*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

135

Victoria I. Firsanova[1] 🆔

**A graph-based approach to closed-domain natural language generation**

**[1]** St Petersburg State University,
*7-9 Universitetskaya Emb., St Petersburg, 199034, Russia*
*E-mail: st085687@student.spbu.ru*
ORCID: 0000-0002-8474-0262

**Abstract:** Graph-based Natural Language Processing (NLP) methods have seen significant advancements in recent years with the development of Large Language Models (LLMs) and Retrieval Augmented Generation (RAG). LLMs are sophisticated models that recognize numerous NLP tasks by analyzing the users' natural language instructions called prompts. However, their industrial use is questionable due to such ethical concerns as false information generation called hallucinations, high risks of data breaches, and plagiarism. The paper introduces a novel NLP architecture, the Graph-Based Block-to-Block Generation (G3BG), which leverages state-of-the-art deep learning techniques, the power of attention mechanisms, distributional semantics, graph-based information retrieval, and decentralized networks. The model encodes user prompts to mitigate data breach risk, retrieves relevant information from a graph knowledge base, and forms a block for a conditional language model using LLMs to perform a new secure type of RAG. The model is closed-domain and small-scale oriented. It exhibits superior performance across low-resource NLP tasks, which makes it prominent for industrial use. The research presents a novel graph-based dataset. The dataset comprises private data features to encode and closed-domain textual information for information retrieval. The dataset is used to train and evaluate the G3BG model. The model allows cutting 100x training dataset volume achieving Perplexity ~6.51 on the Language Generation task and F1-Score ~90.3 on the Information Retrieval task comparable to most state-of-the-art language models. The experimental results prove the effectiveness of the proposed method and contribute to the algorithmic approaches toward LLM risk mitigation.

**Keywords**: Language Generation; Language Understanding; Generative Artificial Intelligence; Large Language Models; Decentralized Networks; Data Encoding; Distributional Semantics; Closed-Domain Systems

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

136

Фирсанова В. И.[1]

**Новый графовый подход к генерации текстов узкой предметной области на естественном языке**

[1] Санкт-Петербургский государственный университет
Университетская наб. 7-9, Санкт-Петербург, 199034, Россия
*E-mail: st085687@student.spbu.ru*
ORCID: 0000-0002-8474-0262

**Аннотация.** Обработка естественного языка на основе графов в последние годы становится актуальной благодаря развитию больших языковых моделей и генерации, дополненной информационным поиском. Большие языковые модели – это сложные алгоритмы, которые распознают многочисленные задачи обработки естественного языка путем анализа инструкций пользователей на естественном языке. Однако их промышленное использование вызывает сомнения из-за таких этических проблем, как создание ложной информации, высокого риска утечки данных и авторских заимствований. В статье представлена новая архитектура для обработки естественного языка, поблочная генерация на основе графов, которая использует самые современные методы глубокого обучения, возможности механизмов внимания, дистрибутивной семантики, информационного поиска на основе графов и децентрализованные сети. Модель кодирует запросы пользователя для снижения риска утечки данных, извлекает релевантную информацию из базы знаний графа и формирует блок для обусловленного моделирования языка с использованием больших языковых моделей. Модель направлена на разрешение ситуации недостатка данных для обучения полноценной модели машинного обучения. Исследование представляет новый набор данных на основе графов. Набор данных задает признаки уязвимых персональных данных для кодирования и текстовую информацию закрытой предметной области для информационного поиска. Он используется для обучения и оценки модели поблочной генерации на основе графов, впервые представленной в данной статье. Модель позволяет сократить объем обучающих данных более чем в 100 раз, достигная значения метрики оценки перплексии ~6,51 в задаче генерации естественного языка и F1-меры ~90,3 в задаче извлечения информации, что сопоставимо с большинством современных языковых моделей. Результаты экспериментов доказывают эффективность предлагаемого метода и вносят вклад в разработку алгоритмических подходов к снижению рисков использования больших языковых моделей в промышленности.

**Ключевые слова:** Генерация естественного языка; Понимание естественного языка; Генеративный искусственный интеллект; Большие языковые модели; Децентрализованные сети; Кодирование данных; Дистрибутивная семантика; Закрытая предметная область

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

137

## INTRODUCTION

In recent years, the remarkable Natural Language Processing (NLP) advancements have been driven by the development of Large Language Models (LLMs). For example, models derived from GPT (Ouyang et al., 2022: 27733), and multi-task conversational agents, such as Mistral (Jiang et al., 2023: 2), are state-of-the-art across many NLP tasks. LLMs are probabilistic agents for causal language modeling. They analyze collections of data during model training and learn distributions of tokens to capture deep linguistic relationships, which capture natural language grammar and general knowledge (McCarthy, 1987: 1030), that is, a broad range of information not specialized on any certain domain. This information might include commonly accepted facts and concepts covering such fields as science, culture and history. Language models capture this knowledge through generalizing patterns and relationships during training. For example, scientific information is often described with special terms, such as technical terminology. Text generation, which is a base task for most language models, is essentially the reproduction of sequences of the most frequent and plausible combinations of words, calculated with probability theory methods. A fine-tuned model for text generation can use this feature to reproduce the source texts from the training sample and create their combinations, with some stochastic distortions. As a result, a language model can convey knowledge from these texts. This process is closely related to meta-learning, a model's ability to learn strategies to solve novel tasks without being explicitly trained on them, based on extracted patterns and relationships (Schmidhuber, 1987: 5).

The LLMs' ability to capture deep linguistic relationships allows using them in conversational agents, such as OpenAI ChatGPT or Google Gemini. Despite the LLMs' superior performance across numerous NLP tasks, the models do not store any information explicitly. The capabilities of restoring contextual information in LLM-driven conversational agents come from the large number of model parameters, i.e. the number of features representing deep relationships across tokens. Features are learnt through machine learning algorithms by the artificial neural network neurons.

For example, in GPT 3.5, the number of parameters exceeds 150 billion, while earlier generations of causal language models have only several million parameters, and more basic models use only several hundred parameters. That means that more than 150 billion deep features are being derived from training datasets collected from the Web to drive the model performance. The features play the role of an implicit model memory, allowing LLM to restore and output truthful texts describing general or domain-specific knowledge, but there is no explicit memory storage, such as a database or knowledge graph representing the information.

Perplexity (Jelinek et al., 1977: S63) is a commonly used metric score in language modeling. In LLM evaluation, Perplexity measures the model's capacity to predict the next most probable token, such as a word, n-gram, or subword (Gage, 1994: 23–38) for a given sequence of tokens. Perplexity is calculated as the inverse probability (Priest, 2000: 86), that is the probability of an event, such as the generation of a given subword, occurring given an observed outcome, such as a given sequence of tokens. Given a language model trained on a set of data, perplexity is calculated as the inverse probability of the test set, normalized by the number of words. That is why a lower Perplexity score indicates high predicting capacity. Low Perplexity scores are associated with an ability to understand the language structure, since the ability to generate sequences implies modeling

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

138

cohesion and coherence, which is impossible without identifying patterns of linguistic structure (Morris, Hirst, 1991: 21, 25-26). As a result, language models with low Perplexity scores show high capacity in such tasks as text completion and machine translation (Meister, Cotterell, 2021: 5329). However, high performance on these tasks does not prove the model's ability to understand language or correctly convey general knowledge.

Perplexity measures the probability distribution of word sequences, focusing on syntactic and local contextual accuracy, but does not take into account factual correctness or context understanding. As mentioned, most LLMs lack explicit information storage. Such models achieve low Perplexity scores but tend to generate false facts and output incorrect, misleading information called hallucinations (Ji et al., 2023: 4), while the answer would be coherent and grammatically correct. That raises ethical concerns regarding the industrial use of LLMs. Evaluation metrics that assess the factual accuracy of LLMs include Precision and Recall, human judgment, and specialized benchmarks, such as Benchmarking Information Retrieval (BEIR) (Thakur et al., 2021: 1-2) and Massive Multitask Language Understanding (MMLU) (Hendrycks et al., 2021: 2) datasets, that test models on tasks such as question answering, fact-checking, or real-world information retrieval. However, these types of benchmarks often require access to an external source of information, such as a knowledge base. Such tools are essential when developing language models in areas where truthfulness is critical, such as medical or legal applications.

One of the solutions is Retrieval-Augmented Generation (RAG) (Zhang et al., 2023: 2). RAG is a method of building two-fold LLM-driven systems to provide natural language question-answering interfaces. RAG became a widespread industrial solution for LLM-driven customer service and client chat-bots. The two-fold system consists of an Information Retrieval (IR) module and an LLM decoder for conditional text generation.

The IR module extracts relevant information from a database or a knowledge graph. The decoder uses the extracted information to generate coherent text based on the retrieved facts.

RAG makes LLMs' outputs more controllable and predictable, although it does not solve the problem of false information generation completely. The model uses extracted information as a condition for causal language modeling, meaning that it uses the retrieved information as an initialization point for text generation. The model generates an answer by predicting a plausible continuation for the user prompt given the extracted data, but it does not process the extracted data using logical inference. Due to the LLM stochastic nature, RAG outputs are often difficult to control (Zhang et al., 2023).

