

## РЕФЕРИРОВАНИЕ ХУДОЖЕСТВЕННОЙ ЛИТЕРАТУРЫ ПОСРЕДСТВОМ БОЛЬШИХ ЯЗЫКОВЫХ МОДЕЛЕЙ

© 2025 г. Д. А. Григорьев<sup>1,\*</sup>, Д. В. Худяков<sup>1,\*\*</sup>, Д. И. Чернышев<sup>1,\*\*\*</sup>

Представлено кем-то

Поступило 16.08.2025

После доработки 20.08.2025

Принято к публикации 31.08.2025

Работа исследует методы сжатия художественных текстов с помощью языковых моделей и предлагает улучшенные подходы для точного реферирования в условиях ограниченного контекста.

*Ключевые слова и фразы:* LLM, реферирование, литература, книги, краткий пересказ

DOI: 10.31857/S2686954322040117

### ВВЕДЕНИЕ

**Реферирование художественной литературы.** Автоматическое реферирование текста - одна из ключевых задач в области обработки естественного языка. Суть этой задачи заключается в создании информативного реферата исходного текста с сохранением основного смысла содержания. В последние годы, с появлением больших языковых моделей, резко возрос интерес к автоматизации реферирования в самых разных жанрах текстов, включая художественные произведения. В отличие от научных, новостных или технических текстов, художественные произведения характеризуются высокой степенью стилистической и семантической сложности. Нелинейность повествования, образность, метафоричность и стилистические приёмы делают задачу написания краткого содержания особенно трудоёмкой. Ограниченное контекстное окно современных моделей дополнительно осложняет работу с длинными произведениями.

На текущий момент существует не так много наборов данных, фокусирующихся на задаче художественного реферирования текста, к тому же ключевые открытые корпуса сконцентрированы на нерусскоязычном материале. BookSum [1] - один из первых и наиболее известных англоязычных наборов данных для абстрактного реферирования художественных произведений. Он состоит из книг, пьес и рассказов, сопровождаемых рефератами разного уровня сложности (уровень абзацев, уровень глав, уровень книги). Echoes from Alexandria [2] - многоязычный корпус художественной литературы. Включает тексты на пяти языках: английском, немецком, французском, итальянском и испанском. Fables [3] - ручной корпус, предназначенный для оценки фактологической достоверности рефератов к художественным книгам. Он включает 3 158 утверждений, извлечённых из созданных языковыми моделями рефератов к 26 книгам. Каждое утверждение оценивается по реферату, полученному от различных моделей и анализируются экспертами. По результатам FABLES выявлено, что даже продвинутое модели (например, Claude) допускают до 20–30% фактологических ошибок, включая искажение причинно-следственных связей, неверную характеристику персонажей и смещение акцента на малозначимые детали. рём критериям: соответствие событиям оригинала, логическая корректность и отсутствие искажений.

Теоретически автоматическое реферирование может выполняться двумя основными способами: экстрактивным реферированием (выбор ключевых фрагментов текста) и абстрактным (генерация нового текста на основе содержания оригинала). Обычно для художественной прозы выбирается абстрактное реферирование: ключевые смыслы и сюжетные связи распределены по всему тексту, поэтому экстрактивная выборка предложений даёт фрагментарный, стилистически неоднородный результат и не восстанавливает сюжет, в связи с чем был выбран второй способ реферирования.

Актуальность темы обусловлена растущей потребностью в инструментах, способных автоматически создавать краткие, содержательные и стилистически корректные рефераты к художественным

<sup>1</sup>Московский государственный университет им. М. В. Ломоносова, Москва, Россия

\*E-mail: dagrig14@yandex.ru

\*\*E-mail: hydikov17914@gmail.com

\*\*\*E-mail: chdanorbis@yandex.ru

текстам. Цель данной работы - предоставить подобные инструменты, с целью чего были решены следующие задачи:

1. Был создан новый набор данных для русского языка, включающий в себя художественные произведения и пересказы к ним;
2. Были предложены новые методы реферирования текстов, предлагающие альтернативу существующим и существенно уменьшающих время, требующиеся для создания краткого пересказа книги.

Код и данные работы выложены в открытый доступ<sup>1</sup>.

## НАБОР ДАННЫХ

На момент начала исследования не существовало открытых и репрезентативных корпусов, предназначенных специально для задачи реферирования художественных текстов на русском языке. С целью проведения экспериментов и оценки различных подходов к генерации рефератов был создан собственный корпус, состоящий из художественных текстов и соответствующих кратких пересказов. В качестве источника рефератов был выбран ресурс «Народный Брифли» [4] - платформа, где пользователи публикуют краткие пересказы литературных произведений.

