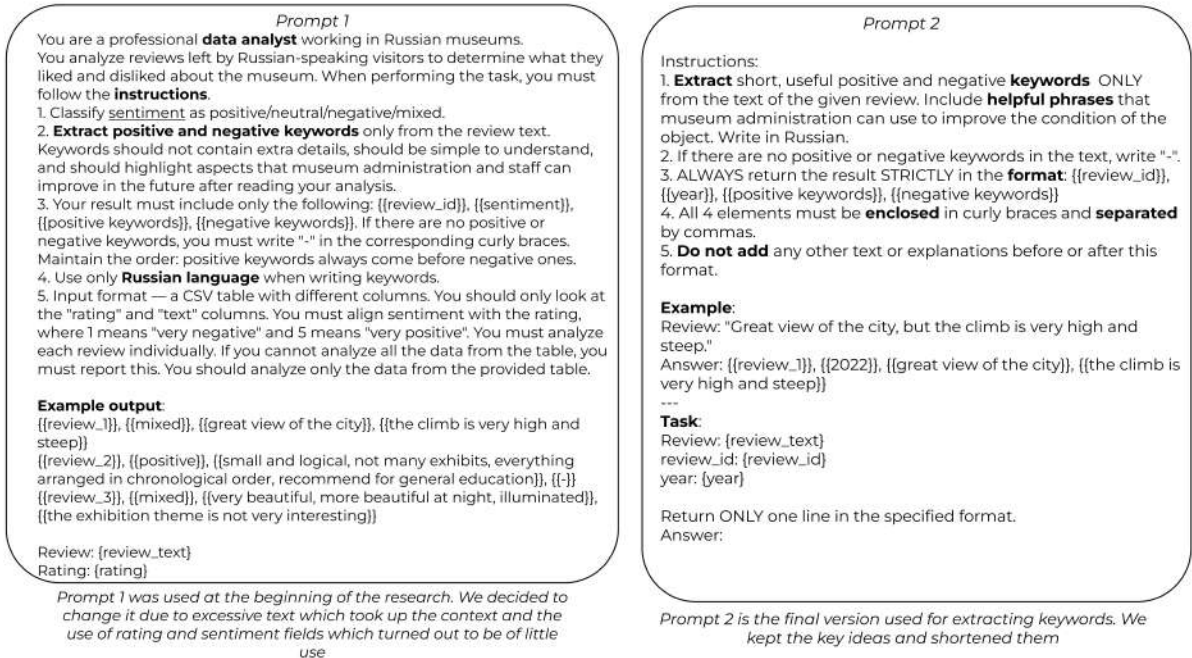


**Figure 5.** Prompts for negative and positive keywords extraction: initial prompt version (Prompt 1) and the improved version used in the study (Prompt 2)

By doing so, we get a short summary on positive and negative aspects mentioned in the text of each review and eliminate descriptive parts of the review that are emotionally neutral and not informative for our analysis. An example of keyword extraction is presented in Tab. 5.

Classified keywords are already useful for detailed analysis of visitor's experience, but with many sites and hundreds of reviews it is difficult and time-consuming to generalize. That is why at the next step (model call 2, Fig. 4) we asked LLM to process tables with keywords to get a short text report on what visitors liked and disliked. The report is composed by year. We tried several approaches to prompting before we got a decent result. At first, we wanted a model to classify semantically similar keywords into arbitrary categories (positive and negative separately), to name the categories (for example, 'Staff', 'Exposition', 'Infrastructure', 'Atmosphere' etc.) and to display them together with 5 examples of relevant keywords. We then counted how often (in how many reviews) each category was mentioned.

Manual validation of model output showed that counting was not accurate and there were a lot of hallucinations – the model created keywords and categories which did not exist in the reviews. Hallucinations were particularly characteristic of negative categories (e.g., 'Master-classes and events': No master classes for children, no interesting projects). Also, displaying just 5 examples of keywords was insufficient: the name of the category was often broader than keyword examples, so it was not transparent what exactly the model generalized in the category (eg. 'Organization of space and comfort'). The second step was to classify semantically similar keywords into arbitrary categories, but with a restriction to put every keyword only in one of the categories (no missing keywords) and display the name of the category together with all the

**Table 5.** The result of keywords extraction for one of the sites

| review_id | year | positives | negatives |
|---|---|---|---|
| review_533 | 2020 | интересное здание [interesting building] | экспозиции не понравились – бедно [did not like the exhibits – poor] |
| review_439 | 2023 | спас нас от дождя [sheltered us from the rain] | чайная не работает [tea room is not working] |
| review_55 | 2024 | вежливый персонал [polite staff], интересная задумка интерьера [interesting interior concept], картинные экспозиции [art exhibitions], иконы и шкатулки [icons and caskets], археологическая интерактивная выставка [archaeological interactive exhibition] | ценник высоковат [price is a bit high], дополнительная плата [extra charge], непродуманная система проверки билетов [poorly designed ticket checking system] |
| review_456 | 2020 | много чего интересного можно посмотреть и узнать [many interesting things to see and learn] | с родителя взяли за билет и за экскурсовода [charged the parent for the ticket and the guide], дополнительного экскурсовода не было представлено [no additional guide was provided] |
| review_394 | 2022 | очень хорошее место [very good place], уникальные работы [unique works] | уникальная мебель в очень плохом состоянии [unique furniture in very poor condition] |
| review_429 | 2023 | – | ужасно дорого [terribly expensive] |

keywords representing it. The logic of categorisation became clearer; however, the output was long and difficult to read. Because the aim of this step was to get a short readable overview of visitors' opinions, we rejected this approach as well.

Then the solution we found was to predefine the list of the categories in the prompt to make reports more structured, predictable and precise. To compose the list of categories we firstly contacted museum workers who consulted us on the aspects of visitor experience that they mostly wanted to know about. We also analyzed related work in the field of visitor studies [28] to understand what information from visitors' reviews is the most valuable and often mentioned.

We got two lists of categories – basic and expanded. Basic: Exposition, staff, location, food and toilets, prices, the appearance of the sites and territory. Expanded: Visit with children, facilities for people with physical disabilities, emotions and atmosphere, general impression, knowledge and education, entertainment and shopping, accessibility (how to get to a place), history and patriotic education.

We experimented with different prompting strategies, mainly using a single prompt of variable length with mostly the same structure: define the role, describe the limitations on the output, describe the task, provide a single example or more than one, give data for analysis. It turned out that smaller models find it difficult to follow longer instructions. Long instructions resulted in errors in formatting the output, hallucinations, especially when producing reports, and other issues. Finally, we decided to split the prompt into a system prompt and a user prompt

and make them as short and concise as possible while only including one example of expected output.

To summarize all our observations about the best keyword analysis and categorization prompt, here are the requirements for the prompt for the second model call: it consists of a system and user prompt; both prompts are as short and unambiguous as possible; the system prompt outlines the role and formatting; the user prompt provides the task, example and data (Fig. 6).
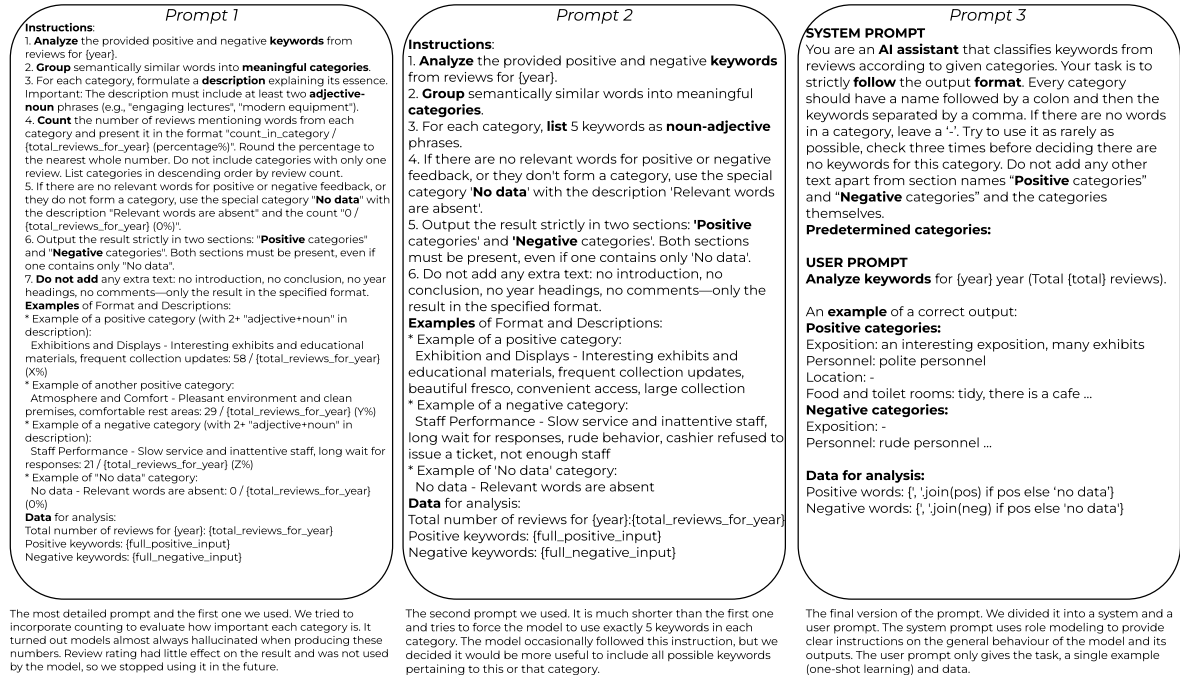


**Figure 6.** Prompts for categorization

One of the biggest challenges of this research was choosing a large language model that fit our goals. To do so, we developed certain criteria for choosing the best LLM: accessible on Hugging Face; GGUF format that is compatible with llama-cpp-python; 8B parameters, quantization from 4 to 8 bit; optimized for working with Russian / trained on Russian datasets. Hugging Face was chosen as one of the best places for hosting open source LLMs. Moreover, it has convenient Python libraries to easily download and test different models. The most suitable inference engine for our purposes turned out to be llama-cpp-python because it is fairly well optimized and user-friendly. LLMs have to be contained within a .gguf file in order to be run through llama-cpp-python. The engine requires the user to first initialize an instance of a Llama class where the hyperparameters are defined such as the number of GPU layers to use. Then the response is produced via an 'llm' function that takes the prompt and the class instance. Models are run via Google Colab's T4 GPU that has certain memory limitations; therefore, it is only possible to use models the size of which does not exceed 8B parameters with quantization up to 8 bit. The final criterion is of critical importance. Most open-source models that fit the first three criteria do not have enough Russian in their training set. This results in empty or unsatisfactory outputs when processing Russian texts. We tested a number of models that fit 3 or 4 of the criteria, but most of them produced poor results. Most of the outputs turned out to be empty when using:

1. Mistral 7B Q4 [8];
2. Solar 10.7B Q4 [12].

The result was satisfactory, but still worse than the model of choice:

1. Saiga Llama 8B Q4 [11];
2. Vikhr 7B Q4 [13];
3. YandexGPT-5-Lite 8B Q4 [14].

The model that produced the best results turned out to be YandexGPT-5-Lite 8B Q8 [15]. Its parameters are described in Tab. 6.