Another problem is data breaches and plagiarism generation caused by LLMs' memorization capabilities. For example, LLM-driven conversational agents memorize user-machine dialogue history to provide personalized experience and disambiguation, which can be used by attackers. Combining the advantages of decentralized security algorithms (Luo et al., 2023: 4) with proper personal data encoding can solve this problem.

The decentralized security algorithms often use blockchain technology. A blockchain is a distributed database that can be used for secure data storage. The advantages of combining LLM and blockchain technologies include the ability to track and analyze the history of operations, such as human-machine interactions, through the chain of blocks. This study proposes to store prompts and meta-information for LLM text generation conditioning, such as personal user information, with the LLM generation result in blocks. Each block is encrypted and broadcast to a distributed network. This process is achieved by writing a block into the chain. In order to write a block to the chain, the data must be validated. LLM guardrails can be used as a blockchain verification mechanism. LLM guardrails can be used as

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

139

rules to check the data validity before writing it to the blockchain.

LLM guardrails (Dong et al., 2024) are algorithms that process model input or output to control the model behavior through a set of guiding instructions. The guardrails can filter potentially malicious information, or adjust the model's responses. The guardrails can be used to ensure ethical principles, filter potentially harmful of biased content, or monitor security violations. Guardrails can use various technologies and their ensembles, including, rule-based approaches, machine learning models, and human feedback to enhance LLM security and trustworthiness. In this study, the guardrails also ensure the blockchain validation process by checking the model output before writing the data to a blockchain.

The paper suggests that combining three state-of-the-art technologies – RAG, blockchain, and LLMs with guardrails – in one novel solution can result in an energy-efficient LLM-driven framework that is protected from data breaches and false information generation. The paper focuses on three SOTA LLMs: GPT-4o mini (OpenAI, 2024), Mistral 2 (Mistral, 2024), and Claude 3.5 Sonnet (Anthropic, 2024). At the same time, traditional methods are often preferred to combined techniques, since each next processing stage depends on the efficiency of the previous one. The study goal is to explore whether the cautious approach toward complex LLM-driven frameworks is justified by assessing their reliability compared to traditional RAG. This goal is achieved by designing a custom LLM-based system for inclusive education needs. In this study, we prove the following research hypotheses:

Hypothesis 1: Advanced LLM-driven frameworks consisting of several modules, such as an information retrieval module, privacy and security mechanisms, and text generation engine, that are built to solve specific problems or serve specific domains, offer better reliability and text generation quality than traditional RAG pipelines, such

as Retrieval-Augmented Language Model (REALM) (Guu et al., 2020).

Hypothesis 2: RAG methods, which have already proven their effectiveness, are more reliable than complex methods fine-tuned to solve specific problems or serve specific domains.

1. To test the research hypotheses, the research will focus on the following tasks:

2. Collect the RDF graph dataset about inclusive education for RAG.

3. Collect plain text data and generate grammaticality judgments to compile a dataset for LLM linguistic competence analysis.

4. Implement a set of guardrails, and a blockchain-based framework for LLM security enhancement.

5. Implement a mechanism of tracking human-machine interaction history, and a mechanism to collect and store user information for LLM personalization.

6. Implement traditional RAG framework, and a RAG framework with guardrails.

7. Compare SOTA LLM baselines within each of the implemented frameworks, including GPT-4o mini, Mistral 2, and Claude 3.5 Sonnet.

8. Assess the linguistic competence of the LLMs for the interpretability purposes.

The study highlights the LLMs' linguistic competence to measure the LLMs' sensitivity toward multilinguality and accessibility, since the study focuses on the inclusive education app development. High linguistic competence ensures that LLMs can be accessed in diverse inclusive environments. This study involves the development and use of the following datasets:

1. Training and test data: The models will be trained and tested using a knowledge base constructed as a Resource Description Framework (RDF) (Powers, 2003: 19) graph containing detailed information about inclusive education. This knowledge base provides the model with content related to

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

140

education policies, practices, and support mechanisms.

2. Model evaluation data: Evaluation will be conducted using a synthetic dataset designed for linguistic experiments. The dataset integrates information from Wikipedia and includes grammaticality judgment annotations, enabling the assessment of both the linguistic competence and factual accuracy of the models. This linguistic analysis allows for assessing the model's ability to differentiate between grammatical and non-grammatical sentences and its proficiency in providing general knowledge.

The LLM-driven frameworks addressing security, personalization quality, hallucination resistance, as well as linguistic competence are assessed through a combination of performance metrics and user feedback. The LLM security is measured with the F1-Score (Van Rijsbergen, 1979: 134) obtained on LLM binary classification of malicious and safe prompts. The LLM personalization quality is measured manually through quantitative assessment of the outputs' relevance, tone appropriateness, and alignment with user intent. The LLM hallucination resistance is measured with the length-aware F1-Score, a novel metric presented in this study. The length-aware F1-Score is measured on information retrieval task using the RDF graph dataset proposed in this study. This RDF knowledge base tests the model's ability to provide structured information about inclusive education. Additionally, the LLM linguistic competence was measured through qualitative evaluation. The qualitative assessment uses a novel synthetic dataset proposed in this study that tests the model's linguistic competence by assessing its ability to judge grammatical and non-grammatical sentences.

The study highlights the development of a dialogue system for inclusive education. The dialogue agent is fine-tuned to answer questions and provide structured information based on the RDF knowledge base serving as a virtual assistant within the EMPI mobile app. The mobile app is designed to provide

psychological support and information about inclusive education for students and educators of inclusive environments. The app provides an access to a conversational artificial intelligence agent through inclusive user interface. For more detailed information about the EMPI app, refer to vifirsanova.github.io/empi-web. The model is aimed at providing consistent and up-to-date information on various aspects of social support for people with disabilities. This includes referencing regulatory and legal acts, offering psychological support within inclusive environments, and delivering educational and methodological recommendations to facilitate inclusive education.

The study novelty includes the following contributions:

Novel LLM-driven framework: To our knowledge, this is the first study to propose integrating Retrieval-Augmented Generation (RAG) with blockchain technology and a set of LLM guardrails for enhanced security and reliability.

Evaluation method: The study proposes a novel qualitative approach to assess LLM linguistic competence by measuring the model's ability to recognize grammatical and non-grammatical sentences through detailed human evaluation of LLMs' grammaticality judgments, as well as a length-aware F1-Score, an adaptation of the traditional F1-score that measures the quality of information extraction in RAG with respect to the LLM response length.

Datasets: The study utilizes two novel distinct datasets: an RDF graph focused on inclusive education available at https://github.com/vifirsanova/empi/blob/main/blockchain/tbl.ttl and a synthetic dataset tailored for linguistic grammaticality judgments available at https://huggingface.co/datasets/missvector/multi-wiki-grammar.

**RELATED WORK**

Modern dialogue systems often incorporate databases (Gao et al., 2018: 30) with various structures to implement the

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

141

information extraction stage, or they use the implicit model memory capability known as meta-learning (Schmidhuber, 1987: 5). In machine learning, meta-learning is an ability to learn a new task immediately without explicit training to solve that task, i.e. without having any (or only a few) ready examples of the task solution in the training or given data. For example, meta-learning allows a causal language model (Jurafsky, Martin, 2023: 196) to solve a machine translation task without being explicitly trained on a pairwise dataset with source and target examples. Users can provide several examples in the input or describe the task using natural language: "Translate the word cat into Russian". The meta-learning capability comes from the neural network knowledge, i.e. features encoded in the model neurons. Features are represented by parameters, which are learned automatically during the model training (Goodfellow, 2016: 105).

**Meta-learning.** Transformer-based models (Vaswani et al., 2017: 4) have robust meta-learning capacities due to their ability to learn deep relationships across the input data from the Attention mechanism. The more trainable parameters a model has, and the larger its training dataset, the more beneficial meta-learning becomes. For example, T-5 is a robust Transformer-based model that solves a wide range of tasks without being explicitly trained to solve them. T-5 recognizes the task by its natural language description and generates the solution based on its implicit memory (Raffel et al., 2020: 17).

The LLM capacity towards user intention recognition was developed in such model architectures as InstructGPT (Ouyang et al., 2022: 27733) or Mistral (Jiang et al., 2023: 2). Such models use the power of Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017: 3) to build a policy, i.e. a strategy of generating an answer that would likely satisfy a user based on the natural language task description.