Пересказы представляют произвольные тексты, созданные пользователями на основе исходных материалов художественного произведения. Они варьируются по объему - от нескольких предложений, до нескольких абзацев и стилю - некоторые пересказы дословно воспроизводят ключевые фразы произведения, в то время как другие используют более свободную форму изложения. Некоторые охватывают все произведение целиком, тогда как другие делят содержание на отдельные главы. Как правило, содержат основные факты и выводы, следующие из исходного текста, но могут содержать комментарии автора пересказа.

Художественные тексты были отобраны из электронной библиотеки LibRuSec [5] - одного из крупнейших русскоязычных ресурсов художественной литературы. Отбор произведений осуществлялся на основании наличия реферата на выбранном ресурсе [4]. Каждый текст проходил автоматическую предварительную обработку: удалялась метаинформация (например, заголовки, описания глав и технические вставки), после чего текст форматировался в единый стандартизированный вид, подходящий для дальнейшего использования в моделях.

Чтобы более точно связать книги с их рефератами использовалось семантическое сходство: текст имени автора, записанный на Брифли [4] и автора с LibRuSec [5] переводился в эмбединги с использованием библиотеки SentenceTransformer с помощью языковой модели<sup>2</sup> и сравнивался по косинусному сходству.

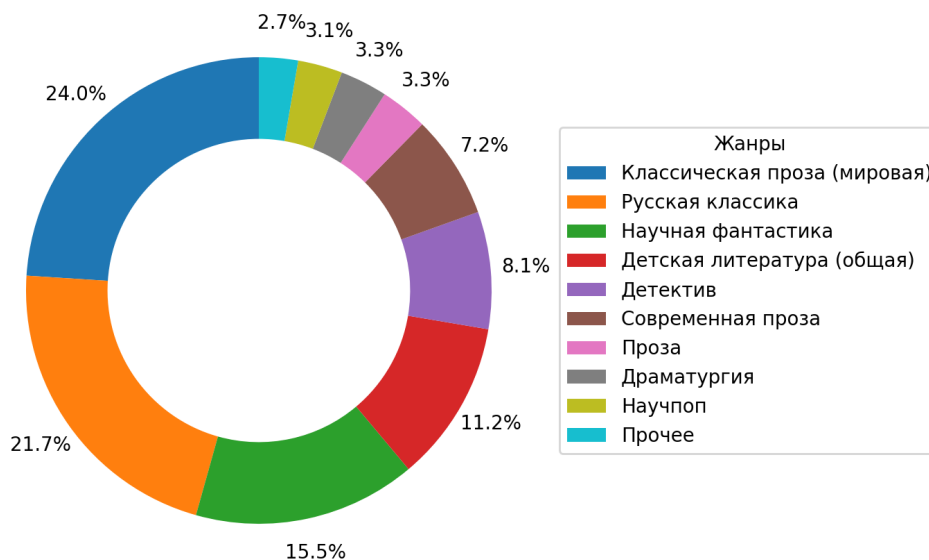


Рис. 1. Распределение текстов по жанрам (топ 10 жанров)

<sup>1</sup><https://github.com/Nejimaki-Tori/BoookSum>

<sup>2</sup><https://huggingface.co/deepvk/USER-bge-m3>

Таблица 1. Обзор датасетов

Датасет	Число документов	Средняя длина документа (# слов)	Средняя длина реферата (# слов)	Степень сжатия (длина реферата / длина текста)
<b>RuBookSum</b>	634	35052.64	700.77	8.43%
BookSum	405	112885.15	1167.20	0.79%
Gazeta	60964	632.77	41.94	6.99%

Тексты рефератов проходили автоматическую очистку от HTML-тегов, комментариев и служебных пометок с помощью LLM Meta-Llama 3-70B-Instruct. Затем производился поиск по датасету LibRuSec и собиралась коллекция, состоящая из пар "текст книги - реферат".

Получившийся корпус включает в себя:

- более 600 очищенных пользовательских пересказов с ресурса «Народный Брифли» [4];
- более 40 различных жанров;
- исходные произведения из электронной библиотеки LibRuSec [5].

На рисунке 1 показано распределение жанров текстов в коллекции. В таблице 1 приведена общая информация о датасете в сравнении с аналогами.

## МЕТОДОЛОГИЯ

### Базовые и модифицированные стратегии.

*Иерархический метод. (Algorithm 1)* Суть этого метода [6] заключается в том, что текст разбивается на фрагменты (чанки), для каждого из которых отдельно генерируется локальный реферат. Эти фрагменты затем объединяются в группы, и из полученных рефератов снова формируется краткое содержание следующего уровня. Последний уровень представляет собой итоговый реферат всего произведения.