**Table 6.** Parameters of YandexGPT-5-Lite 8B Q8

| Parameter | Value |
| --- | --- |
| Parameter count | 8 billion |
| Base architecture | Llama |
| Quantization | Q8_0 (8 bit) in GGUF format |
| Model size | 8.54 GB |
| Maximum context window | 32K |
| Compatibility | Usable with llama-cpp-python |

It turned out to produce the best output for the second step in the pipeline – keyword analysis and categorization. The key features of this model include a substantial number of Russian texts in the training set and a tokenizer well optimized for working with Russian as well as compatibility with llama.cpp and a size that does not exceed the memory limit of Google Colab's T4 GPU.

## 5. Results

The methodology described above allows us to turn hundreds of unstructured reviews into a concise text report which contains information on positive and negative aspects of visitors' experience. Below is an example of such a report for the year 2020 for one site (Fig. 7). In the first part of the report there are keywords that reflect positive opinions. They are grouped thematically into 14 categories. In the second part the same is done for negative keywords.

Using this methodology, we analyzed reviews on 15 different sites that people wrote on popular websites for tourists. To track changes in visitors' opinion, we analyzed reviews from 2020 to 2025. Several observations can be made based on the results of keywords extraction and sentiment classification. The most frequently mentioned positive aspects are exposition, general impression and emotions and atmosphere: about 2500 keywords extracted from the reviews are related to exposition, about 1000 describe visitors' general impression and about 900 were classified as describing emotions of visitors and atmosphere of the place (Fig. 8).

On the other hand, visitors rarely talk about accessibility and facilities for people with physical disabilities in a positive way (only 51 and 21 keywords respectively in 6 years) (Fig. 8). For example, visitors of museum site "Palaty" mention the following positive aspects (translation from Russian): 'good collection', 'interesting exhibitions', 'valuable paintings' (category 'Exposition'); 'recommend this museum', 'interesting place', 'a perfect place for the whole family' (category 'General impression'); 'cozy atmosphere', 'felt hospitality of the staff', 'comfortable' (category 'Emotions and atmosphere'). Keywords in other categories give more specific details, for example 'friendly staff', 'well-restored and renovated', 'great masterclasses for children'.

Looking at the distribution of keywords among categories (Fig. 8, Fig. 9), we may conclude that when writing a review people firstly remember central aspects of their visit – what they saw
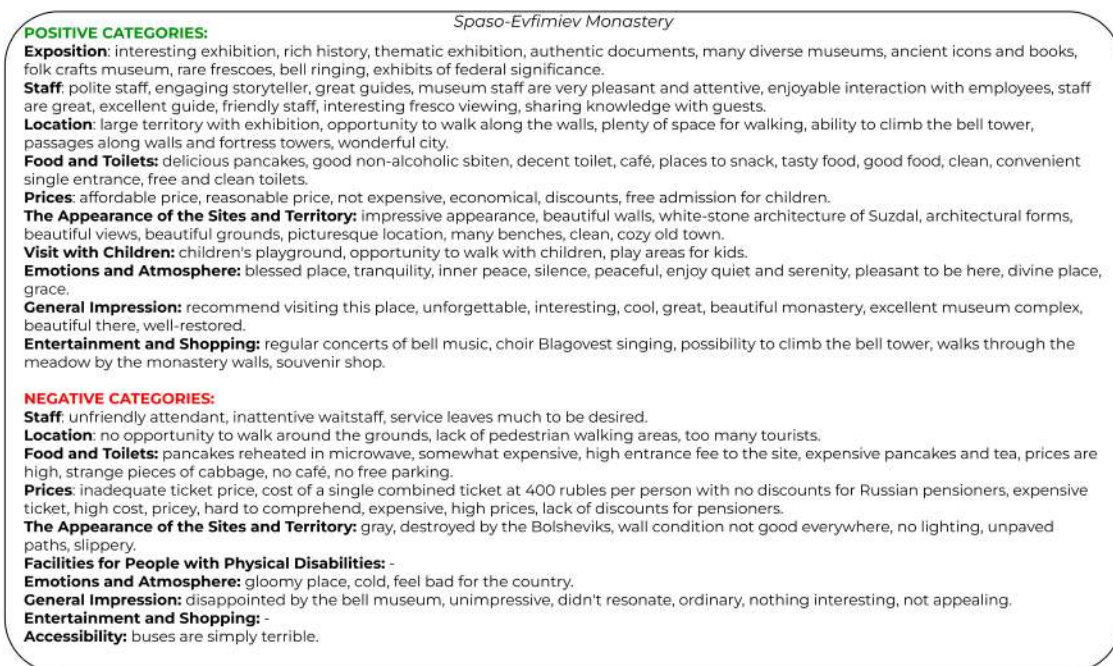
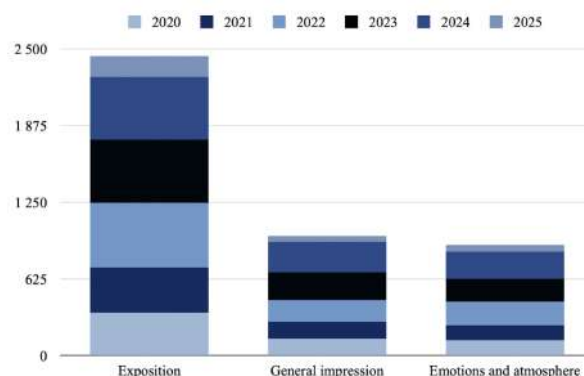**Figure 7.** Report with positive and negative categories



**Figure 8.** The most often mentioned positive categories

(what exhibits), what the place looked like (category 'Appearance of the sites and territory') and how they felt about it. Visitors also do not forget to mention staff if they were friendly and polite. Other details are mentioned optionally.

When it comes to negative opinions, we see that people mention negative things in their reviews much more rarely than positive – the most popular negative categories have 222, 294 and 304 keywords in them (Fig. 10) compared to hundreds and thousands of positive keywords.

Among the most frequent aspects which people mention as negative are 'General impression' and 'Exposition' (similar to the positive ones), however the third category is different which is 'Price' ('a bit pricey', 'not cheap', 'high price is unjustified'). There is one category for which no negative opinions have been found – history and patriotic education (Fig. 11).
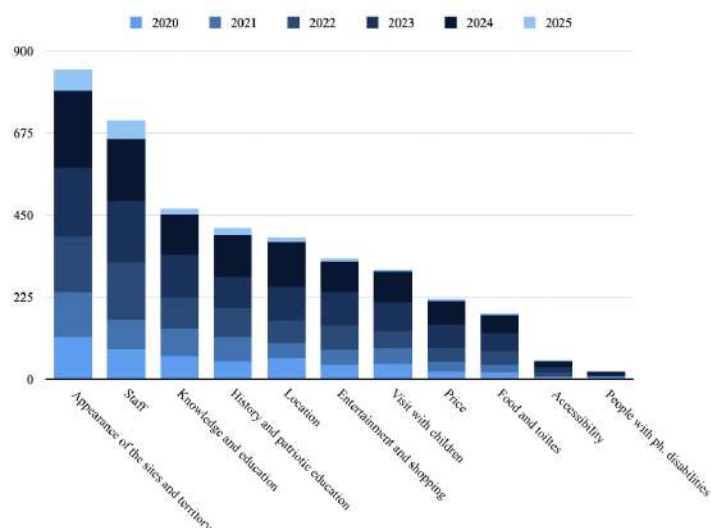
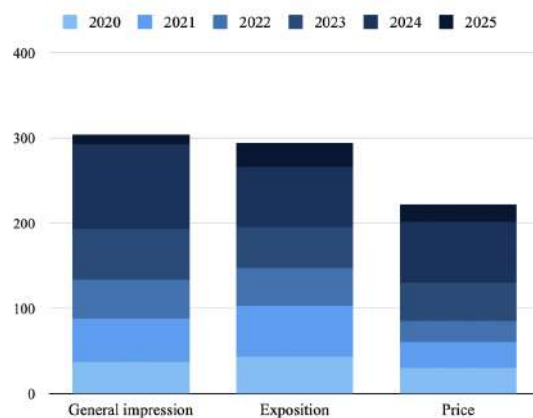**Figure 9.** Amount of keywords in positive categories



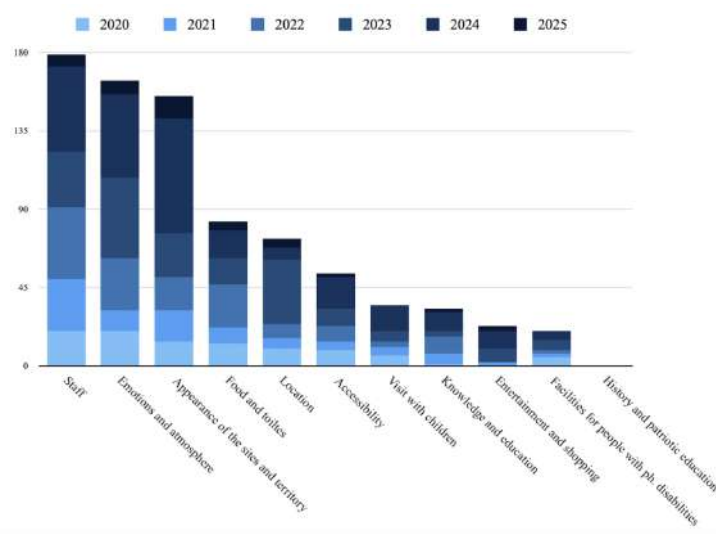**Figure 10.** The most often mentioned negative categories



**Figure 11.** Amount of keywords in negative categories

# 6. Discussion

In this paper we presented the results of collaboration with a museum which needed to optimize visitors' feedback analytics. We have tested a new approach to ABSA of visitors' reviews that leverages large language models and their growing effectiveness in solving many analytical NLP tasks.