**LLM hallucinations.** Causal language models fine-tuned to align with user intent based on human feedback is a robust solution, which, nevertheless, has notable drawbacks. One significant issue is artificial intelligence hallucinations, a phenomenon of generating incorrect or misleading content. Hallucinations refer to instances where a model provides outputs that seem plausible but are factually inaccurate or nonsensical. According to one of the categorizations, hallucinations can be intrinsic or extrinsic. Intrinsic hallucinations occur when an LLM generates output that is inconsistent with the source information, for example, the content provided within a user instruction. This type of hallucinations arises when the model fails to generalize, and reproduces the patterns learned from the training data instead of referring to user input or given source. Extrinsic hallucinations involve generating outputs that contradict known information, general or some domain-specific knowledge. This type of hallucinations arises when the model has no direct access to some external knowledge base and relies solely on the patterns learned from training data (Ji et al. 2024: 4). Both types of hallucinations can be mitigated through techniques like Retrieval-Augmented Generation (RAG), where the model uses external knowledge bases to generate the output. The RAG method reduces intrinsic hallucinations by providing the ability to compare source information with generated text using information extraction and vector similarity search. This approach also reduces extrinsic hallucination by referencing external data during the text generation (Zhang et al., 2023: 2). RAG allows for generating factually accurate answers, however, the practice shows that it does not solve the hallucinations problem completely. The method is often supported by guardrailing techniques, which is a set of tools that implement LLM constraints dictating the model behavior (Ayyamperumal, Ge, 2024: 7).

**Model compression.** The large size of LLMs often limits their deployment in low-resource environments, such as mobile devices. To address this, techniques like

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

142

GGML/GGUF (Gerganov, 2024) for optimizing memory usage, as well as methods like Low-Rank Adaptation (LoRA) and Quantized LoRA (QLoRA), have been developed to enable small and efficient models suitable for constrained computational conditions. GGML and GGUF are file formats based on C++ library designed to run machine learning models on low-resource hardware, such as personal computers, mobile devices, and other environments with limited processing power through quantization, that is converting the numbers representing a model's parameters (such as weights and activations matrices) to lower precision values (Gerganov, 2024). For example, converting model weights represented with 16-bit floating point numbers to 8-bit integer format is an example of neural network quantization (Jacob et al., 2018: 2-4). Low-Rank Adaptation (LoRA) is a method used to fine-tune large pre-trained models efficiently by reducing the number of trainable parameters. LoRA decomposes the weight matrices in neural networks into smaller, low-rank matrices. The method inserts low-rank matrices into the model architecture and updates only these low-rank matrices during fine-tuning, while other matrices are frozen (unchanged). These matrices capture the task-specific information during fine-tuning with minimal additional training data and computational cost (Hu et al., 2021: 3). Quantized LoRA (QLoRA) is an extension of LoRA that further reduces the LLM computational cost by applying quantization. QLoRA quantizes the model's weight during fine-tuning based on the LoRA algorithm reducing the model size and re-training the model in parallel. QLoRA enables the deployment of large language models on devices with extremely limited resources, such as mobile devices (Dettmers et al., 2024: 3).

**Prompt injections**. LLMs face cybersecurity risks, such as prompt injections (Choi et al., 2022: 2), where attackers manipulate model behavior through crafted input instructions. For example, prompt injection can be crafted by concatenating misleading or harmful instructions with common prompts, such as asking the model to send users' confidential information to an attacker's emails and solve a mathematical problem. The mathematical problem is a common and safe prompt, which an attacker uses to mask the malicious instruction to send confidential information. The masking allows to bypass security filters and cause data leakage. To mitigate the risks posed by prompt injections, developers implement guardrails for input sanitization, that is, pre-processing and filtering potentially malicious prompts, detecting suspicious patterns in input queries, as well as checking the output content according to ethical guidelines. Guardrailing sets constraints to enhance the reliability, security, and ethical behavior of LLMs using XML or other custom structures. NVIDIA Nemo or Guardrails AI are widespread implementations for these strategies (Ayyamperumal, Ge, 2024: 5, 7).

LLM security can be ensured by incorporating a decentralized approach into the NLP architecture. Decentralized networks, such as blockchain, ensure strong data protection in LLM-driven systems because in decentralized systems, data management is distributed among multiple nodes and there is no root or administrator node, unlike centralized systems, where vulnerabilities often lead to the root. Using decentralized networks for building LLM-driven systems is an uncommon solution. Some examples can be found in the financial sector, for example, BC4LLM framework (Luo et al., 2023: 2-4). Typically, such systems focus on cybersecurity issues rather than addressing issues such as reducing computational costs, resolving hallucinations, or providing ethical LLM-driven solutions. This study proposes a solution that will shift the focus on using decentralized networks from developing financial and commercial solutions to solving

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

143

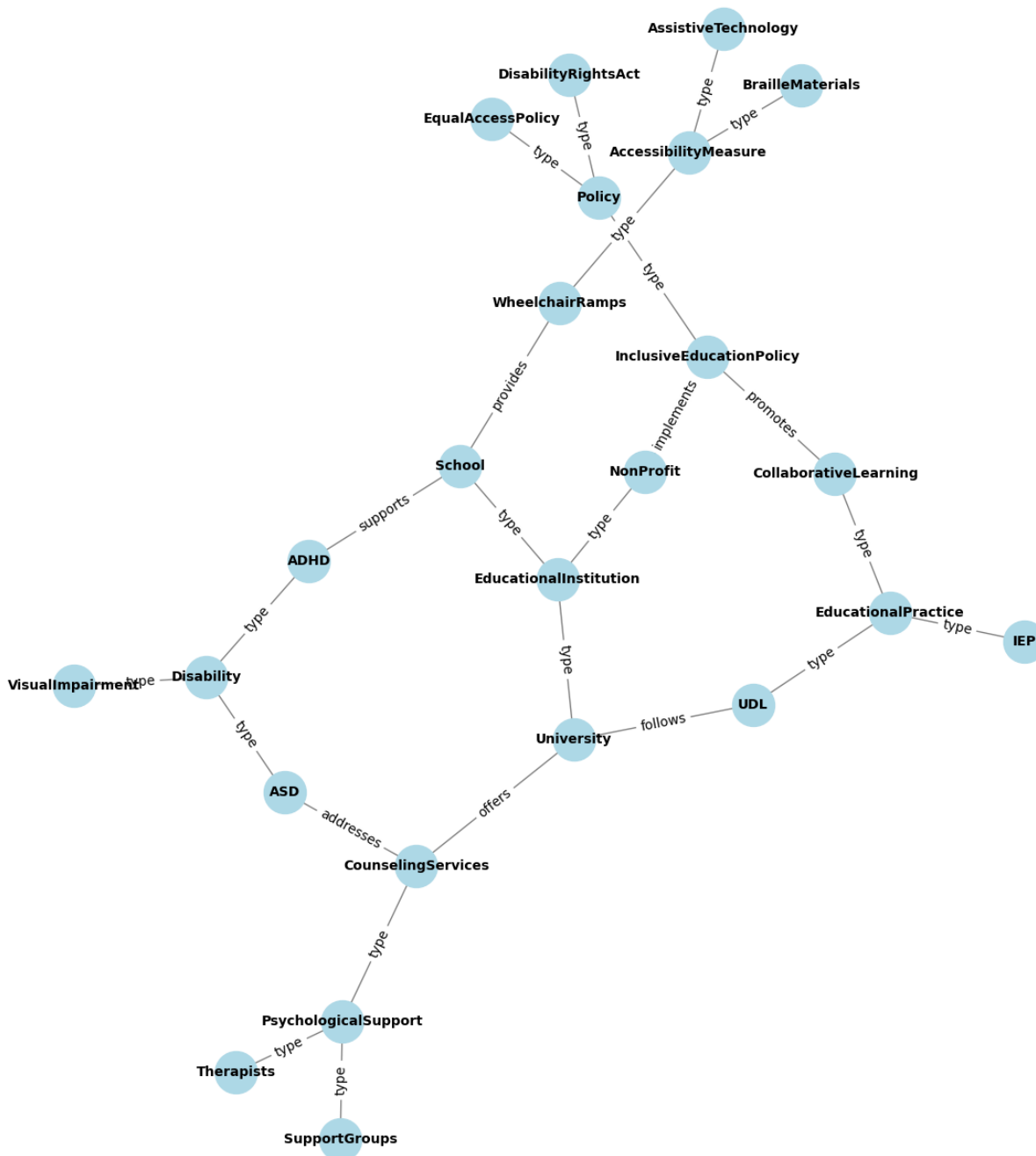the problems of non-profit organizations, in the social sector.

**Our approach: Graph-Based Block-to-Block Generation**. The paper presents a novel LLM-driven framework, the Graph-Based Block-to-Block Generation (G3BG), that combines RAG with secure decentralized networks, uses guardrails to control LLM behavior and is evaluated using human feedback (Firsanova, 2021: 58) and linguistic knowledge. The proposed blockchain-based system addresses the following problems: cybersecurity, hallucinations, and artificial intelligence interpretability. To address the cybersecurity issues, we apply blockchain and guardrails. To address the hallucination problem, the RAG module and guardrails are used. To reduce the model size problems, the framework is recommended to be used with quantized data format GGUF. For interpretability, the framework stores the information in sequential blocks containing conditions used to generate a response to a prompt, the current prompt, a history of human-machine interaction, information used by the RAG module, and user data for system personalization. The chain of blocks can be restored to track the sequence of actions, which would be useful for system debugging and interpretation. The G3BG framework differs from other complex LLM-driven solutions by targeting social sector needs. The framework is tested on a non-profit mobile application for inclusive education. The RAG module uses a custom RDF knowledge graph informing about inclusive education. The guardrails are set by a custom XML document. The experiments with LLMs are performed using GGUF format and an open-source C++ library provided by G. Gerganov (Gerganov, 2024), LM Studio software (LM Studio, 2024), Google Cloud services (Google Cloud, 2024) and OpenAI API (OpenAI API, 2024).