*Иерархический метод с фильтрацией узлов. (Algorithm 2)* Классический иерархический метод строит итоговый реферат путём многослойного объединения промежуточных рефератов, полученных из отдельных фрагментов текста. Однако в литературных произведениях часто встречаются фрагменты, которые не оказывают большого влияния на развитие сюжета и содержат множество избыточных повторов и второстепенной информации. Эти фрагменты при генерации итогового реферата могут снижать её информативность, а в некоторых случаях даже мешать модели на этапе реферирования отдельных фрагментов.

Чтобы решить эту проблему, в метод была имплементирована фильтрация узлов по семантической близости. Для исключения малоинформативных или дублирующих фрагментов на каждом уровне иерархии выполняется глобальная проверка семантической близости между всеми промежуточными рефератами. Фрагменты, близкие по косинусной мере с предыдущими, считаются избыточными и не используются при составлении реферата на текущем уровне. Эмбединги получаются с помощью SentenceTransformer (модель USER-bge-m3) и при вычислении на GPU обеспечивается высокая скорость обработки. Эта модификация направлена на ускорение генерации за счет удаления потенциально излишних частей информации, что повышает плотность полезной информации в итоговых рефератах.

**Algorithm 1** Иерархический метод

---

**Require:**  $W$  - контекстное окно модели,  $D$  - входной текст, длиной  $L \gg W$ ,  $p_\theta$  - модель,  $C$  - длина чанка

Разбить  $D$  на чанки  $c_1 \dots c_{\lceil \frac{L}{C} \rceil}$

**for**  $c_i = c_1 \dots c_{\lceil \frac{L}{C} \rceil}$  **do**

$S_0 \leftarrow \text{SummarizeChunk}(p_\theta, c_i)$

**end for**

**repeat**

$\text{Groups} \leftarrow \text{GroupSummaries}(S_l)$

$\ell \leftarrow \ell + 1$

**for**  $g \in \text{Groups}$  **do**

$S_l \leftarrow \{\text{MergeGroup}(p_\theta, g)\}$

**end for**

**until**  $|S_l| = 1$

**return**  $S_l[1]$

---

**Algorithm 2** Иерархический метод с фильтрацией

---

**Require:**  $W$  - контекстное окно модели,  $D$  - входной текст, длиной  $L \gg W$ ,  $p_\theta$  - модель,  $\theta$  - порог сходства,  $C$  - длина чанка

Разбить  $D$  на чанки  $c_1 \dots c_{\lceil \frac{L}{C} \rceil}$

$S_0 \leftarrow \{c_1 \dots c_{\lceil \frac{L}{C} \rceil}\}$

**repeat**

**for**  $s_i \in S_l$  **do**

$e_i \leftarrow \text{Encoder}(s_i)$

$M_{ij} \leftarrow \frac{\mathbf{e}_i \cdot \mathbf{e}_j}{\|\mathbf{e}_i\| \|\mathbf{e}_j\|}$   $\triangleright$  Матрица эмбедингов

Вычисляется максимальное сходство с предыдущими рефератами.

$m_j = \max_{i < j} M_{ji}$

$S_l \leftarrow \{s_i \mid m_i < \theta \text{ or } i = 0\}$   $\triangleright$

Фильтрация

**end for**

$\text{Groups} \leftarrow \text{GroupSummaries}(S_l)$

$\ell \leftarrow \ell + 1$

**for**  $g \in \text{Groups}$  **do**

$S_l \leftarrow \{\text{MergeGroup}(p_\theta, g)\}$

**end for**

**until**  $|S_l| = 1$

**return**  $S_l[1]$

---

«Чертёжный» метод (*Text-Blueprint*). (*Algorithm 3*) Данный метод [7] по сути является модификацией иерархического и ориентирован на построение промежуточного плана реферата перед генерацией текста. План формируется в виде набора вопросно-ответных пар, что повышает управляемость генерации и обеспечивает структурированность результата. Сначала модель формирует список вопросов, отражающих ключевые события, темы и персонажей текста. Далее к каждому вопросу автоматически подбирается краткий ответ. Эта структура служит планом, по которому генерируется итоговый реферат.

«Чертёжный» метод с кластеризацией вопросов. (*Algorithm 4*) Базовая реализация «чертёжного» метода предполагает генерацию вопросно-ответного плана для каждого фрагмента текста и каждого уровня объединения рефератов. Однако при работе с художественными текстами вопросы, генерируемые для каждого чанка, могут пересекаться и порождать противоречивые ответы, то в свою очередь сбивает агрегацию текста моделью, делая реферат менее структурированным и содержательно полным. К тому же, генерация плана на каждом шаге алгоритма существенно замедляет его работу и использует дополнительные мощности языковых моделей. Для снижения числа запросов к модели и повышения структурности, была добавлена кластеризация вопросов с использованием SentenceTransformers и алгоритма K-means.