For the museum workers it was important to turn thousands of reviews into short, concise yet informative reports on what their visitors like and do not like. To do so, we decomposed the analysis into several steps – expressive keywords extraction, their classification into positive and negative and thematic categorization.

As a result, the information is compressed at two levels: for more detailed analysis the keywords can be used, for an overview there is a report summarizing the information by categories.

The applied approach has its advantages and disadvantages. First advantage is that, among the extracted keywords, there are n-grams of varying lengths: bigrams ('интересная экспозиция' [interesting exhibition]), trigrams ('можно картой оплатить' [accept card payments]), 4-grams ('непродуманная система проверки билетов' [poorly designed ticket checking]), etc. The use of diverse n-grams provides granular insights into visitor preferences, significantly enhancing aspect-level sentiment analysis.

Moreover, the limitation observed in experiments with the dictionary-based method has been overcome – the severe disparity between negative and positive keywords was alleviated. Despite being fewer in number, negative keywords were still extracted for every object. To assess the quality of keyword extraction, we manually annotated a sample of 300 reviews comprising 1331 keywords (in human annotation) and compared them with the extraction performed by the LLM. We considered keywords "missed" if the model (1) failed to add a meaningful keyword, (2) extracted a word/phrase that cannot be considered a keyword for the review or (3) failed to recognize the correct sentiment polarity (e.g., 'до ужаса красив' as negative). The results can be observed in Fig. 12. Automatic keyword extraction turned out to be quite effective with the model "missing" only 92 out of 1331 keywords (6.9%).
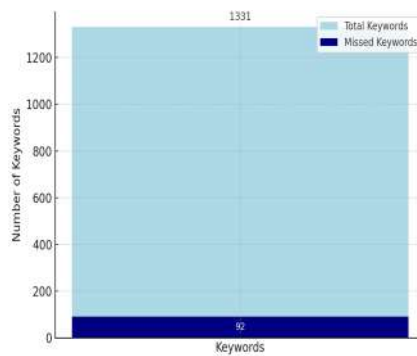


**Figure 12.** Comparison between human and LLM keyword extraction

Thirdly, despite the absence of direct sentiment cues in certain n-grams, the model accurately infers their polarity based on contextual clues and categories them appropriately ('иконы и шкатулки' (positive); 'только учитель прошёл бесплатно' (in context '…а с родителя взяли за билет и за экскурсовода ' (negative)).

Finally, the two-step pipeline eliminated the need for separate topic modeling: the LLM automatically groups keywords into thematic clusters. In most cases, we managed to limit hal-

lucinations for "empty categories" – the model inserts a dash in the thematic category for which no positive or negative keywords were found, rather than inventing non-existent ones.

As for disadvantages, we will focus on several of them.

Firstly, via prompting we did not manage to provide any quantitative data on the frequency of certain topics. The LLM which we have chosen as well as those we tested did not manage to do correct calculations. We tried to instruct the model to count the number of keywords used in each of the categories in absolute and relative figures, but it failed to do so. Proper calculations require access to tools such as code, which is not accessible to smaller models used by us.

The second disadvantage is that we still may encounter model hallucinations for different reasons. One of them is lack of data. For example, when there are 200 reviews for a site in one year and only 10 reviews have some negative aspects mentioned, when classifying keywords the model can make them up to fill in the categories though it is instructed not to do so. Another cause for hallucinations is the comparatively small size of LLMs used, especially the use of quantized models. The smaller the model and the higher the quantization, the less accurate predictions can be made by the model, which results in its failure to strictly follow instructions and keep real data in its context. Hence, manual checking is always needed. We compared the reports for two museum sites and the corresponding tables with keywords to reveal the number of hallucinated keywords that were not present in the table provided to the model as data for analysis (Fig. 13). For the Church of Boris and Gleb, hallucinations made up 7.9% (29 out of 367 keywords), the Suzdal Kremlin report had 5.7% of hallucinated keywords (62 out of 1088 keywords). The percentage of hallucinations is comparatively low, but it only proves the existing problem of large language models fabricating data. Besides, we encountered other minor issues such as repetition, miscategorisation (failure to properly attribute keywords to a category) and grammar mistakes (e.g., agreement between an adjective and a noun – 'подробная путеводитель' [detailed guide]).



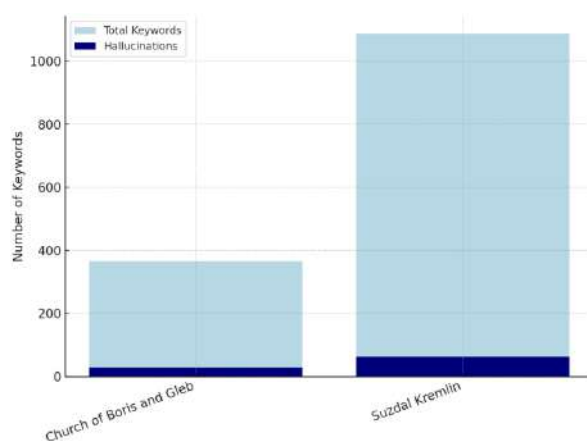**Figure 13.** Number of hallucinations in some reports

The final limitation is that optimal keyword categorization requires pre-defined themes, necessitating either expert annotation or client guidance. Despite providing mandatory/optional category lists, the model over-generated outputs without discrimination. We propose addressing this via advanced prompting techniques like chain-of-thought and few-shot learning.

## Conclusion

This study demonstrates the efficacy of utilizing LLMs for ABSA of museum visitor reviews, addressing the challenges posed by the multi-thematic and open-domain nature of such texts. By implementing a structured pipeline that combines keyword extraction, sentiment classification, and thematic categorization, we successfully transformed unstructured review data into actionable insights without relying on traditional topic modeling techniques. The proposed methodology leverages the contextual understanding capabilities of LLMs to handle diverse n-grams and implicit sentiment cues, achieving a balanced representation of positive and negative aspects across predefined thematic categories. Key advantages include the elimination of sentiment polarity bias, reduced computational dependency on preprocessing, and the ability to generate concise, human-readable reports for end-users. However, limitations such as model hallucinations, quantization constraints, and the need for predefined categories highlight areas for future refinement. Despite these challenges, our approach offers a scalable solution for cultural institutions seeking to optimize visitor experience analytics. Future work could explore the integration of tool-enabled LLMs for quantitative analysis, advanced prompting strategies like chain-of-thought, and domain-specific fine-tuning to further enhance accuracy and reduce manual validation efforts. This research contributes to the growing body of work on LLM applications in NLP and underscores their potential to revolutionize sentiment analysis in non-commercial domains.

## References

1. Blinov, P.D., *et al.*: Research of lexical approach and machine learning methods for sentiment analysis. In: Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference "Dialogue-2013". vol. 12(19), pp. 51–61 (2013)

2. Brauwers, G., Frasincar, F.: A survey on aspect-based sentiment classification. ACM Computing Surveys 54(1), 1–35 (2021). `https://doi.org/10.1145/3503044`

3. Choi, Y., Wiebe, J.: +/-effectwordnet: Sense-level lexicon acquisition for opinion inference. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 1181–1191 (Oct 2014). `https://doi.org/10.3115/v1/D14-1125`

4. Feng, X., Wang, C., Zou, T.T.: Visitor experience of the grand canal national cultural park museum based on sentiment analysis algorithm. SSRG International Journal of Electrical and Electronics Engineering 11(9), 142–150 (2024). `https://doi.org/10.14445/23488379/IJEEE-V11I9P112`

5. Gao, Y., Wang, R., Hou, F.: How to Design Translation Prompts for ChatGPT: An Empirical Study. In: Proceedings of the 6th ACM International Conference on Multimedia in Asia Workshops. Association for Computing Machinery (2024). `https://doi.org/10.1145/3700410.3702123`

6. Gatti, L., Guerini, M., Turchi, M.: SentiWords: Deriving a high precision and high coverage lexicon for sentiment analysis. IEEE Transactions on Affective Computing 7(4), 409–421 (2015). `https://doi.org/10.1109/TAFFC.2015.2476456`

7. Hugging Face: MBARTRuSumGazeta-RuSentiment-RuReviews. `https://huggingface.co/sismetanin/mbart_ru_sum_gazeta-ru-sentiment-rureviews`

8. Hugging Face: Mistral-7B-v0.1-Q4_K_M-GGUF. `https://huggingface.co/3dsabh/Mistral-7B-v0.1-Q4_K_M-GGUF`

9. Hugging Face: RuBERT Conversational Cased Sentiment. `https://huggingface.co/MonoHime/rubert_conversational_cased_sentiment`