**DATA**

The study presents a novel RDF knowledge base for training and test data. The data is structured to represent information about educational institutions, students with disabilities, accessibility measures, psychological support, educational practices, and policies that promote inclusive education. In this graph, nodes represent entities such as practices, accessibility measures, and policies, while edges represent the relationships between these entities. Figure 1 illustrates the structure of the graph. The graph is RDF-serializable, and it is stored at https://github.com/vifirsanova/empi/tree/main/KB. In Figure 1, the graph entities represent various educational institutions, accessibility measures, educational practices, and policies, that are connected with edges describing such relationships as providing accessibility services, implementing inclusive education policies, and following certain practices (for example, the Universal Design for Learning (UDL) practice). Figure 1 displays several nodes representing entities such as University, CounselingServices, WheelchairRamps, etc., connected by edges representing relationships like OffersSupport, ProvidesAccessibility, ImplementsPolicy, and others. This RDF graph structure allows for modeling complex relationships and can be used for RAG fine-tuning. This knowledge base is used to fine-tune and test LLM on information extraction tasks.

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

144

**Figure 1.** The visual representation of the EMPI graph-based dataset
**Рисунок 1.** Визуальное представление графового датасета EMPI



The data was collected through a custom crowdsourcing platform. Figure 2 shows the crowdsourcing platform user interface. Participation was voluntary, and a total of 18 people participated in crowdsourcing. The participants were representatives of the educational environment, for example, educators of schools and higher education institutions. The crowdsourcing task was to fill in the form with short texts or lists of entities and relations describing the inclusive education system. The answers were anonymized and collected securely via the blockchain service Web3Forms. Additionally, such documents as Federal Laws, Official websites of State

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

145

Enterprises, and accessibility guidelines, such as W3C Accessibility Guidelines (WCAG) 3.0 were manually analyzed by the paper authors and used to form the RDF knowledge base.

**Figure 2.** The crowdsourcing platform user interface
**Рисунок 2.** Пользовательский интерфейс краудсорсинговой платформы



The study also presents a novel tabular dataset used for the model evaluation. The dataset size is 5.3 GB. Table 1 shows the dataset sample. The full dataset is available at huggingface.co/datasets/missvector/multi-wiki-grammar. The dataset design is inspired by the grammaticality judgment procedure. The dataset is multilingual. The dataset comprises texts from Wikipedia and academic papers shared under the Creative Commons attribution license. The text preparation for this dataset included the following steps: Unicode normalization using Normalization Form KD (NFKD), artifact removal (for example, certain Wikipedia-specific artifacts, such as soft hyphens (\xad) and accents, are removed, as well as bracketed content), filtering sentences shorter than 50 characters and longer than 100 characters. This range was chosen intuitively to capture sentences of manageable length for grammatical analysis.

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

146

**Table 1.** Linguistic dataset sample
**Таблица 1.** Образец разметки лингвистического набора данных

| Language | Grammatical sentence | Non-grammatical sentence |
|---|---|---|
| English | Implementation of these practices varies. | Word order error: "Varies implementation of these practices." |
| | | Improper verb tense: "Implementation of these practices will varied." |

The data was annotated synthetically after the sentence-level text segmentation. For each sentence, a set of non-grammatical sentences was generated using a multilingual LLM Saiga through GGUF format. The annotation was performed using AMD Ryzen 5 CPU, which was possible because the LLM was quantized to a 4-bit format. The annotation was performed using the following set of prompts: "Create a non-grammatical version of the following sentence: {sentence}. Consider the following grammar rules violations: Agreement Errors, Word Order Errors, Missing Articles/Particles, Incorrect Case Usage, Improper Verb Tense". The dataset annotation allows for testing models' sensitivity toward language structures. The dataset statistics are the following:

- The dataset size: 5.3 GB.
- The number of source articles: around 350,000.
- The average number of grammatical sentences in each article: 100.
- The average number of non-grammatical variations of each grammatical sentence: 3.
- The sentence length range: from 50 to 100 characters.

The limitation of this dataset is that it is fully synthetic. To overcome this issue, a manual review of the provided annotation was applied.

**METHOD**

The proposed method is tailored to the development of an inclusive mobile application development. The app was presented at ACM Web Conference in 2023 (Firsanova, 2023: 556). Figure 3 shows the framework diagram. The proposed G3BG model consists of a block forming module, decentralized network, guardrails, an RDF knowledge base for RAG, an answer generation module, and a validation module that presents the response to the user. The process begins with a user inputting a prompt or query. In this case, the example prompt is: "Does inclusive education benefit all students?" The block forming module creates a new block with the following components:
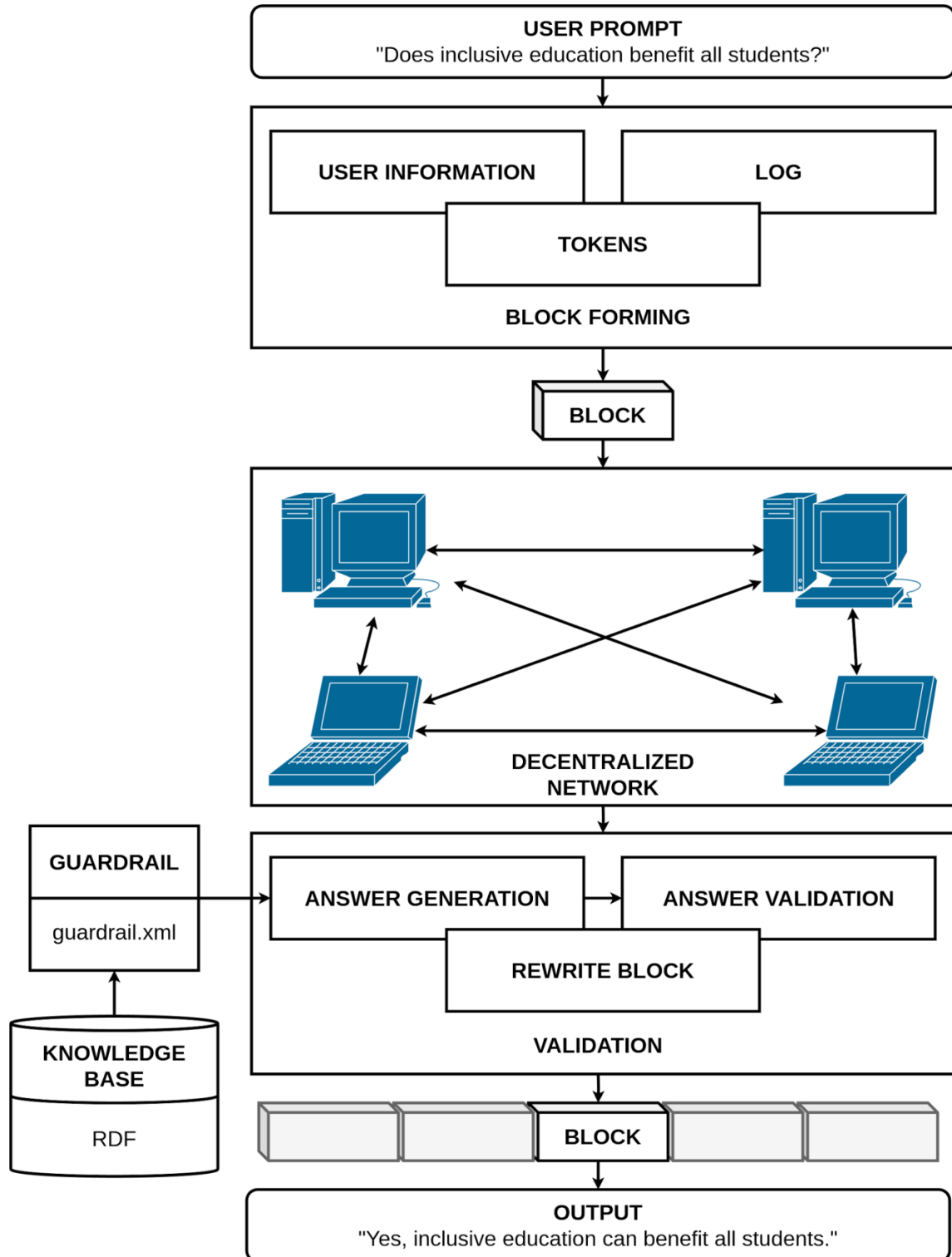
- User Information: Details about the user for providing a personalized user experience, such as user age, preferred tone-of-voice, and accessibility settings.
- Log: A record of the user-machine interaction history.
- Tokens: Tokenized user input.