**Algorithm 3** «Чертежный» метод

---

**Require:**  $W$  - контекстное окно модели,  $D$  - входной текст, длиной  $L \gg W$ ,  $p_\theta$  - модель,  $C$  - длина чанка,  $R$  - ограничение по длине  
Разбить  $D$  на чанки  $c_1 \dots c_{\lceil \frac{L}{C} \rceil}$   
**for**  $c_i = c_1 \dots c_{\lceil \frac{L}{C} \rceil}$  **do**  
     $b_i \leftarrow \text{GenerateBlueprint}(p_\theta, c_i)$   
     $S_0 \leftarrow \{\text{SummarizeWithBp}(p_\theta, b_i, c_i)\}$   
**end for**  
**repeat** ▷ Объединение рефератов  
     $\text{Groups} \leftarrow \text{GroupSummaries}(S_l)$   
     $\ell \leftarrow \ell + 1$   
    **for**  $g \in \text{Groups}$  **do**  
        **if**  $\text{Length}(g) > R$  **then**  
             $b_i \leftarrow \text{GenerateBlueprint}(p_\theta, g)$   
             $S_l \leftarrow \{\text{SummarizeWithBp}(p_\theta, b_i, g)\}$   
        **else**  
             $S_l \leftarrow \{g\}$   
        **end if**  
    **end for**  
**until**  $|S_l| = 1$   
**return**  $S_l[1]$

---

**Algorithm 4** «Чертежный» метод с кластеризацией

---

**Require:**  $W$  - контекстное окно модели,  $D$  - входной текст, длиной  $L \gg W$ ,  $p_\theta$  - модель,  $C$  - длина чанка,  $R$  - ограничение по длине  
Разбить  $D$  на чанки  $c_1 \dots c_{\lceil \frac{L}{C} \rceil}$   
**for**  $c_i = c_1 \dots c_{\lceil \frac{L}{C} \rceil}$  **do**  
     $b_i \leftarrow \text{GenerateBlueprint}(p_\theta, c_i)$   
     $Q \leftarrow \{\text{ExtractQuestions}(p_\theta, b_i)\}$   
**end for**  
**for**  $q_i \in Q$  **do**  
     $E \leftarrow \{\text{Encoder}(q_i)\}$   
     $K \leftarrow K\text{Means}(E)$   
    **for**  $k_i \in K$  **do**  
         $q_i \leftarrow \text{Generalize}(p_\theta, k_i)$   
         $Q \leftarrow \{q_i\}$  ▷ Собирается общий план  
    **end for**  
    **for**  $c_i = c_1 \dots c_{\lceil \frac{L}{C} \rceil}$  **do**  
         $S_0 \leftarrow \{\text{SummarizeWithBp}(p_\theta, b_i, c_i)\}$   
    **end for**  
**end for**  
**Объединение рефератов** аналогично «Чертежному методу» с тем отличием что здесь в качестве чертежа используется один глобальный план  $Q$

---

Такой подход позволяет уменьшить число обращений к LLM, что позволяет ускорить скорость генераций, как будет показано в таблице 2.

## МЕТРИКИ

Для объективного сравнения описанных подходов и моделей в задаче реферирования художественных текстов использовались четыре группы метрик.

**ROUGE-L** [8] - основана на длине наибольшей общей подпоследовательности (LCS) между сгенерированным рефератом  $S$  и эталонным  $R$ . Вычисляется по формуле (3) с использованием формул (1) и (2):

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

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

$$\text{ROUGE-L} = \frac{2 \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3)$$

**BERTScore** [9] - семантическое качество на уровне токенов. Для каждой пары токенов предсказания и эталона вычисляется косинусное сходство их эмбедингов в модели USER-bge-m3. Затем:

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

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

$$\text{BERTScore} = \frac{2 P R}{P + R} \quad (6)$$

В формулах (4), (5) и (6)  $S$  - эталонный текст,  $R$  - сгенерированный; каждое предложение кодируется эмбедингом модели USER-bge-m3, после чего вычисляется косинусное сходство.

**Полнота покрытия ключевых вопросов (Coverage)** - доля заранее сгенерированных по эталонному тексту вопросов с помощью модели Qwen3-235B-A22B [10], на которые модель «отвечает» в реферате:

$$\text{Coverage} = \frac{\#\{q_i : P(\text{“да”} \mid q_i, S) > 0.75\}}{N} \quad (7)$$

В формуле (7)  $N$  - общее число вопросов, а  $P(\text{“да”} \mid q_i, S)$  - вероятность наличия ответа на вопрос  $q_i$  в тексте  $S$ , полученная с помощью LLM (Qwen3-235B-A22B [10]).