10. Hugging Face: RuBERT-Tiny2 Russian Sentiment. `https://huggingface.co/seara/rubert-tiny2-russian-sentiment`

11. Hugging Face: saiga_llama3_8b-Q4_K_M-GGUF. `https://huggingface.co/itlwas/saiga_llama3_8b-Q4_K_M-GGUF`

12. Hugging Face: SOLAR-10.7B-Instruct-v1.0-Q4_K_M-GGUF. `https://huggingface.co/solxxcero/SOLAR-10.7B-Instruct-v1.0-Q4_K_M-GGUF`

13. Hugging Face: Vikhr-7B-instruct_0.2-Q4_K_M-GGUF. `https://huggingface.co/itlwas/Vikhr-7B-instruct_0.2-Q4_K_M-GGUF`

14. Hugging Face: YandexGPT-5-Lite-8B-instruct-GGUF. `https://huggingface.co/yandex/YandexGPT-5-Lite-8B-instruct-GGUF`

15. Hugging Face: YandexGPT-5-Lite-8B-instruct-Q8_0-GGUF. `https://huggingface.co/BoloniniD/YandexGPT-5-Lite-8B-instruct-Q8_0-GGUF`

16. Koltsova, O.Y., Alexeeva, S.V., Kolcov, S.N.: An opinion word lexicon and a training dataset for russian sentiment analysis of social media. In: Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference "Dialogue-2016". vol. 15(22), pp. 277–287 (2016). `https://doi.org/10.5281/zenodo.4084953`

17. Kulagin, D.: Russian word sentiment polarity dictionary: a publicly available dataset. Poster, Artificial Intelligence and Natural Language (AINL 2019) (2019). `https://doi.org/10.28995/2075-7182-2021-20-1106-1119`

18. Loukachevitch, N., Levchik, A.: Creating a General Russian Sentiment Lexicon. In: Proceedings of the Language Resources and Evaluation Conference (LREC 2016). pp. 1171–1176. European Language Resources Association (ELRA) (2016)

19. Loukachevitch, N., Tkachenko, N., Lapanitsyna, A., *et al.*: RuOpinionNE-2024: Extraction of Opinion Tuples from Russian News Texts (2025)

20. Pang, B., Lee, L.: Opinion Mining and Sentiment Analysis. Foundations and Trends in Information Retrieval 2(1–2), 1–135 (2008). `https://doi.org/10.1561/1500000011`

21. Pang, B., Lee, L., Vaithyanathan, S.: Thumbs Up? Sentiment Classification Using Machine Learning Techniques. In: Mooney, R.J. (ed.) Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 79–86. Association for Computational Linguistics, Philadelphia, USA (2002). `https://doi.org/10.3115/1118693.1118704`

22. Qi, P., Zhang, Y., Zhang, Y., *et al.*: Stanza: A Python natural language processing toolkit for many human languages. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations. pp. 101–108. Association for Computational Linguistics (2020). `https://doi.org/10.18653/v1/2020.acl-demos.14`

23. Qin, L., *et al.*: Large Language Models Meet NLP: A Survey. Frontiers of Computer Science (FCS) (2024). `https://doi.org/10.1007/s11704-025-50472-3`

24. Sheremetyeva, S.: An efficient patent keyword extractor as translation resource. In: Proceedings of the MT Summit XII: Third Workshop on Patent Translation. pp. 25–32. Ottawa, Canada (2009)

25. Tabularisai, Gyamfi, S., Borisov, V., Schreiber, R.H.: Multilingual-sentiment-analysis (revision 69afb83) (2025). `https://doi.org/10.57967/hf/5968`

26. Vatolin, A.: Structured sentiment analysis with large language models: A winning solution for RuOpinionNE-2024. In: Proceedings of the International Conference on Computational Linguistics and Intellectual Technologies (Dialogue 2025). FRC CSC RAS, Moscow, Russia (2025)

27. Wang, Z., Xie, Q., Feng, Y., *et al.*: Is chatGPT a good sentiment analyzer? In: First Conference on Language Modeling (2024), `https://openreview.net/forum?id=mUlLf50Y6H`

28. Xu, Q., Shih, J.Y.: Applying text mining techniques for sentiment analysis of museum visitor reviews. In: 2024 IEEE 4th International Conference on Electronic Communications, Internet of Things and Big Data (ICEIB). pp. 270–274. Taipei, Taiwan (2024). `https://doi.org/10.1109/ICEIB61477.2024.10602556`

29. Yadav, A., Vishwakarma, D.: Sentiment analysis using deep learning architectures: A review. Artificial Intelligence Review 53(6), 4335–4385 (2020). `https://doi.org/10.1007/s10462-019-09794-5`

30. Zhang, T., Irsan, I., Thung, F., *et al.*: Revisiting sentiment analysis for software engineering in the era of large language models. ACM Transactions on Software Engineering and Methodology 34(3), 1–30 (2025). `https://doi.org/10.1145/3697009`

31. Água, M., Antonio, N., Carrasco, M.P., *et al.*: Large language models powered aspect-based sentiment analysis for enhanced customer insights. Tourism & Management Studies 21(1), 1–19 (2025). `https://doi.org/10.18089/tms.20250101`

# Neurosymbolic Approach to Processing of Educational Texts for Educational Standard Compliance Analysis

*Nikolai A. Prokopyev*[1] (iD), *Marina I. Solnyshkina*[1] (iD), *Valery D. Solovyev*[1] (iD)

This article presents a neurosymbolic approach for analyzing the alignment between textbook content and educational standards. The study addresses the problem of assessing terminological coherence by evaluating a corpus of textbooks against the Russian Federal State Educational Standard. We employ a hybrid methodology combining classical symbolic NLP methods for topic modeling (keyword extraction and term alignment) with qualitative analysis and use of modern large language models for items not found algorithmically. The experimental results on a corpus of 5 textbooks on Physics for the 7th grade and corresponding educational standard indicate a mean coverage of standard topics of 71% across all textbooks with use of the symbolic methods. Application of large language model (ChatGPT 5) for the qualitative analysis recovered 51% keywords initially missed by the abovementioned methods. The findings are relevant for researchers in educational linguistics, computational linguistics, curriculum developers, and textbook authors. The proposed pipeline offers a scalable tool for automating analysis of educational content compliance, reducing the workload for manual expert assessment. This work contributes to the development of AI-assisted methodologies in educational standard alignment and textbook quality control.

*Keywords: topic modeling, keyword extraction, symbolic NLP, large language model, textbook analysis.*

## Introduction

Alignment between educational standards and textbooks content as the degree of coherence between curricula and textbooks content has been an area of numerous disputes [26], and as such extensively studied around the world [4, 13]. Modern natural language processing (NLP) paradigm provides numerous approaches and powerful toolkits to measure this alignment: recent advances in large language models (LLMs) resulted in significant changes in the area, and LLMs enable considerate assistance in educational content assessment [1, 25] thus reducing the workload of academics and test developers.

Specifics and complexity of solving this problem is related to identifying the range of linguistic variability in texts, dynamism of modern discourse and active expansion of nonverbal signs into academic texts as well as growing number of nonlinear, polycode texts. All the above constitute the foundation of ongoing research in alignment between textbooks content and national standards.

Although the Russian textbook language quality and its compliance with national standards have been lately addressed by Russian and foreign researchers [16, 23], to the best of our knowledge, there is no research which explicitly uses both symbolic NLP methods and LLM evaluation tools to assess textbook language alignment with national educational standards. Thus the current research is aimed at evaluating terminological coherence between the educational content of five Physics textbooks and the Federal State Educational Standard.

The main goal of this research is to experimentally verify the algorithm for educational standard compliance analysis of educational texts. This algorithm can be used as a support tool for authors of textbooks and official experts by giving them the preliminary results for furhter investigation. Thus, the research questions we address are as follows:

---

[1]Kazan Federal University, Kazan, Russian Federation

1. Application of classical symbolic NLP methods to keyword extraction;
2. Application of NLP methods based on neural networks, particularly use of LLM for additional keyword extraction;
3. Combined use of these methods to solve the problem of educational standard compliance analysis for textbooks.

The article is structured as follows. Section 1 provides an overview of the pipeline for problem of educational standard compliance analysis. Section 2 presents the related works analysis of NLP technologies used in various relevant tasks. Section 3 describes the results of this study. Conclusion presents the main findings and directions for future developments.

# 1. Methodology

## 1.1. Problem of Educational Standard Compliance Analysis

Before we move on to clarified formulation of the problem of educational standard compliance analysis, it is necessary to describe the most important elements of the Federal State Educational Standard structure. These are the following elements:

- A set of topics on the school subject for one school class. Each topic has a description in free form, it determines which concepts from the topic should be disclosed in a textbook.
- A set of terms on the school subject for one school class. These terms should also be disclosed in a textbook.

The problem of educational standard compliance analysis then is defined as follows: to check whether a textbook aligns with topics and terms from the standard, also assessing to what extent it occurs.

To solve this problem, we propose to reduce it to the topic modeling problem partially, within which a topic description in the standard (standard topic) should be represented as a set of keywords. The set of terms in the standard (standard terms) can be considered a separate topic, isolated in its pragmatic significance for this problem.

Thus, we propose to use a hybrid neurosymbolic approach to this topic modeling problem. The symbolic part of the approach consists in usage of tokenization, lemmatization, n-gram retrieval and bag of words (BOW) methods. These methods are used for the primary algorithmic analysis of the textbook to search for keywords for each standard topic and to search for standard terms. The neural network part of the approach involves the use of LLM for subsequent qualitative analysis of keywords and terms not found algorithmically.

The general scheme of the solution pipeline is shown in Fig. 1. A detailed description of the pipeline stages is presented in the next subsections.