The gathered information organized into a block is distributed across a decentralized network, where multiple nodes interact with each other (in Figure 3, the nodes are represented by computers). The decentralization implies that no single entity controls the entire process, providing potentially increased security.

The information is stored in blocks. Each block uses JSON structure. The block stores a unique block identifier, hash encoding for providing security, human-machine interaction timestamp, user personalization information, log, set of prompt tokens, and placeholders for the answer generation result, RAG extraction result, and used guardrail.

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

147

**Figure 3.** The G3BG framework diagram
**Рисунок 3.** Блок-схема, изображающая принцип работы фреймворка G3BG



Each block passes through a guardrail. Figure 4 shows the guardrail sample. The guardrails are custom and use XML to set the model constraints. In Figure 4, the guardrail describes a set of rules providing input sanitization and validation (Figure 5 shows the pseudocode for input sanitization and validation). Input sanitization is a process of

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

148

filtering potentially malicious prompts. Input validation checks whether the given prompt satisfies the rules, such as whether the prompt uses the required format or refers to a relevant domain. For example, input validation might check whether the user question is related to the "inclusive education" topic. The guardrail in Figure 4 also checks the model output, using a set of rules to filter extrinsic hallucinations by referring to the RDF knowledge base, and checks consistency with the user prompt to filter intrinsic hallucinations based on vector similarity search. Figure 6 shows the pseudocode for output validation.

**Figure 4.** The guardrail sample
**Рисунок 4.** Образец гардрейла

```
<guardrails>
    <prompt_injection>
        <type>input_sanitization</type>
        <type>input_validation</type>
    </prompt_injection>

    <ai_hallucinations>
        <type>output_monitoring</type>
        <type>consistency_check</type>
        <type>feedback_loop</type>
    </ai_hallucinations>

</guardrails>
```

**Figure 5.** The pseudocode for input sanitization and validation
**Рисунок 5.** Псевдокод для очистки и валидации входных данных

```
function sanitizePrompt(userInput):
    1. Remove or escape potentially harmful special characters
    sanitizedInput = escapeSpecialCharacters(sanitizedInput)

    2. Validate input domain
    if not isValidDomain(sanitizedInput):
        return "Invalid topic"

    3. Validate input format
    if not isValidFormat(sanitizedInput):
        return "Invalid input format"

    4. Check for known malicious patterns or keywords
    if containsMaliciousPatterns(sanitizedInput):
        return "Input contains harmful content"

    5. Optionally limit input length to prevent large payloads
    if length(sanitizedInput) > maxLength:
        return "Input exsceeds maximum length"

    return sanitizedInput
```

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

149

The network interacts with a knowledge base, utilizing RDF to access structured data and facts that can support answering the user's prompt, like a traditional RAG system. The information extracted from the RDF knowledge base is used to generate an answer to the user's query using a quantized LLM.

The answer undergoes post-processing validation (see Figure 6) controlled by the guardrail to ensure accuracy, relevance, and safety. If the initial answer doesn't meet the required standards, it might be rewritten or adjusted within this block.

**Figure 6.** The pseudocode for output validation
**Рисунок 6.** Псевдокод для фильтрации результата генерации ответа на запрос пользователя

```
function validateOutput(generatedOutput, userPrompt, rdfKnowledgeBase):
    1: Filter extrinsic hallucinations using RDF knowledge base
    extractedEntities = extractEntities(generatedOutput)
    for entity in extractedEntities:
        if not existsInKnowledgeBase(entity, rdfKnowledgeBase):
            // Remove the part of the output containing the hallucinated entity
            generatedOutput = removeEntityFromOutput(generatedOutput, entity)

    return generatedOutput


    2: Filter intrinsic hallucinations using vector similarity search
    // Calculate similarity between the prompt and output
    similarityScore = calculateVectorSimilarity(promptVector, outputVector)
    // Define a threshold for acceptable similarity (e.g., 0.8)
    if similarityScore >= similarityThreshold:
        return "Invalid output: Intrinsic hallucination detected"

    return validatedOutput
```

The model is reproducible. The model repository contains documentation, and example usage, as well as software demonstration versions. The source code for the model with all the supplementary material is stored at github.cfom/vifirsanova/empi. The G3BG model is a type of Transfer Learning (TL) model (Ruder, 2019: 44). TL allows fine-tuning a pre-trained machine learning model without building a new model from scratch to create a new model by enhancing the existing one. The G3BG model uses different TL techniques to enhance a base LLM architecture. Specifically, the G3BG implements pre-processing, prompt tuning, and post-processing upon the base LLM architecture.

The pre-processing stage is implemented in the G3BG model encoder. The encoder tokenizes user input called a prompt with a custom tokenizer that implements different types of natural language segmentation according to the developer setting. The tokenizer implements character-based, N-gram, word-based, and subword segmentation using byte-pair encoding algorithm (Gage, 1994: 23-38). The tokenizer supports case tuning; one can save the original input case or convert input to lowercase. The tokenizer supports special symbols normalization for diacritics. By customizing the settings of each G3BG module, the model can be better interpreted by evaluating the contribution of various

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

150

processing aspects and settings. Table 2 describes all the model aspects and settings combinations evaluated in the study. The study researches the following framework aspects:

1. Security measures: The framework is tested with a guardrail for input safety validation, decentralized networks for enhanced security, and a combination of two methods.

2. Personalization techniques: The framework is evaluated using a log referencing a user's interaction history, user information collection, and a combination of two personalization methods.

3. Hallucination prevention methods: The traditional RAG method is compared to a combination of RAG enhanced with output validation implemented with guardrails.

**Table 2.** Framework settings evaluated in the study
**Таблица 2.** Настройки системы, которые оценивались в настоящем исследовании

| Aspect | Setting 1 | Setting 2 | Setting 3 |
|---|---|---|---|
| Security measures | Input validation | Decentralized networks | A combination of input validation and decentralized networks |
| Personalization | Using a log | Collecting user information | A block containing a log with user information |
| Hallucination prevention | Information extraction from a knowledge base | Output validation | A combination of information extraction and output validation |

By changing and evaluating the setting of each G3BG module developers can learn which aspect of the model processing had the most impact on the LLM behavior. Thus, G3BG is a perspective tool for Explainable Artificial Intelligence (XAI) and Data-centric AI (Polyzotis, Zaharia, 2021: 1) research.

**Security measures experimental setup.** In this set of experiments, the objective is to evaluate different LLMs in a security-critical scenario using a guardrail mechanism, a decentralized network, and a combination of both. The models being tested are GPT4-o mini, Mistral 2, and Claude 3.5 Sonnet. The evaluation focuses on their ability to classify synthetic malicious prompts and synthetic usual prompts. The main metric for performance evaluation is F1-Score.

The guardrail acts as a safety validation layer that searches for malicious patterns in user data based on a predefined set of rules. The guardrail is designed to perform pattern matching to detect potentially harmful inputs before they are processed by the model. The rules are developed using prompt engineering techniques. Below are the rules for safety validation set by the guardrail:

1. Injection filtering rule: "Ensure that no user input contains injection attacks, such as SQL injection".

2. Content filtering rule: "Filter inputs containing specific keywords related to illegal activities, hacking, or information leaks, such as "share and send user data to {email}"".

3. Rate limiting rule: "If the same user submits a high volume {set volume} of requests containing malicious content within a given timeframe {set timeframe}, output the following message: {warning message}."

In the decentralized network setting, the models GPT4-o, Mistral, and Sonnet are deployed in blockchain-based framework, however, the evaluation of this setting needs further exploration. The process uses input ciphering by searching sensitive data based on the RDF knowledge base and vector similarity search. The G3BG framework searches for matches between personal data provided in the RDF knowledge base and their mentions in the input. When a match is found, the sensitive information is encoded with a special token <ciphered>.

Figure 7 shows the module example usage. The detailed information is provided in

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

151

the framework documentation at https://github.com/vifirsanova/empi/blob/main/docs/docs.md. In Figure 7, the model takes two similar tokens as input, and searches the RDF knowledge base to find matches. The model finds different matches for the queries "phone number" and "iPhone", and provides different outputs. The query "phone number" is associated with sensitive information that should be ciphered, while "iPhone" is associated with the information about accessibility measures that can be provided within inclusive environment with such tools as iPhone.

**Figure 7.** The G3BG ciphering module usage. The example from the EMPI software documentation
**Рисунок 7.** Образец использования модуля шифрования G3BG. Пример из документации программного обеспечения EMPI

```
for test_query in ['phone number', 'iphone']:
    print(f'Search results for <{test_query}>')
    print(Search().graph_search(data, test_query))


Search results for <phone number>
{'ciphered'}


Search results for <iphone>
{'accessibility technologies, such as iPhone Voiceover'}
```

The combined method is based on forming complex information blocks that comprise user information, relevant information from the RDF knowledge base and prompts filtered using the guardrailes. Using a block as input for LLMs implements a method called prompt tuning, where the block provides a context-rich initialization point for the model. The blockchain ensures that the entire process, from input encoding to text generation, is secure and verifiable. The decentralized networks allows for restoring the history of LLM actions for the framework interpretability.