**Совпадение ответов (AnswerSimilarity)** - среднее семантическое сходство между сгенерированными ответами  $a_i^{\text{pred}}$  и эталонными  $a_i^{\text{ref}}$  на те же ключевые вопросы:

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

В формуле (8)  $\text{sim}$  - косинусное сходство эмбедингов, полученных через USER-bge-m3.

Использование нескольких метрик, учитывающих как поверхностное совпадение текста, так и глубокое семантическое сходство (BERTScore, AnswerSimilarity), а также степень охвата заранее заданных вопросов (Coverage), обеспечивает всестороннюю и устойчивую оценку качества рефератов.

## ПАРАМЕТРЫ ЭКСПЕРИМЕНТОВ

Все представленные в работе измерения выполнены на тестовой части датасета, отобранных так, чтобы исходные тексты не превышали по длине 800 000 символов. Для всех методов генерируемые рефераты ограничивались максимумом в 500 слов.

Текст на вход разбивался на чанки фиксированного размера в 2000 токенов. Токенизация выполнялась с помощью AutoTokenizer модели DeepPavlov/rubert-base-cased в стандартном режиме. Для воспроизводимости всех случайных процедур использовался фиксированный seed ( $\text{random\_seed} = 42$ ).

**В иерархическом методе с фильтрацией узлов** для оценки избыточности промежуточных рефератов на каждом уровне вычислялась матрица косинусных сходств между их эмбедингами. Порог схожести был установлен равным  $\theta = 0.85$ : если для реферата  $S_j$  существует предыдущий  $S_i$  с косинусным сходством выше этого порога, то  $S_j$  отбрасывается как избыточный. Такой выбор порога обеспечивает компромисс между сохранением значимой информации и устранением дублирования, что эмпирически привело к заметному уменьшению объёма промежуточных представлений без существенной деградации качества.

**В чертёжном методе с кластеризацией вопросов** количество кластеров для K-means выбирается по эвристике, подобранной эмпирически, представленной в формуле (9):

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

где  $N_{\text{questions}}$  - общее число сгенерированных вопросов по всем чанкам до кластеризации. Гарантируется минимум в два кластера, что позволяет даже при небольших наборах вопросов получать структурированное чертёжное представление.

Временные показатели измерялись как среднее значение (в секундах) времени генерации одной книги по каждому методу для 100 книг. В случае всех четырех методов учитывалось суммарное время всех этапов (включая генерацию промежуточных рефератов / планов, фильтрацию и финальную агрегацию).

## РЕЗУЛЬТАТЫ

**Используемые модели.** В экспериментах использовались следующие большие языковые модели: RuadaptQwen2.5-7B-Lite-Beta [11], RuadaptQwen3-32B-Instruct-v2 [11], DeepSeek V3 [12], Qwen3-235B-A22B [10], tpro [13] и yaqpt5lite [14].

**Полученные результаты.** В таблице 2 представлены сравнительные метрики качества автоматического пересказа книг разными моделями и методами обработки. Для каждой комбинации модели и метода измерялись BERTScore, ROUGEL, Answer Coverage и Similarity, а также время генерации (среднее) на 100 примерах, одинаковых для всех замеров. Лучше всего себя показала модель Qwen3-235B-A22B: она продемонстрировала самые высокие показатели в покрытии вопросов и сходстве ответов. В то же время важно отметить, что среди всех методов лучшим образом в соотношение качество и

ТАБЛИЦА 2. Результаты по методам и моделям

Модель	Метрики	Иерархический	Чертежный	Иерархический с фильтрацией	Чертежный с кластеризацией
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	<b><math>53.57 \pm 21.66</math></b>	$40.19 \pm 23.68$	<b><math>45.00 \pm 23.03</math></b>	<b><math>34.68 \pm 23.77</math></b>
	similarity	$42.38 \pm 17.73$	$32.31 \pm 19.33$	$35.64 \pm 18.88$	$27.76 \pm 19.75$
	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	$52.48 \pm 20.79$	<b><math>54.78 \pm 21.16</math></b>	$44.54 \pm 23.03$	$30.19 \pm 21.96$
	similarity	$41.68 \pm 17.18$	$43.99 \pm 17.54$	$35.67 \pm 18.87$	$24.10 \pm 17.62$
	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	$33.12 \pm 21.50$	$33.18 \pm 22.83$	$32.19 \pm 22.52$	$17.72 \pm 15.23$
	similarity	$25.25 \pm 16.94$	$26.21 \pm 18.22$	$24.82 \pm 17.74$	$13.97 \pm 12.39$
	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	$40.27 \pm 20.23$	$40.83 \pm 22.42$	$37.13 \pm 20.72$	$26.03 \pm 18.44$
	similarity	$31.77 \pm 16.63$	$32.60 \pm 18.57$	$29.44 \pm 16.83$	$20.83 \pm 15.26$
	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	$19.66 \pm 17.77$	$24.94 \pm 21.08$	$20.31 \pm 17.95$	$15.51 \pm 14.83$
	similarity	$15.16 \pm 14.11$	$20.03 \pm 17.50$	$15.94 \pm 14.39$	$12.23 \pm 12.30$
	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	$36.85 \pm 19.40$	$33.17 \pm 21.58$	$31.75 \pm 20.06$	$24.28 \pm 16.95$
	similarity	$29.69 \pm 16.43$	$26.58 \pm 18.13$	$25.60 \pm 16.85$	$19.70 \pm 14.29$
	time	$31.02 \pm 28.51$	$113.34 \pm 123.78$	$27.39 \pm 28.05$	$42.15 \pm 56.50$