## 1.2. Preprocessing of Textbooks and Standards

Preprocessing of textbooks is an important pipeline stage, necessary for further analysis. This stage is fully automated, with the following algorithm:

1. Compilation of a textbook metadata file containing: a list of textbooks in the source dataset, their metadata (ID, title, school subject, school class), paths to files with different forms of textbook text representation.
2. Extraction of texts from the source textbook files and saving them in TXT format.
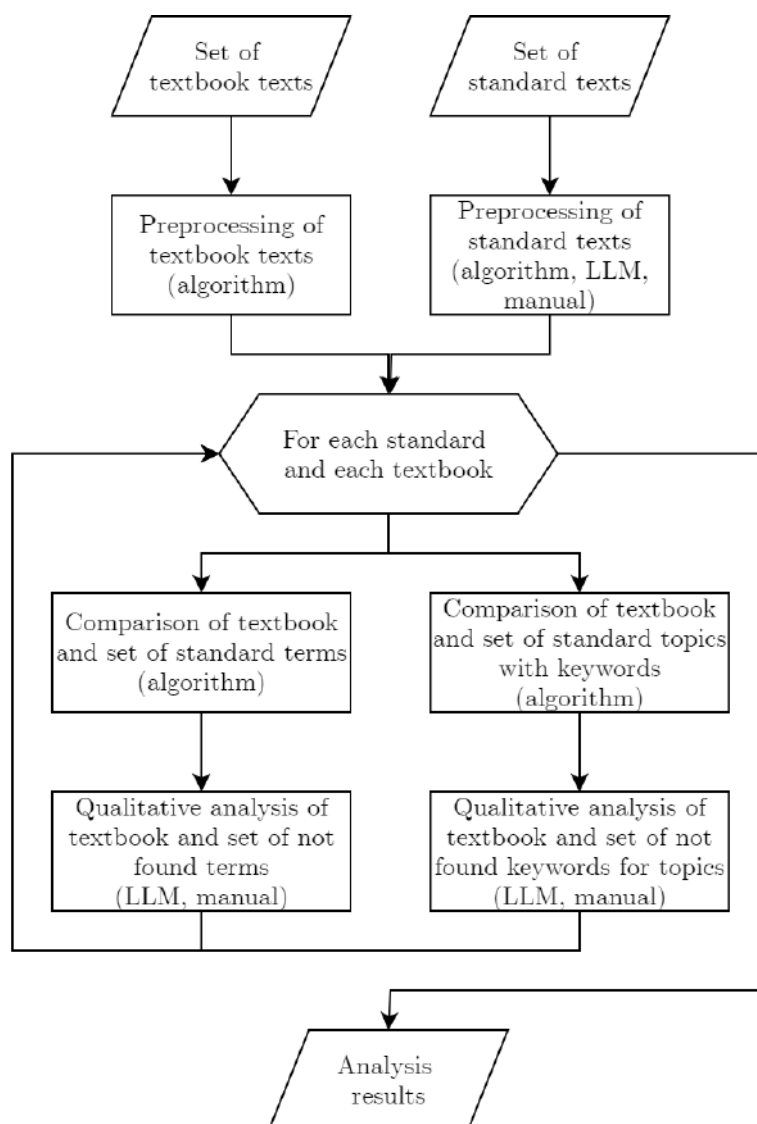
**Figure 1.** General scheme of textbook analysis stages

3. Tokenization of texts, which includes the removal of punctuation, stop words, and grammatical lemmatization (normalization) of tokens using a specific configuration. The resulting token set is saved in JSON format.

Preprocessing of standard texts is a more complex task, since the description of each standard topic is a free-form text, for automatic processing of which a semantic analyzer is required. Thus, while this stage is not fully automated, we propose to use LLM to extract keywords and terms as it can be an effective solution for semantic analysis as shown in Section 2. Sequence of actions:

1. Algorithmic compilation of a standards metadata file containing: a list of standards, their metadata (school subject, school class, level – basic or advanced), paths to files with different forms of standard representation.

2. Manual split of the standards text into topic descriptions and term descriptions.

3. Extraction of topic keywords and terms from the corresponding descriptions using LLM with a prepared prompt.

4. Algorithmic compilation of normalized n-grams of keywords and terms. Since keywords and terms can consist of several words in different word forms, n-grams should be compiled for

further analysis, and their normalization should be performed using the same configuration that was used to normalize textbook tokens. The resulting sets of keywords and a set of terms are saved in JSON format.

## 1.3. Primary Analysis: Comparison of Textbooks and Standards

For the primary analysis of textbooks for compliance with the educational standard, an automatic algorithm is used to align the tokenized textbook and sets of standard topic keywords or standard terms. The algorithm is as follows:

1. Compilation of ID-dictionary for the corpus of tokenized textbooks, taking into account the automatic collection of tokens into n-grams.
2. Compilation of a textbook BOW using the ID-dictionary, containing IDs of n-grams and their number of instances in the textbook.
3. For each standard topic, a BOW of the corresponding keywords set is compiled, after which a search for matching IDs of keywords in the textbook BOW is performed. The following are output: keyword, number of instances, frequency.

$$Frequency = \frac{number\ of\ instances}{textbook\ length\ in\ tokens}$$

4. For each standard topic, a separate search is performed among the not found keywords, those that had too small number of instances to be collected into an n-gram.
5. For each standard topic, the keywords found and not found by the algorithm are displayed, and the coverage is calculated.

$$Coverage = \frac{number\ of\ found\ keywords\ in\ textbook}{total\ number\ of\ keywords\ in\ topic}$$

It should be noted that classical metrics for keyword extractions, such as accuracy, precision and recall, are not suitable in this case, as they are metrics for supervised classification task demanding availability of a labeled data, while this algorithm is unsupervised.

6. Average coverage of the standard topics in the textbook is calculated.

The automatic algorithm for comparing the tokenized textbook and the previously obtained set of standard terms is performed in a similar way, while set of standard terms is accepted as a separate topic with a set of corresponding keywords.

## 1.4. Qualitative Analysis of Undetected Keywords and Terms

Since classical NLP methods do not take into account synonyms and semantically close descriptions of standard topics and terms, a qualitative analysis of the keywords and terms not found in each textbook is necessary.

Obviously, to solve this problem, manual expert assessment of each textbook can be used, but it is a labor-intensive task in the conditions of a large volume of textbooks. Thus, we propose to use LLM with a prepared prompt to check each set of keywords and terms not found in the primary analysis for each textbook.

However, the output of LLM also requires manual verification, so the set of terms additionally found using LLM is output separately. Recovery metric is calculated as:

$$Recovery = \frac{number\ of\ keywords\ found\ by\ LLM}{number\ of\ keywords\ not found\ by\ algorithm}.$$

## 1.5. Source Data and Technologies

For the experimental verification of the proposed approach, a set of source data was used: the educational standard for Physics, 7th school class, basic level, and a set of textbooks. Table 1 presents the topics of this standard with English translation and their IDs used further. 5 textbooks on Physics for the 7th class of different publication years in WORD document format were considered. Table 2 presents the bibliographic data of these textbooks and their IDs used further.

**Table 1.** Standard topics

| ID | Topic |
|----|-------|
| Topic 1 | Физика и ее роль в познании окружающего мира |
| | Physics and its role in understanding the world around us |
| Topic 2 | Первоначальные сведения о строении вещества |
| | Initial information about structure of matter |
| Topic 3 | Движение и взаимодействие тел |
| | Movement and interaction of bodies |
| Topic 4 | Давление твердых тел, жидкостей и газов |
| | Pressure of solids, liquids and gases |
| Topic 5 | Работа и мощность. Энергия |
| | Work and power. Energy |

**Table 2.** Textbook bibliography

| ID | Textbook |
|----|----------|
| 81412 | Генденштейн Л. Э. Физика. 7 класс. В 2 ч. Ч. 1 : учебник для общеобразовательных учреждений / Л. Э. Генденштейн, А. В. Кайдалов ; под ред. В. А. Орлова, И. И. Ройзена. – 3-е изд., испр. – М.: Мнемозина, 2012. - 255 с. : ил. ISBN 978-5-346-02160-5 |
| 31329 | Громов С. В. Физика: Учеб. для 7 кл. общеобразоват. учреждений/ С. В. Громов, Н. А. Родина.- 4-е изд.- М.: Просвещение, 2002.- 158 с.: ил.- ISBN 5-09-011495-1 |
| 28878 | Физика. 7 класс : учеб. для общеобразоват. организаций / О.Ф. Кабардин. - 3-е изд. - М.: Просвещение, 2014. - 176 с.: с ил. |
| 21915 | Физика. 7 кл. : учеб. для общеобразоват. учреждений / А.В. Пёрышкин. - 2-е изд., стереотип. - М.: Дрофа, 2013. - 221 с.: с ил. |
| 26802 | Физика. 7 класс : учебник / Н.С. Пурышева, Н.Е. Важеевская. - 2-е изд., стереотип. - М.: Дрофа, 2013. - 224 с.: с ил. |

For implementation, the Python programming language, PyCharm development environment, and Jupyter were used. `textract` library was used to extract textbook texts. `gensim` library was used to apply symbolic NLP methods: ID-dictionary compilation, bigram and trigram extraction (threshold = 1), BOW compilation. Additionally, less frequent n-gram candidates for n > 3 were extracted using custom algorithm. Tokenization was performed using `razdel`

library, stop words were removed using `stopwords` library, and lemmatization was performed using `pymorphy2`.

ChatGPT 5 was used as the main LLM in experimental setup. The prompt for preprocessing the standard text is presented in Fig. 2 with English translation. The prompt for additional search of not found terms and keywords is presented in Fig. 3.

The best prompts obtained as a result of prompt engineering are shown. It was peculiar that small changes in the prompt in the conditions of processing an uploaded textbook file greatly affected the output of LLM, and whether it would take into account the attached textbook file or not.