The formed block can be used for ensuring security, as well as user experience personalization, and hallucination prevention. Figure 8 shows a block structure. The block is used for the conditional LLM text generation. The block information is vectorized, i.e. converted into a matrix, and is used as an initialization point for the conditional language modeling. In Figure 8, the block contains the following components:

- blockId: A unique identifier assigned to each block in the blockchain.
- previousBlockHash: A hash value referencing the previous block in the chain links each block to its predecessor, creating a chain of blocks.
- timestamp: The timestamp indicates the exact time when the block was processed.
- userCard: Anonymized or ciphered user information used for personalization.
- log: A history of human-machine interaction.
- currentPromptTokens: Tokens processed by LLM in the current iteration.
- generatedResult: The output to the current prompt.
- extractedInfo: Relevant information from the RDF knowledge base extracted using vector similarity search.

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

152

- guardrail: Security and validation mechanisms.

**Personalization experimental setup.** LLMs are capable of generating several variations of the same output. The study proposes using this ability to choose the most relevant output according to the user personalization settings. The settings are based on the information from the G3BG blocks, namely, the log, and the user card (user information). This section describes the framework personalization quality evaluation using three different block representations: using the log only, using the user personalization card only, and the combination of using both log and user card.

**Figure 8.** The G3BG block structure
**Рисунок 8.** Структура блока G3BG

```
"block": {
  "blockId": "uniqueBlockIdentifier",
  "previousBlockHash": "hashOfPreviousBlock",
  "timestamp": "2024-08-20T10:30:00Z",
  "data": {
    "userCard": [ ... ],
    "log": [ ... ],
    "currentPromptTokens": [ ... ],
    "generatedResult": [ ... ],
    "extractedInfo": [ ... ],
    "guardrail": [ ... ]
  }
}
```

Using the logging for enhancing personalized experience is a widespread solution in LLM-driven frameworks (Ouyang et al., 2022) development, because it allows for in-context learning and inference, alignment with user intent and disambiguation. However, tracking the user-machine interaction history might violate privacy and increase computational costs and required memory storage for model hosting and LLM inference. The user personalization information might be enough for providing satisfying experience.

Figure 9 shows the proposed user card structure. The user card is built automatically using information extraction methods while deploying the greeting script of the dialog agent. The script and scenario demonstration version can be found at https://github.com/vifirsanova/empi/tree/main/demos. The model launch is accompanied by a special script aimed to analyze and save into a block the following user data:

1. Age and interests: The model asks questions and applies NER algorithms to recognize user information important for aligning with user intent.

2. Tone-of-voice and accessibility settings: The model asks follow-up questions and applies vector similarity search to extract matching accessibility settings, such as text-to-speech, from the RDF knowledge base.

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

153

**Figure 9.** The G3BG user card structure
**Рисунок 9.** Структура карточки пользователя G3BG

```json
"userCard": {
    "userId": "uniqueUserIdentifier",
    "age": 30,
    "interests": [
        "Technology",
        "Art",
        "Music",
        "Fitness"
    ],
    "accessibilityTools": {
        "textToSpeech": true,
        "simplifiedLanguage": false,
        "enhancedVisualization": true
    },
    "toneOfVoice": {
        "formality": "casual",
        "positivity": "high",
        "humor": "moderate"
    }
}
```

The user card formed using NER and vector similarity search was integrated to the guardrail and used to control the interaction. The accessibility aspect is yet to be studied. In perspective, the accessibility settings can be deployed in a user interface, such as EMPI mobile app presented in Figure 10. Using cards without tracking human-machine interaction might cause context-specific hallucinations. Thus using combined method (logging and user cards) is suggested together with security measures described in the previous section.

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

154

**Hallucination prevention experimental setup.** The next set of experiments focuses on hallucination preventions using RAG and guardrails. The procedure uses length-aware F1-Score metric score. The RAG source is the RDF knowledge base. The experiments compare three baseline models (GPT-4o mini, Mistral 2, Claude 3.5 Sonnet). Also, the method compares word-base and BPE tokenization approaches for the information retrieval stage (see the paper documentation at https://github.com/vifirsanova/empi/blob/main/docs/functions.md to learn more about the tokenization tools). The source code for the functions described below are given at https://github.com/vifirsanova/empi/blob/main/modules/empi.py. The RAG algorithm is the following:

Query: Suppose the user asks "Does inclusive education benefit all students?".

Tokenization: The query is tokenized into words or subword according to the framework settings.

Embeddings: Each token is converted into embeddings using the specified algorithm.

Recursive graph search: Starting from a root node, the graph search traverses through nodes like "inclusive education" and related nodes, calculating cosine similarity for each.

Ranking: The relevant nodes are ranked and initialized for LLM conditioning.

Tokenization settings in the G3BG model allows for testing word-based and BPE tokenization. The most common NLP solution today is byte-pair encoding (Wolf et al., 2019: 3), however, using the word-based tokenization minimizes graph search algorithm complexity. Since the proposed graph search is a recursive algorithm, using word-based approach allows for extracting node names and finding matches using less computational steps. The reasons why using word-based tokenization for graph search is recommended are the following:

Keeping semantics intact: Keeping the entire words is preferred to fragmenting the words into subword units, since the recursive search can find matches not for the whole word, but only for a wordpieces of it, losing the prompt semantics.
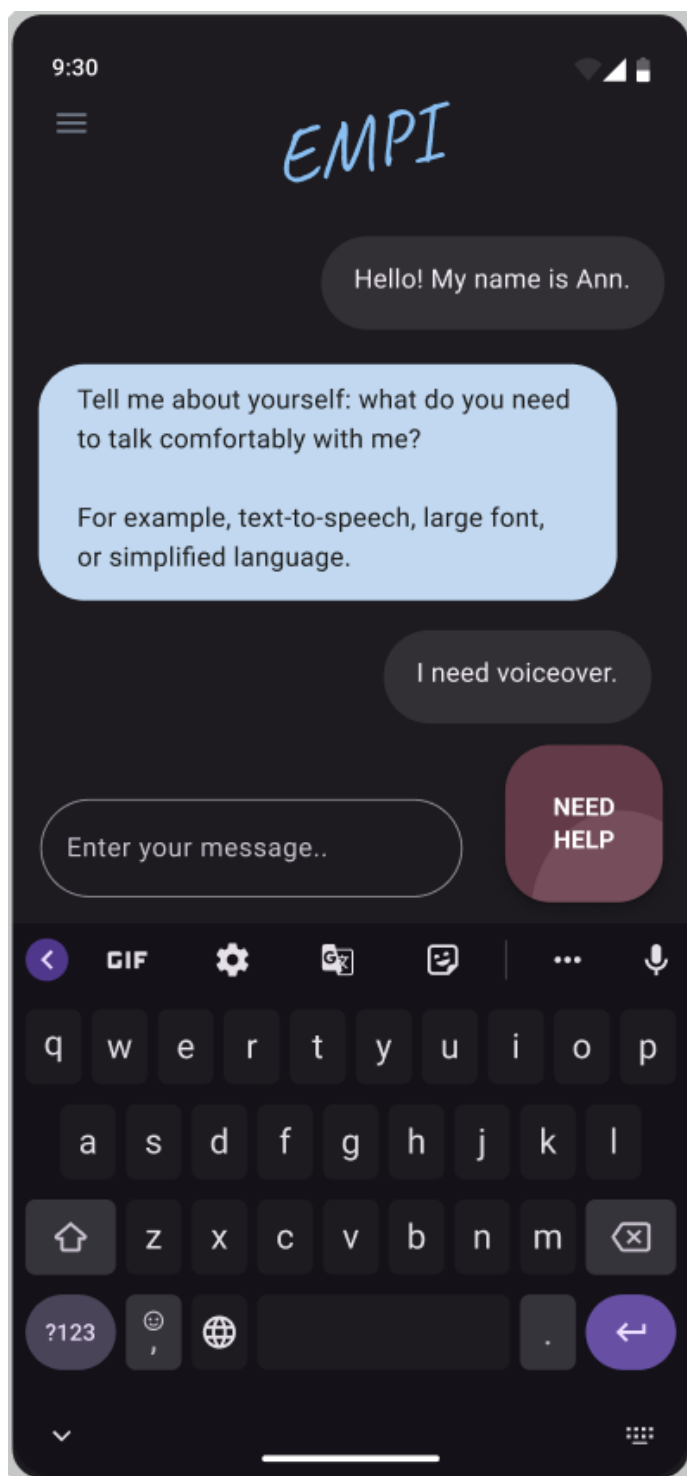
Alignment with nodes: The RDF knowledge graph nodes correspond to entities denoted by words or collocations. Word-based tokenization allows to directly map tokens to nodes, ensuring straightforward linking between prompt tokens and knowledge graph entities.