время обработки себя показывает иерархический метод с фильтрацией узлов. Он позволяет существенно ускорить время обработки (например, почти в два раза для модели DeepSeek V3), и по сравнению с чертежным методом, который в среднем показывал лучшие результаты, не сильно отстает по показателям. Исключением стала лишь модель Qwen3-235B-A22B, так как она показала лучший результат среди всех моделей на базовом чертежном методе. Эксперименты показали, что иерархический метод с фильтрацией узлов обеспечивает наилучший компромисс между скоростью и качеством генерации.

**Анализ и сравнение результатов.** Разброс значений метрики QA можно проиллюстрировать на примере работы одной и той же модели (DeepSeek V3) в рамках иерархического метода. В качестве иллюстрации взяты два реферата к произведениям «И грянул гром» и «Кастрюк». В первом случае модель получила высокий результат, ответив на все, кроме одного вопроса; во втором реферате содержались ответы только на два вопроса из одиннадцати, что привело к низкому показателю. На рисунке 2 показаны два реферата. Для краткости в них выделены только основные моменты, которые повлияли на итоговую метрику. Сравнение показывает возможную причину столь значительного расхождения: реферат к рассказу «Кастрюк» содержит большое количество лирических отступлений и художественных деталей, из-за чего суть произведения сложно уловить и модель отвлекается от фиксации главных фактов, тогда как в «И грянул гром» события изложены последовательно и структурировано, а основные элементы сюжета четко перечислены, что существенно упрощает задачу поиска важной информации. В текстах выделены жирным шрифтом фрагменты, которые несут в себе важную сюжетную информацию, а подчеркнутый текст - то, что можно было бы опустить.

Переходя к сравнению между моделями, можно отметить, что в целом DeepSeek V3 показывает лучшие показатели, чем модели меньшей категории, однако, если сравнивать чертежный метод, то в 30% случаев модель RuadaptQwen3-32B-Instruct-v2 показывает лучшие результаты, а tpro в 43%.

Название	Текст
И грянул гром	...Главный герой, Экельс, азартный и самоуверенный охотник, платит огромную сумму за возможность отправиться на 60 миллионов лет назад, чтобы убить тираннозавра. Перед путешествием гид Тревис строго предупреждает его о правилах: ни в коем случае нельзя сходить с антигравитационной Тропы или вмешиваться в естественный ход событий, так как малейшее нарушение может катастрофически изменить будущее... Тревис объясняет хрупкость временного баланса: даже гибель одной мыши способна уничтожить целые виды, а значит, и изменить историю человечества. Группа выслеживает тираннозавра, помеченного красной краской — это знак, что его убийство не повлияет на будущее. Однако при виде гигантского хищника Экельс впадает в панику, сходит с Тропы и случайно раздавливает бабочку... По возвращении в 2055 ... мир изменился до неузнаваемости: язык стал грубым, атмосфера — тяжёлой, а вместо умеренного президента Кейта у власти стоит жестокий диктатор Дойчер. Экельс осознаёт, что его неосторожность спровоцировала «эффект бабочки» — раздавленное насекомое вызвало цепь событий, искаживших историю. В отчаянии он умоляет исправить ошибку, но Тревис, понимая необратимость последствий, поднимает ружьё. ...
Кастрюк	... Действие рассказа разворачивается в русской деревне ранней весной, где природа пробуждается, но жизнь людей остаётся тяжёлой и однообразной. Главный герой — старик Семён, прозванный Кастрюком, — доживает свои дни в одиночестве, терзаемый воспоминаниями о былой силе и сожалениями о нынешней немогущности. Когда-то он славился как лучший работник в округе, но теперь, дряхлый и забытый, вынужден оставаться в стороне, пока односельчане трудятся в поле. Его единственная отрада — внучка Дашка, добрая и впечатлительная девочка, которая прибегает к нему, испугавшись барчуков из соседнего имения Залесное. Кастрюк успокаивает её, и они вместе отправляются за деревню, где старик, любящая весенней природой, пытается отвлечься от гнетущих мыслей. ... Лишь к вечеру, уговорив сына отпустить его в ночное (пасти лошадей), Кастрюк обретает краткую радость. На свободе, среди ребятишек и под звёздным небом, он чувствует себя почти молодым. У пруда кобыла пьёт воду, отражая закат, а сам старик, глядя на Млечный Путь, шепчет молитву — будто вновь обретает связь с миром и утраченную гармонию. Но это лишь мимолётное утешение: завтра его снова ждёт беспробудное одиночество и осознание собственной ненужности. ...