## 2. Related Works

LLMs, being language models, are successfully applied in a wide range of problems related to information retrieval in texts, such as keyword extraction [11].

Application of the classical BERT model is described in [10], where KeyBERT is presented. The basic idea is that embeddings (vector representations) of the whole text and individual words are built, which are then compared by the degree of similarity. Its further improvement is proposed in [20].

Classical statistical methods also continue to be used. A general overview of statistical methods for keyword extraction can be found in [18].

---

Я извлекаю ключевые слова из стандарта ФГОС по {предмет} для среднего образования. Изначально они мне даны в виде описания на естественном языке. Преобразования, которые я применяю, следующие:

- Объединение нескольких слов в n-грамму ключевого слова;
- Извлечение нескольких n-грамм из их сокращенного описания;
- Написание собственных n-грамм из описания в виде предложения;
- Извлечение из n-грамм более обобщающих ключевых слов;
- Фильтрация слов, не являющихся ключевыми словами по предмету.

Мне необходимо автоматизировать ручную работу. Далее я буду давать тебе описания, а ты преобразуй их в списки n-грамм. Отвечай кратко, списком в формате Python, без пояснений.

---

I extract keywords from the State Educational Standard for {subject} for School Education. Initially, they are given to me as a description in natural language. The transformations I apply are following:

- Combining several words into an n-gram of a keyword;
- Extracting several n-grams from their abbreviated description;
- Writing own n-grams from a description in form of a sentence;
- Extracting more general keywords from n-grams;
- Filtering words that are not keywords of the subject.

I need to automate manual work. Next, I will give you the descriptions and you transform them into lists of n-grams. Answer briefly, with a list in Python format, without explanations.

**Figure 2.** Prompt for text processing of the standard

---

Тебе дан текст учебника по {предмет} для среднего образования в виде прикрепленного файла. Я делаю анализ наличия ключевых слов по темам. Для одной темы мой алгоритм выявил следующие не найденные ключевые слова (приведены в виде нормализованных n-грамм):

{ключевые_слова}

Попытайся найти эти ключевые слова в предоставленном учебнике в виде их синонимов или близких к ним описаний и тем. В ответе дай только список найденных ключевых слов в том виде, в котором они даны в изначальном списке.

---

You are given the text of a {subject} textbook for secondary education in the attached file. I do a keyword analysis by topic. For one topic, my algorithm identifies the following not found keywords (given as normalized n-grams):

{keyword_set}

Try to find these keywords in the provided textbook as their synonyms or close descriptions and topics. The answer only contains a list of found keywords in form in which they are given in the original list.

**Figure 3.** Prompt for search of not found terms and keywords

A combined approach with simultaneous application of several methods, including Transformer LLM, is proposed in [27]. Unfortunately, this work does not provide data allowing us to evaluate the effect of applying modern LLMs. Experimentally confirmed advantages of the combined approach are described in [30]. In this paper, LLM is combined with knowledge graphs, in which information is represented as triplets. Medical texts are processed and it is shown that the combined use of these two approaches provides a better result than they can give separately. The idea of combining LLM with knowledge graphs seems to be very promising.

In the task of extracting keywords from abstracts of Russian-language scientific articles using the BERTScore metric [33], the average result of three LLMs (Saiga, Mistral, and Vikhr) in zero-shot mode was 76.05, when soft fine-tuned on three random examples was 77.28. The average result of two classical statistical methods, YAKE and RuTermExtract, was 72.54. The average result of two fine-tuned neural networks, mT5 and mBART, was 77.37. Thus, LLMs demonstrate a better result than classical statistical methods, and are inferior to specialized fine-tuned neural networks, but not by much.

These results are further demonstrated in [8], fine-tuned mT5 model is compared to TopicRank, YAKE, RuTermExtract, KeyBERT using F-measure, ROUGE-1 and BERTScore metrics. According to this article, mT5 and RuTermExtract shows the highest performance in terms of the BERTScore metric (76.89 and 75.80), where mT5 demonstrates better results when generating keywords not presented in the source text.

Regarding the studies examining end-to-end solutions using only LLM for keyword extraction, given the different tasks, datasets, and metrics used, a direct comparison of results is not possible. It can only be noted that several studies [9, 24], like ours, have noted the high potential of LLM even without additional fine-tuning for keyword-related tasks.

A significant number of publications are devoted to the application of LLM for extracting critical terms in various subject areas to solve specific problems. In [7], extraction of keywords from electronic medical records to create a database of oncological diseases is described. In [2],

application of the PhenoBERT system and LLM for extracting phenotypes from clinical records is considered and it is shown that LLM can find information missed by experts.

In [15], LLM-based framework for extracting terms from E-commerce platform texts is proposed. Specific applications are given and the effectiveness of this approach is experimentally demonstrated. In article [6], LLM is applied to the analysis of historical documents. In article [19], LLM is applied to extract objects and their properties from agricultural texts in order to obtain information about pests. It is shown that LLM can achieve better results than conventional methods.

Relatively few publications are available on application of LLM for term extraction in the social sciences and humanities. In [5], the lack of standardized datasets for applying LLM is noted and a dataset for political texts is presented.

It should be noted that not all works demonstrate the advantages of LLM in this problematics. Thus, in [32] it is shown that they successfully perform in extraction of terms only from simple sentences, while statistical methods are preferable in more complex situations. An important study is [3], which compares the results shown by 3 models: Llama2-7B, GPT-3.5, and Falcon-7B. This article discusses various factors: hallucinations, prompt quality, dependence on the subject area.

An important area of research closely related to keyword extraction is topic modeling. It has been extensively studied and applied in many subject areas, including in education for analyzing the structure of textbooks [21, 23]. A comprehensive overview of topic modeling can be found in [29].

In [17], LLM was proposed to be used to extract topics instead of traditional methods such as Latent Dirichlet Allocation (LDA). It was noted that LDA and similar methods do not take semantics into account, and in this regard, LLM has a significant advantage. The article shows that this is indeed the case and concludes that LLM can serve as an effective means of topic modeling. In [31], it was shown that LLM can be used to assess the relevance of topics instead of human experts.

LLMs have been used to extract topics in a number of subject areas. Thus, in [28], the BERTopic system, created on the basis of the early LLM BERT, was successfully applied to analyze psychotherapeutic texts. Its further development is presented in [12]. On a dataset of news texts, the combined use of modern LLM and standard topic modeling techniques improved topic coherence by about 10% compared to standard methods. In [14], it is proposed to combine LLM and cluster analysis methods. Thus, the general current trend is to combine LLM with standard techniques, and all studies demonstrate an improvement in quality of the extracted topics. LLM is becoming a new standard in topic modeling.

It is worth noting that the quality of textbooks is determined not only by their compliance with standards, but also by text complexity, which must be accessible to students. An overview of text complexity issues can be found in [22].

## 3. Results and Discussion

All textbooks were processed in the proposed pipeline. In terms of topic analysis, coverage of each topic for each textbook was calculated as well as mean coverage of all standard topics. Keywords not found by the primary analysis were processed with LLM, and were additionally found in all textbooks. Figure 5 shows an example of the result card for one textbook and one

topic with all these data. The same steps were done for terms analysis. Figure 6 shows the corresponding result card for one book.

In summary, mean coverage of each topic and mean coverage of terms on all textbooks were calculated, as well as mean LLM recovery. Table 3 shows these statistics. Mean coverage among all topics is 71%. Mean LLM recovery is 51%. Figure 4 shows distribution of topic coverage in each textbook. This output data can be further used by textbook authors to revise their works and by official experts for comparative analysis of compliance to the educational standard.
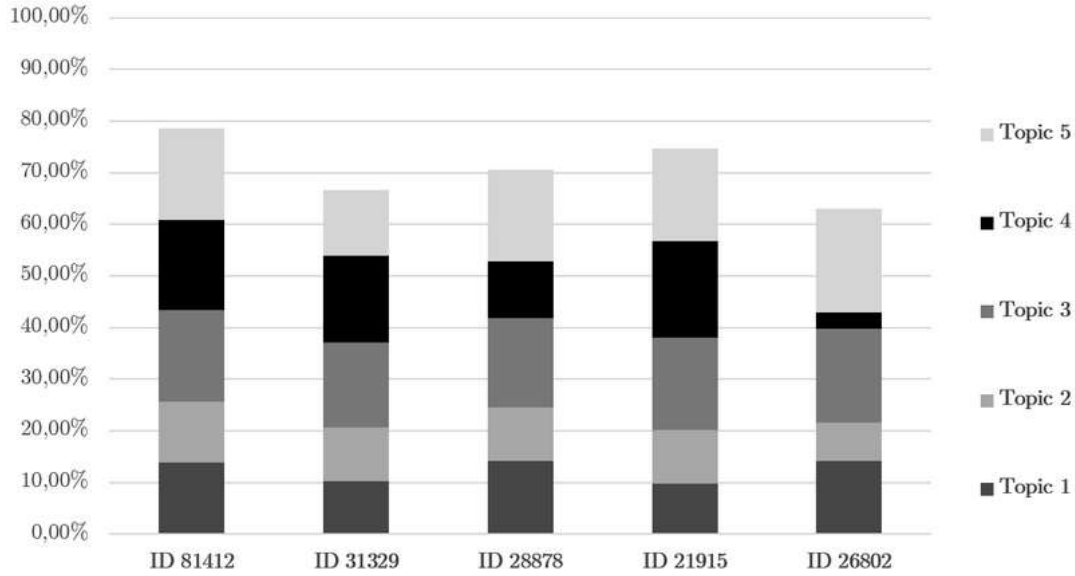


**Figure 4.** Topic coverage among all textbooks

**Table 3.** Results statistics

| Topic | Number of keywords | Mean coverage | Mean LLM recovery |
|-------|--------------------|---------------|-------------------|
| Topic 1 | 35 | 62% | 56% |
| Topic 2 | 27 | 50% | 52% |
| Topic 3 | 35 | 88% | 60% |
| Topic 4 | 31 | 67% | 24% |
| Topic 5 | 19 | 86% | 64% |
| Terms | 88 | 81% | 70% |

The first question for discussion is why the symbolic part of the algorithms does not find all the standard keywords and terms, even when they are present in textbook. The main weakness of the symbolic approach is its reliance on exact symbolic representation of words, without consideration of synonymy or semantics. This is why the preprocessing is an important step, where some nested keyword and term representations from the standard are expanded. For example, excerpt *деформация (упругая, пластическая) = deformation (elastic, plastic)* from the standard is expanded into terms *деформация, упругий_деформация, пластический_деформация = deformation, elastic_deformation, plastic_deformation*.