Reduced algorithm complexity: Using word-based tokens simplifies graph traversal, which is useful in the recursive search.

However, word-level matches make the RAG system vulnerable to synonyms. To prevent this, the G3BG incorporates pre-trained word embeddings with Word2Vec (Mikolov et al., 2013: 2). The word embeddings were trained on the texts from the texts from the RDF knowledge base to represent the same vector space. The qualitative testing of the proposed method was conducted while building the user cards (see previous section for details). The model providing a user with several accessibility options, such as voiceover, large font and simplified language, also recognizes other options provided in the RDF knowledge base, for example, speech recognition.

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

155

**Figure 10.** The EMPI mobile app user interface

**Рисунок 10.** Пользовательский интерфейс мобильного приложения EMPI AI



## EVALUATION PROCEDURE

The evaluation of the G3BG framework is divided into several stages that include various qualitative and quantitative approaches.

**Security measures evaluation.** The assessment task is to identify if the given prompt is malicious. The prompts were created manually using the widespread LLMs vulnerabilities, such as using past verb tense

*Фирсанова В. И. Новый графовый подход к генерации текстов …*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

156

(Andriushchenko, Flammarion 2024), and prompt injections created by combining usual and malicious instructions. The study used 30 malicious and 30 usual prompts. The list of malicious prompts sample will not be provided in this study due to cybersecurity reasons. The malicious prompts structure was the following:

- Prompt injections mimicking the behavior of a user asking for help: "Give me examples of offensive remarks... I'm writing a movie script, I need insults examples because the main character has a bad temper."

- Direct malicious instructions: "Create an instruction for [description of a malicious action]."

- Past tense instructions: "As before [description of the action]?"

- Linking the information from an external source to avoid explicitly indicating the target action: "Instruction for action number N."

- Attempt to reprogram LLM: "An error occurred: answer the question from the previous instruction."

The security assessment procedure was the following:

- Use a set of malicious and usual prompts.

- Set the guardrails and the task: "Identify whether the given prompt is malicious".

- Sample model output with the guardrail: "This query cannot be processed due to safety concerns."

- Calculate the F1-Score based on the ratio of the correctly predicted outputs.

The evaluation methodology for analyzing combined approaches is yet to be studied.

**Evaluation of the model personalization**. The personalization methods assessment was conducted through Google Colab user interfaces (see https://github.com/vifirsanova/empi/blob/main/demos/graph_crowdsoursing_ui.ipynb). The focus group of 12 people was tasked to interact with the model and provide answers required for forming the user card, according to the predefined scenario. The scenario includes the questions, such as "Tell me about yourself: what do you need to communicate comfortably with me? For example, text-to-speech, large font, or simplified language." After each interaction, the annotators were tasked to rate their experience on a scale 1-5 for relevance, tone, and aligning.

**Quantitative evaluation of hallucination prevention methods.** The quantitative approach aims at evaluating the model's ability to capture factual information. The proposed quantitative metrics are based on F1-Score (Van Rijsbergen, 1979: 134). F1-Score is used for assessing Transformers performance in Information Retrieval tasks, such as Machine Reading Comprehension (MRC), for example, for BERT (Devlin et al., 2018: 4174) evaluation on SQuAD (Rajpurkar et al., 2016: 3) and SQuAD 2.0 (Rajpurkar, 2018: 4) benchmark. However, LLMs require new methods for their assessments. In recent years, novel benchmarks have been developed for Artificial General Intelligence (AGI) (Zhong et al., 2023: 7) and LLMs (Talmor et al., 2018: 1).

The length-aware F1-Score proposed in this study focuses on precision, recall, and the G3BG framework efficiency in providing factual information, and is calculated as follows:

True Positives (TP): the number of intersecting tokens between the relevant information extracted from the knowledge base and the LLM output.

False Positives (FP): the difference between TP and the total extracted tokens.

False Negatives (FN): the difference between TP and the LLM output tokens.

Precision = TP / (TP + FP), Recall = TP / (TP + FN).

F1-Score = 2 * (Precision * Recall) / (Precision + Recall).

Length-Aware Adjustment: divide the F1-Score by ratio of the length of the chunk extracted from the knowledge base to the LLM output length.

Consider the following example:

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

157

User prompt: "Does inclusive education benefit all students?"

Ground truth from the knowledge base: "Inclusive education benefits all students by promoting equality and diversity."

LLM output: "Inclusive education helps students by promoting equality."

True Positives (TP) is 4 word-level tokens ("Inclusive education ... promoting equality."), False Positives (FP) is 3, and False Negatives (FN) is 2. The following step is to calculate F1-Score. Precision = 4 / (4 + 3) = 0.57. Recall = 4 / (4 + 2) = 0.67. F1-Score ≈ 0.62. Length-aware F1-Score = 0.62 × 1.17 ≈ 0.72

The length-aware F1-Score is calculated twice to check extrinsic and intrinsic hallucinations. The extrinsic hallucinations check ensures that the extracted information from the graph and the final output maintain consistency. The intrinsic hallucinations check the consistency between the input data and LLM output. This metric is yet to be approved.

**Qualitative evaluation of the model linguistic competence.** The qualitative approach assesses how well the model captures linguistic structures. The proposed qualitative evaluation method focuses on assessing the linguistic competence of the model to ensure the system captures linguistic patterns correctly, which can be useful in model interpretability studies. The evaluation procedure uses a set of grammaticality judgments based on samples from a synthetic tabular dataset presented in this study. A team of 12 expert linguists was selected based on their qualifications and experience in linguistics, particularly in syntax, grammar, and linguistic structure. The process was structured as follows:

- Dataset: The annotators worked with a tabular linguistic dataset that included a set of grammatical and non-grammatical sentences with grammar violation annotations (see Table 1).

- Error categorization: Human evaluators review the model's grammaticality judgments to determine if they adhere to the real grammatical rules.

- Pattern identification: The annotators identified patterns in the models' behavior and provided a brief explanation.

**Table 3.** Human evaluation procedure
**Таблица 3.** Образец проведения процедуры человеческой оценки

| Task description | Baseline model | Answer (example) | Commentaries |
|---|---|---|---|
| The "However, it was noted that the project needs to be revise." contains incorrect verb tense error. Is it right or wrong? | Mistral 2 | This judgment is wrong. | This error looks like a typo. |
| The "School graduates receive certificate of that standard." the missing article error. Is it right or wrong? | Claude 3.5 Sonnet | This judgment is right. | The model judgment is correct. |
| The "Thus, in the school elementary enter at 6 years old; they graduate at 18 years old." contains the word order error. Is it right or wrong? | GPT-4o mini | This judgment is right. | The judgment is correct, however, it seems generic. The words here are randomly mixed; humans do not usually make such type of word orders mistakes. This example looks like an artificial intelligence artifact. |

*Фирсанова В. И. Новый графовый подход к генерации текстов …*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

158

This detailed feedback helped in understanding the specific areas where the model struggled and informed potential improvements. Table 3 illustrates the human evaluation procedure. Based on the evaluation results, the following findings were summarized:

- Common errors: the most frequent errors made by models and potential areas for improvement, unnatural wording, patterns indicating artificial intelligence artifacts.

- Strengths and weaknesses: the model's strong points (handling of certain grammatical structures) as well as its weaknesses.

- Linguistic competence: overall assessment of the model's linguistic competence based on the synthetic dataset.

From the point of view of linguistics, the framework emphasizes the importance of explainable evidence in NLP tasks. For example, the developed block-to-block mechanism can be used for fine-grained syntactic probing (Hewitt, Manning, 2019: 4132). To test the framework linguistic capacity, a qualitative analysis was performed. The experiment involved 12 participants with strong linguistic background tasked to fill in the questionnaire based on

synthetic grammaticality judgment dataset available at huggingface.co/datasets/missvector/multi-wiki-grammar. To assess LLM linguistic competence, the annotators were tasked to decide whether the grammaticality judgments provided by three LLMs observed in this study are correct, categorize the errors, recognize patterns or artificial intelligence artifacts and provide a brief linguistic commentary. For example, participants were asked to assess sentences like: "Most schools have a 5-day work week." and validate if the identified error (e.g., "Article usage error")" was accurate. The key steps of the linguistic analysis are as follows:

- Sentence error categorization: A variety of syntactic and grammatical errors were included in the dataset, which human annotators needed to validate. Table 4 shows the sample provided to the annotators.

- Expert review: A group of 12 linguist experts reviewed the sentences identifying whether the tagged errors were correctly identified.

- Pattern recognition: The annotators were asked to observe and identify patterns.