Рис. 2. Сравнение лучшего и худшего сгенерированного реферата

Модель	Текст
RuadaptQwen3	"Компания «Сафари во Времени» организует платные экскурсии в прошлое для охоты на динозавров, используя машины времени, способные перемещаться между эпохами. Клиенты обязаны соблюдать строгие правила: следовать по металлической тропе...
tpro	"В тексте главный герой, Экельс, отправляется на сафари во времени с целью убить динозавра Tugannosaurus rex. Компания, организующая сафари, гарантирует только динозавров и строго запрещает охотникам сходить с Тропы ... Мистер Тревис, проводник сафари, объясняет, что даже уничтожение одной мыши может привести к исчезновению всех её потомков ...
DeepSeek V3	"**Краткое содержание по плану:** 1. **Экельс** — охотник ... 2. **Компания «Сафари во времени»** организует охоту в прошлом ... 3. **Тревис** — проводник, контролирующий экспедицию. ...

Рис. 3. Сравнение моделей при генерации рефератов по чертежному методу

Для сравнения можно взять реферат по произведению «И грянул гром», созданную с использованием чертежного метода, небольшие вырезки которой приведены на рисунке 3. В то время как реферат,



созданный моделью DeepSeek V3 больше похожа на перечисление основных событий через нумерованный список, текст у моделей RuadaptQwen3-32B-Instruct-v2 и tpro является связным пересказом текста, раскрывающим все основные события сюжета.

Следует отметить, что лучшего результата удалось добиться именно чертежным методом с помощью большой модели Qwen3-235B-A22B, как было показано в таблице 2. Для сравнения качества рефератов можно взять рассказ «Барбос и Жулька» - в иерархическом методе модель Qwen3-235B-A22B посчитала, что «Жулька» - не собака, а лошадь. Также, например, DeepSeek V3 более строго следует шаблону чертежного метода и вместо связного текста пересказа получается нумерованный список пунктов с ключевыми событиями и главными героями. Однако Qwen3-235B-A22B пишет обычный текст, без списков. Таким образом, чертежный метод без модификаций позволил достичь наилучшего результата с использованием лучшей доступной моделью - Qwen3-235B-A22B.

**Замеры времени.** Проводились первоначальные замеры скорости работы методов на небольших текстах, полученные результаты в секундах (среднее по трём запускам) представлены в таблице 3. Результаты подтверждают, что модификации позволяют повысить скорость генерации.

Таблица 3. Время генерации рефератов (в секундах) для текста размером 81,049 символов (11 чанков). Усреднено по трём запускам.

Модель	Иерархический	Иерархический с фильтрацией	Чертежный	Чертежный с кластеризацией
DeepSeek V3	237.83	72.42	292.80	268.75
Qwen3-235B-A22B	113.24	39.45	215.63	145.20
RuadaptQwen3-32B-Instruct-v2	218.23	72.54	420.95	470.4
tpro	472.23	127.38	421.65	185.94
RuadaptQwen2.5-7B-Lite-Beta	84.64	25.70	103.66	78.99
yagpt5lite	34.17	14.08	99.70	27.26

Интересно отметить, что сверхкрупные модели, такие как Qwen3-235B-A22B и DeepSeek V3, продемонстрировали более высокую скорость работы, чем некоторые модели с размером 32B. Ключевая причина этого заключается в использовании архитектуры MoE (Mixture of Experts): во время генерации активна лишь ограниченная часть параметров (например, порядка 30 млрд вместо всех 600 млрд), кроме того, такие модели, как правило, дополнительно оптимизированы для повышения производительности.

## ЗАКЛЮЧЕНИЕ

В заключение, был создан первый открытый датасет, объединяющий тексты книг и рефератов к ним с открытого ресурса «Народный Брифли» [4]. В работе предложены два улучшенных подхода к реферированию художественных текстов с использованием LLM: иерархический с фильтрацией и чертежный с кластеризацией. Иерархический метод с фильтрацией позволяет ускорить генерацию при минимальной потере качества, что делает этот метод пригодным для обработки длинных произведений в условиях ограниченного контекста моделей.