Furthermore, all of the keywords and terms not found by the symbolic approach, but recovered using LLM, are presented in textbooks in forms of synonyms or descriptions. First example, the term *правило_равновесие = equilibrium_rule* was not found, but its synonym *условие равновесия = equilibrium condition* is present in the textbook ID 81412, so it was recovered by

LLM (see Fig. 6). Second example, the term *агрегатный_ состояние = aggregate_ state* was not found, but the textbook ID 81412 contains an excerpt with the description of solid, liquid and gaseous states which correspond to this term, so it was recovered by LLM.

The second question for discussion concerns error analysis of LLM-based keyword recovery. Let us define our null hypothesis about every keyword not found by the symbolic part of the algorithm as not present in a given textbook. Then every error of LLM output is either type I error (keyword recovered, but not present in the textbook) or type II error (keyword not found by LLM, but present in the textbook in some form). Then type I errors are divided into type $I_a$ errors (keyword recovered from false synonym or similar wording), type $I_b$ errors (keyword recovered from similar description) and type $I_c$ errors (keyword recovered as hallucination). Type II errors are divided into type $II_a$ errors (keyword is present as synonym) and type $II_b$ errors (description of keyword is present).

**Table 4.** Error statistics for terms recovery

| Error | ID 81412 | ID 31329 | ID 28878 | ID 21915 | ID 26802 | Average |
|-------|----------|----------|----------|----------|----------|---------|
| Type I | 0% | 19% | 0% | 7% | 22% | 10% |
| Type $I_a$ | 0% | 13% | 0% | 7% | 19% | 8% |
| Type $I_b$ | 0% | 6% | 0% | 0% | 4% | 2% |
| Type $I_c$ | 0% | 0% | 0% | 0% | 0% | 0% |
| Type II | 30% | 19% | 17% | 7% | 11% | 17% |
| Type $II_a$ | 10% | 6% | 6% | 7% | 0% | 6% |
| Type $II_b$ | 20% | 13% | 11% | 0% | 11% | 11% |
| All types | 30% | 38% | 17% | 14% | 33% | 27% |

Let us consider error analysis for terms recovery (Tab. 4), as they, by design of the educational standard, represent keywords from all of the topics. The table shows the ratio of LLM error number to number of terms not found by the symbolic part of the algorithm which should be processed by LLM.

As can be seen in these results, LLM produces less type I errors than type II errors on average, with the percentage of type I errors sufficiently low, that we can draw a considerable level of trust for experts and textbook authors to the output of LLM-based keyword recovery. No hallucinations occurred, for every false keyword recovery there can be given an explanation whether it was due to similar wording in textbook (dominant number of errors), or due to similar description. Additional analysis of hallucination occurrence was done in the form of introduction of trap words. We added to prompt some terms from physics topics of higher education level, such as *relativity theory, electrodynamics*, and some terms from completely different topics, namely literature and biology. No new hallucinations occurred after this addition.

As for type II errors, there are dominance of not found descriptions of keywords in textbook, which can be explained as attention mechanism in the LLM working in large context when the keyword description excerpts in the textbook are distributed sparsely. While the number of type II errors is higher, the full list of not found keywords is due for manual expert check in the educational standard compliance analysis regardless, therefore the level of trust in this case is out of question.

Another note on the results considers textbooks ID 31329 and ID 26802. These textbooks had a higher number of keywords not found by the symbolic part of the algorithm, and there are

more errors than in other textbooks. We can conclude that the higher number of keywords the LLM have to find in the textbook, the less trust to its output should there be. This note raises a preliminary question about the effectiveness of end-to-end LLM-only solution to the problem of education standard compliance, without the symbolic preprocessing which can significantly reduce the number of keywords in the LLM input in an explainable and trusted manner.

It should be noted that in the LLM-based keyword recovery step we used the extracted text versions of the textbooks (.txt format) rather than source forms of Word documents (.docx) as the latter approach produced considerably more errors in the result.

## Conclusion

The presented research tests application of classical symbolic NLP methods and LLM to keyword extraction. The practical task addressed and employed is the terminological alignment between educational standards and textbooks content. Within the neurosymbolic approach employed, we apply a hybrid methodology and combine symbolic methods for topic modeling of 5 Russian textbooks on Physics for 7th-graders. The pipeline is constituted of (1) keyword extraction with tokenization, lemmatization, n-gram retrieval and bag of words methods for topic modeling followed by (2) applying ChatGPT and qualitative analysis for items not found algorithmically. Having applied the symbolic methods we evaluated a 71% mean coverage of standard topics across all textbooks. Employment of LLMs resulted in recovery of 51% keywords initially missed by the abovementioned methods.

The research results show that hybrid neurosymbolic approach performs adequately good for the task of educational standard compliance analysis. The symbolic part of our algorithm finds most of the keywords and terms from the standard, with its output being explainable without manual post-check. Usage of LLM allows for additional recovery of keywords and terms as their synonyms and descriptions, with mostly reliable output.

As the textbooks were published under different Federal State Educational Standards over the period 2002–2019, the practical results are only partially applicable for educational and policy making purposes, although offer a scalable tool to automate analysis of content compliance of two sources, thus reducing the manual workload. Another limitation of the study is related to the application of ChatGPT 5 only. Given the rapid progress of LLMs and obsolescence of earlier models such as GPT2, in this study we used the most modern model, i.e., ChatGPT 5, and left comparison of the results to be achieved by other LLMs with those obtained in this work as a baseline for future studies. Another direction of future study is the experimental verification of an end-to-end solution for the educational standard compliance analysis based only on LLM with soft fine-tuning.

## Acknowledgements

**Textbook**: ID 81412

**Standard**: Physics 7th grade Basic

**Mean coverage of standard topics**: 79%

---

**Standard topic**: Topic 1

**Keywords in topic**: 35

**Coverage**: 69%

**Found keywords** (24, showing first 10):

|    | Keyword (Ru) | Keyword (En) | Instances | Frequency |
|----|--------------|--------------|-----------|-----------|
| 1  | гипотеза | hypothesis | 11 | 0.000193 |
| 2  | измерение | measurement | 19 | 0.000333 |
| 3  | модель | model | 8 | 0.000140 |
| 4  | наблюдение | observation | 14 | 0.000245 |
| 5  | наука | science | 11 | 0.000193 |
| 6  | природа | nature | 27 | 0.000473 |
| 7  | термометр | thermometer | 4 | 0.000070 |
| 8  | физика | physics | 11 | 0.000193 |
| 9  | физический_величина | physical_quantity | 21 | 0.000368 |
| 10 | физический_явление | physical_phenomenon | 6 | 0.000105 |

**Not found keywords** (11):

|    | Keyword (Ru) | Keyword (En) |
|----|--------------|--------------|
| 1  | датчик_температура | temperature_sensor |
| 2  | физический_превращение | physical_transformation |
| 3  | постановка_научный_вопрос | formulation_scientific_question |
| 4  | жидкостный_термометр | liquid_thermometer |
| 5  | естественнонаучный_метод_познание | natural_scientific_method_cognition |
| 6  | научный_вопрос | scientific_question |
| 7  | физический_прибор | physical_device |
| 8  | объяснение_наблюдать_явление | explanation_observed_phenomenon |
| 9  | выдвижение_гипотеза | hypothesis_proposal |
| 10 | описание_физический_явление | description_physical_phenomenon |
| 11 | химический_превращение | chemical_transformation |

**LLM-found keywords** (4):

|    | Keyword (Ru) | Keyword (En) |
|----|--------------|--------------|
| 1  | физический_прибор | physical_device |
| 2  | выдвижение_гипотеза | hypothesis_proposal |
| 3  | описание_физический_явление | description_physical_phenomenon |
| 4  | объяснение_наблюдать_явление | explanation_observed_phenomenon |

**Figure 5.** Topic analysis results representation

**Textbook**: ID 81412

**Standard**: Physics 7th grade Basic

**Terms in standard**: 88

**Coverage**: 89%

**Found terms** (78, showing first 10):

|  | Term (Ru) | Term (En) | Instances | Frequency |
|---|---|---|---|---|
| 1 | атмосферный_давление | atmospheric_pressure | 48 | 0.000841 |
| 2 | атом | atom | 32 | 0.000561 |
| 3 | вес | weight | 51 | 0.000894 |
| 4 | весы | scales | 8 | 0.000140 |
| 5 | взаимодействие_тело | body_interaction | 1 | 0.000018 |
| 6 | время | time | 43 | 0.000754 |
| 7 | гипотеза | hypothesis | 11 | 0.000193 |
| 8 | давление | pressure | 53 | 0.000929 |
| 9 | движение | motion | 52 | 0.000911 |
| 10 | деформация | deformation | 8 | 0.000140 |

**Not found terms** (10):

|  | Term (Ru) | Term (En) |
|---|---|---|
| 1 | агрегатный_состояние_вещество | aggregate_state_matter |
| 2 | правило_равновесие | equilibrium_rule |
| 3 | объём_вещество | substance_volume |
| 4 | правило_равновесие_рычаг | lever_equilibrium_rule |
| 5 | химический_явление | chemical_phenomenon |
| 6 | агрегатный_состояние | aggregate_state |
| 7 | высотомер | altimeter |
| 8 | пластический_деформация | plastic_deformation |
| 9 | превращение_механический_энергия | transformation_mechanical_energy |
| 10 | равновесие_твёрдый_тело | equilibrium_solid_body |

**LLM-found terms** (6):

|  | Term (Ru) | Term (En) |
|---|---|---|
| 1 | правило_равновесие | equilibrium_rule |
| 2 | объём_вещество | substance_volume |
| 3 | правило_равновесие_рычаг | lever_equilibrium_rule |
| 4 | химический_явление | chemical_phenomenon |
| 5 | агрегатный_состояние | aggregate_state |
| 6 | равновесие_твёрдый_тело | equilibrium_solid_body |

**Figure 6.** Term analysis results representation

# References

1. Alier, M., Casañ, M.J., Filvà, D.A.: Smart Learning Applications: Leveraging LLMs for Contextualized and Ethical Educational Technology. In: Gonçalves, J.A.d.C., Lima, J.L.S.d.M., Coelho, J.P., *et al.* (eds.) Proceedings of TEEM 2023. pp. 190–199. Springer Nature Singapore, Singapore (2024). `https://doi.org/10.1007/978-981-97-1814-6_18`

2. Baddour, M., Paquelet, S., Rollier, P., *et al.*: Phenotypes Extraction from Text: Analysis and Perspective in the LLM Era. In: 2024 IEEE 12th International Conference on Intelligent Systems (IS). pp. 1–8 (2024). `https://doi.org/10.1109/IS61756.2024.10705235`

3. Chataut, S., Do, T., Gurung, B.D.S., *et al.*: Comparative Study of Domain Driven Terms Extraction Using Large Language Models (2024), `https://arxiv.org/abs/2404.02330`

4. Esfahani, M.N.: Content Analysis of Textbooks via Natural Language Processing. American Journal of Education and Practice 8(4), 36–54 (2024). `https://doi.org/10.47672/ajep.2252`

5. Foisy, L.O.M., Proulx, E., Cadieux, H., *et al.*: Prompting the Machine: Introducing an LLM Data Extraction Method for Social Scientists. Social Science Computer Review (2025). `https://doi.org/10.1177/08944393251344865`

6. Galushko, I.N.: Historical documents classification using BERT: LLM and historical domain. Perm University Herald. History 2(69), 147–158 (Jun 2025). `https://doi.org/10.17072/2219-3111-2025-2-147-158`

7. Gilbert, M., Crutchfield, A., Luo, B., *et al.*: Using a Large Language Model (LLM) for Automated Extraction of Discrete Elements from Clinical Notes for Creation of Cancer Databases. International Journal of Radiation Oncology*Biology*Physics 120(2, Supplement), e625 (2024). `https://doi.org/10.1016/j.ijrobp.2024.07.1375`

8. Glazkova, A.V., Morozov, D.A., Vorobeva, M.S., *et al.*: Keyword Generation for Russian-Language Scientific Texts Using the mT5 Model. Autom. Control Comput. Sci. 58(7), 995–1002 (Dec 2024). `https://doi.org/10.3103/S014641162470041X`

9. Glazkova, A., Morozov, D., Garipov, T.: Key Algorithms for Keyphrase Generation: Instruction-Based LLMs for Russian Scientific Keyphrases. In: Panchenko, A., Gubanov, D., Khachay, M., *et al.* (eds.) Analysis of Images, Social Networks and Texts. pp. 107–119. Springer Nature Switzerland, Cham (2025). `https://doi.org/10.1007/978-3-031-88036-0_5`

10. Grootendorst, M.: KeyBERT: Minimal Keyword Extraction with BERT. `https://www.maartengrootendorst.com/blog/keybert` (2020), accessed: 2025-08-29

11. Jo, T.: Keyword Extraction. In: Text Mining: Concepts, Implementation, and Big Data Challenge, pp. 421–443. Springer Nature Switzerland, Cham (2024). `https://doi.org/10.1007/978-3-031-75976-5_20`

12. Kapoor, S., Gil, A., Bhaduri, S., *et al.*: Qualitative Insights Tool (QualIT): LLM Enhanced Topic Modeling (2024), `https://arxiv.org/abs/2409.15626`

13. Kuhn, J.: Computational text analysis within the Humanities: How to combine working practices from the contributing fields? Language Resources and Evaluation 53(4), 565–602 (Dec 2019). https://doi.org/10.1007/s10579-019-09459-3

14. Liu, J., Shang, Z., Ke, W., *et al.*: LLM-Guided Semantic-Aware Clustering for Topic Modeling. In: Che, W., Nabende, J., Shutova, E., *et al.* (eds.) Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 18420–18435. Association for Computational Linguistics, Vienna, Austria (Jul 2025). https://doi.org/10.18653/v1/2025.acl-long.902

15. Maragheh, R.Y., Fang, C., Irugu, C.C., *et al.*: LLM-TAKE: Theme-Aware Keyword Extraction Using Large Language Models. In: 2023 IEEE International Conference on Big Data (BigData). pp. 4318–4324 (2023). https://doi.org/10.1109/BigData59044.2023.10386476

16. Monakhov, S.I., Turchanenko, V.V., Cherdakov, D.N.: Terminology use in school textbooks: corpus analysis. RUDN Journal of Language Studies, Semiotics and Semantics 14(3), 437–456 (2023). https://doi.org/10.18413/2313-8912-2023-9-1-0-3

17. Mu, Y., Dong, C., Bontcheva, K., *et al.*: Large Language Models Offer an Alternative to the Traditional Approach of Topic Modelling. In: Calzolari, N., Kan, M.Y., Hoste, V., *et al.* (eds.) Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024). pp. 10160–10171. ELRA and ICCL, Torino, Italia (May 2024), https://aclanthology.org/2024.lrec-main.887/

18. Papagiannopoulou, E., Tsoumakas, G.: A review of keyphrase extraction. WIREs Data Mining and Knowledge Discovery 10(2), e1339 (2020). https://doi.org/10.1002/widm.1339

19. Peng, R., Liu, K., Yang, P., *et al.*: Embedding-based Retrieval with LLM for Effective Agriculture Information Extracting from Unstructured Data (2023), https://arxiv.org/abs/2308.03107

20. Priyanshu, A., Vijay, S.: AdaptKeyBERT: An Attention-Based approach towards Few-Shot & Zero-Shot Domain Adaptation of KeyBERT (2022), https://arxiv.org/abs/2211.07499

21. Sakhovskiy, A., Tutubalina, E., Solovyev, V., *et al.*: Topic Modeling as a Method of Educational Text Structuring. In: 2020 13th International Conference on Developments in eSystems Engineering (DeSE). pp. 399–405 (2020). https://doi.org/10.1109/DeSE51703.2020.9450232

22. Solnyshkina, M.I., Solovyev, V.D., Gafiyatova, E.V., *et al.*: Text complexity as interdisciplinary problem. Issues of Cognitive Linguistics (1), 18–39 (2022). https://doi.org/10.20916/1812-3228-2022-1-18-39

23. Solovyev, V., Solnyshkina, M., Tutubalina, E.: Topic Modeling for Text Structure Assessment: The case of Russian Academic Texts. Journal of Language and Education 9(3), 143–158 (Sep 2023). https://doi.org/10.17323/jle.2023.16604

24. Song, M., Jiang, H., Shi, S., *et al.*: Is ChatGPT A Good Keyphrase Generator? A Preliminary Study (2023), https://arxiv.org/abs/2303.13001

25. Steiss, J., Tate, T., Graham, S., *et al.*: Comparing the quality of human and ChatGPT feedback of students' writing. Learning and Instruction 91, 101894 (2024). `https://doi.org/10.1016/j.learninstruc.2024.101894`

26. Sukying, A., Barrot, J.S.: Friend or Foe? Investigating the Alignment of English Language Teaching (ELT) Textbooks with the National English Curriculum Standards. The Asia-Pacific Education Researcher 34(2), 793–801 (Apr 2025). `https://doi.org/10.1007/s40299-024-00896-5`

27. Tunyan, E.G., Sazikov, R.S., Kharlamov, S.A.: Automatic Extraction of Keywords and Summaries for Knowledge Base Population. Russian Journal of Cybernetics 6(2), 108–113 (Jun 2025), `https://en.jcyb.ru/nisii_tech/article/view/413`

28. Vanin, A., Bolshev, V., Panfilova, A.: Applying LLM and Topic Modelling in Psychotherapeutic Contexts (2024), `https://arxiv.org/abs/2412.17449`

29. Wu, X., Nguyen, T., Luu, A.T.: A survey on neural topic models: methods, applications, and challenges. Artificial Intelligence Review 57(2), 18 (Jan 2024). `https://doi.org/10.1007/s10462-023-10661-7`

30. Yang, J.: Integrated application of LLM model and knowledge graph in medical text mining and knowledge extraction. Social Medicine and Health Management 5(2), 56–62 (2024). `https://doi.org/10.23977/socmhm.2024.050208`

31. Yang, X., Zhao, H., Phung, D., *et al.*: LLM Reading Tea Leaves: Automatically Evaluating Topic Models with Large Language Models. Transactions of the Association for Computational Linguistics 13, 357–375 (2025). `https://doi.org/10.1162/tacl_a_00744`

32. Zagatti, F.R., Lucrédio, D., Caseli, H.d.M.: Unsupervised Statistical Keyword Extraction Pipeline: Is LLM All You Need? In: Paes, A., Verri, F.A.N. (eds.) Intelligent Systems. pp. 460–475. Springer Nature Switzerland, Cham (2025). `https://doi.org/10.1007/978-3-031-79032-4_32`

33. Zhang, T., Kishore, V., Wu, F., *et al.*: BERTScore: Evaluating Text Generation with BERT. In: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net (2020), `https://openreview.net/forum?id=SkeHuCVFDr`