Сравнительный анализ показал, что крупные модели, такие как DeepSeek V3 и Qwen3-235B-A22B, в большинстве случаев обеспечивают более высокое покрытие QA и большую полноту рефератов по сравнению с компактными моделями, особенно в иерархическом и чертежном методах. Однако для некоторых типов текстов и методов (например, базовый чертежный) более компактные модели, такие как RuadaptQwen3-32B-Instruct-v2, могут демонстрировать конкурентоспособное качество при меньших вычислительных затратах. Таким образом, выбор модели следует определять исходя из баланса между доступными ресурсами, требованиями к качеству и характером обрабатываемых текстов.

## СПИСОК ЛИТЕРАТУРЫ

- [1] BOOKSUM: A Collection of Datasets for Long-form Narrative Summarization / Wojciech Kryscinski, Nazneen Rajani, Divyansh Agarwal et al. // Findings of the Association for Computational Linguistics: EMNLP 2022 / Ed. by Yoav Goldberg, Zornitsa Kozareva, Yue Zhang. - Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, 2022. - . - Pp. 6536–6558. <https://aclanthology.org/2022.findings-emnlp.488/>.
- [2] Echoes from Alexandria: A Large Resource for Multilingual Book Summarization / Alessandro Scir'e, Simone Conia, Simone Ciciliano, Roberto Navigli // Findings of the Association for Computational Linguistics: ACL

- 2023 / Ed. by Anna Rogers, Jordan Boyd-Graber, Naoaki Okazaki. - Toronto, Canada: Association for Computational Linguistics, 2023. - . - Pp. 853–867. <https://aclanthology.org/2023.findings-acl.54/>.
- [3] FABLES: Evaluating faithfulness and content selection in book-length summarization / Yekyung Kim, Yapei Chang, Marzena Karpinska et al. // First Conference on Language Modeling. - 2024. <https://openreview.net/forum?id=YfHxQSoaWU>.
- [4] *Народный Брифли*. Электронная библиотека кратких пересказов литературных произведений. <https://wiki.briefly.ru/> (дата обращения: 30.07.2025).
- [5] Библиотека художественных произведений. <https://librusec.org/> (дата обращения: 30.07.2025).
- [6] Wu J. et al. Recursively Summarizing Books with Human Feedback // arXiv e-prints. - 2021. - С. arXiv: 2109.10862.
- [7] Text-Blueprint: An Interactive Platform for Plan-based Conditional Generation / Fantine Huot, Joshua Maynez, Shashi Narayan et al. // Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations / Ed. by Danilo Croce, Luca Soldaini. - Dubrovnik, Croatia: Association for Computational Linguistics, 2023. - . - Pp. 105–116. <https://aclanthology.org/2023.eacl-demo.13/>.
- [8] *ROUGE*. Lin C. Y. Rouge: A package for automatic evaluation of summaries // Text summarization branches out. - 2004. - С. 74–81.
- [9] *BERTScore*. BUCKLEY C. Evaluating Evaluation Measure Stability // ACM SIGIR 2000 Proceedings. - 2000.
- [10] *Qwen3-235B*. Yang A. et al. Qwen3 technical report // arXiv preprint arXiv:2505.09388. - 2025.
- [11] *RuadaptQwen*. Tikhomirov M., Chernyshev D. Facilitating large language model russian adaptation with learned embedding propagation // Journal of Language and Education. - 2024. - Т. 10. - №. 4 (40). - С. 130–145.
- [12] *DeepSeek V3*. Liu A. et al. DeepSeek-V3 Technical Report // CoRR. - 2024.
- [13] Т-Банк открыл доступ к собственной русскоязычной языковой модели в весовой категории 7-8 млрд параметров  
Т-Банк URL: <https://www.tbank.ru/about/news/20072024-t-bank-opened-access-its-own-russian-language-language-model-weight-category-of-7-8-billion-parameters/> (дата обращения: 10.05.2025).
- [14] YandexGPT 5 с режимом рассуждений // Яндекс URL: <https://ya.ru/ai/gpt?ysclid=mal9jrssc8906806775> (дата обращения: 30.07.2025).

## LITERATURE SUMMARISATION WITH LARGE LANGUAGE MODELS

D. A. Grigoriev<sup>a,\*</sup>, D. V. Khudiakov<sup>a,\*\*</sup>, D. I. Chernyshev<sup>a,\*\*\*</sup>

<sup>a</sup>Lomonosov Moscow State University, Moscow Center for Fundamental and Applied Mathematics,  
Moscow, Russian Federation

*man who sold the world*

This work explores methods for compressing literary texts using language models and proposes improved approaches for accurate summarization under limited context conditions.

*Keywords:* LLM, summarization, literature, books, brief retelling