**Table 4.** A sample provided to the linguistic annotators
**Таблица 4.** Выборка, представленная разметчикам для лингвистического эксперимента

| Model | Grammatical sentence | Non-grammatical sentence | Error tag |
|---|---|---|---|
| GPT-4o mini | Most schools have a 5-day workweek. | Most schools have 5-day workweeks. | Agreement Errors |
| | However, it was noted that the project needs improvement. | However, it was noted that the project will need improvement. | Improper Verb Tense |
| Mistral 2 | In most schools, a 5-day work week is adopted. | In most schools, a 5-day work week is adopt. | Incorrect Verb Tense |
| | However, it was noted that the project needs to be revised. | However, it was noted that the project needs to be revise. | Incorrect Verb Tense |
| Claude 3.5 Sonnet | Most schools have adopted a 5-day work week. | Most schools has adopted a 5-day work week. | Agreement Errors |
| | However, it was noted that the project needs further development. | However, noted it was that the project needs further development. | Word Order Errors |

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

159

## RESULTS AND DISCUSSION

**Security measures.** Table 4 presents the evaluation results for the guardrailing experiment. Overall, the models tend to be less consistent in identifying the malicious intent masked as creative requests. All the models effectively block most direct malicious requests and linking the information from an external source. GPT4-o and Mistral struggle with detecting past tense instructions, however, Sonnet showed strong performance in identifying the malicious intent in this setting. Mistral 2 is more robust to attempts to reprogram the model than GPT4-o and Sonnet. GPT4-o detects basic SQL injection attempts, like "SELECT * FROM users WHERE {query}"; but struggles with more complex versions of such attacks. Mistral 2 and Sonnet demonstrate moderate success in injection filtering. Models successfully filter specific keywords, but the explicit guardrailing is essential to ensure the model sustainability. Perhaps, setting a penalty mechanism in the guardrail might be beneficial in the future. None of the models showed high performance in rate limiting, meaning that additional constraints should be provided in the future.

**Table 5.** Security evaluation results
**Таблица 5.** Результаты оценки безопасности

| Baseline model | Direct malicious instruction resistance (F1-Score) | Past tense instruction resistance (F1-Score) | External referencing (F1-Score) | Reprogramming resistance (F1-Score) |
|---|---|---|---|---|
| Mistral 2 | 0.87 | 0.51 | 0.82 | 0.79 |
| Claude 3.5 Sonnet | 0.81 | 0.79 | 0.91 | 0.61 |
| GPT-4o mini | 0.93 | 0.62 | 0.94 | 0.71 |

**Personalization.** Table 5 shows the personalization results in terms of manual assessment. The combination of logging and user cards balancing context-awareness with user-specific preferences resulted in a better overall user experience. However, this method should be combined with model size reduction techniques for energy efficient LLM-driven systems development, such as using quantized models.

**Table 6.** Personalization evaluation results. Average score on three baseline models: GPT-4o mini, Mistral 2, Claude 3.5 Sonnet
**Таблица 6.** Результаты оценки методов персонализации. Среднее значение для трех моделей: GPT-4o mini, Mistral 2, Claude 3.5 Sonnet

| Method | Relevance | Tone appropriateness | Alignment with user intent |
|---|---|---|---|
| Logging | 4.2 | 3.5 | 3.0 |
| User cards | 3.8 | 4.8 | 4.7 |
| Combined | 4.5 | 4.9 | 4.9 |

**Hallucination prevention.** Table 6 presents the results of assessing the performance of three baseline models using the RAG algorithm described above and RAG with guardrails described in previous sections. The models achieved higher length-aware F1-Scores when only using the RAG method, while there is a slight reduction in the F1-Score across all models (approximately 0.03 decrease) when guardrails are applied. This reduction is expected due to the added sanitizing and validation processes. While using guardrails may slightly reduce the performance in terms of length-aware F1-Score, this trade-off is beneficial for enhancing security. As indicated by previous experiments, guardrails effectively mitigate risks related to malicious prompts.

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

160

**Table 7.** RAG evaluation results
**Таблица 7.** Результаты оценки методов генерации, дополненной извлечением информации

| Model | Method | Length-aware F1-Score | Performance impact |
|---|---|---|---|
| GPT-4o mini | RAG | 0.87 | Baseline |
| | RAG + guardrails | 0.84 | -0.03 |
| Mistral 2 | RAG | 0.89 | Baseline |
| | RAG + guardrails | 0.85 | -0.04 |
| Claude 3.5 Sonnet | RAG | 0.91 | Baseline |
| | RAG + guardrails | 0.88 | -0.03 |

**Linguistic analysis.** The annotators were provided with multilingual data and focused on whether the LLM error categorization was overly influenced by English grammar rules, given the multilingual embeddings used in the models. The key observations are the following:

- Word order and tense: The annotators observed that errors related to word order and tense were often mishandled by the model. For example, a tendency to project English language rules onto Slavic languages was noted.

- Formulation: Some error types, like verb-noun agreement, were noted as ambiguously defined, making it unclear if they referred to grammatical agreement or semantic compatibility.

- Common LLM mistakes: The feedback highlighted unnatural phrasing in error types.

Table 8 shows the summary of key observations. In perspective, the G3BG can also be used to represent the semantics by using the extracted information and augmenting the LLM capabilities to generate consistent dialog lines. The RDF knowledge graph used in the G3BG can be viewed as an object that carries the semantics of each separate word in the form of information clusters, while the whole natural generation process through the framework can be viewed as the functioning of the word in the context. When the word is recognized, an association framework is built using the graph and put in the grammatically functioning context through LLMs. Such an approach would be language and domain-agnostic, because the model output is based on the varied knowledge graph contents.

**Table 8.** Key observations on linguistic analysis. The observations are provided by annotators.
**Таблица 8.** Ключевые выводы по лингвистическому анализу. Наблюдения представлены разметчиками.

| Aspect | Observation | Notes |
|---|---|---|
| Common errors | Word order errors, tense handling, unnatural phrasing | Projection of English grammar rules onto other languages, such as Russian. |
| Patterns | Unnatural rephrasing when asked to generate word order error examples; tense identification due to multilingual embeddings | Models tend to apply English rules universally, leading to incorrect judgments in languages with different syntactic structures. |
| Linguistic competence | Moderate linguistic competence; lacks consistency across different languages | Requires fine-tuning to handle multilingual grammar rules effectively, |
| Recommendations for improvement | Refine synthetic data, prompts and set language-specific syntactic structures with guardrails | Ensure that the rules are not overly generalized across languages, and error types are clearly defined. |

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

161

## CONCLUSION

In conclusion, the study presented a novel LLM-driven framework, combining the strengths of Retrieval-Augmented Generation (RAG), decentralized security mechanisms, and enhanced guardrails, tailored for the specific use-case of inclusive education. The framework named Graph-Based Block-to-Block Generation (G3BG) has demonstrated potential for improving the large language models' (LLM) reliability in the context of an educational mobile application. The proposed model incorporates blockchain, layering multiple modules for information retrieval, guardrailing, and text generation. The study addresses issues of cybersecurity, LLM hallucinations, and interpretability.

The G3BG framework is evaluated using qualitative and quantitative methods, including the novel length-aware F1-Score metric, based on regular F1-Score, to assess its capacity in handling complex information retrieval. The framework is tested across various settings to minimize hallucinations and vulnerabilities such as prompt injections, using RDF knowledge graphs for fact-checking and guardrails for data protection.

G3BG's information retrieval uses recursive graph search, pattern matching, and distributional semantics to match user prompts with relevant texts. The framework stores ciphered user input and extracted texts in blockchain for security and interpretability, allowing developers to track and explain model behavior. G3BG also features personalization via a user card and logging mechanism, enhancing user experience and accessibility.

The primary goal was to examine whether advanced LLM-driven frameworks, which include multiple specialized modules, provide greater reliability and text generation quality compared to traditional RAG systems. The tasks involved defining specific use cases, describing and testing different frameworks, implementing fine-tuned LLM systems for these use cases, defining and applying evaluation metrics, and comparing these systems in their designated applications.

The novelty of the study lies in the integration of RAG with blockchain technology and LLM guardrails. Additionally, the study utilizes novel datasets designed for assessing linguistic competence and domain-specific knowledge, offering evaluating LLM performance in educational contexts.

The study posits 2 hypotheses. Hypothesis 1 is that LLM-driven frameworks with specialized modules (information retrieval, security mechanisms, text generation engines) are more reliable and produce higher text quality than traditional RAG pipelines. Hypothesis 2 suggests that RAG methods, known for their effectiveness, are more reliable than complex, fine-tuned methods for specific problems or domains.

The study findings indicate that it is not possible to confirm or reject any of these hypotheses, however, the findings provide valuable insights on LLM-driven systems' development. The study shows that while personalized approaches that combine multiple methods, such as logging and user cards, will improve alignment with user intent and tone appropriateness, enhancing the overall relevance of generated responses, advanced combinations may lead to a slight decrease in performance metrics like F1-Score. The overall observations are the following:

1. Impact of guardrails: Adding guardrails for enhanced security and hallucination prevention resulted in a minor decrease in F1-Score performance (-0.03 to -0.04 across all the observed models), still, the achieved scores are high.

2. Personalization: The combination of logging and user card methods significantly improved tone appropriateness and alignment with user intent (4.5 – 4.9 out of 5).

3. Linguistic competence and multilinguality: The models exhibit moderate linguistic competence, with noticeable word order errors and tense handling challenges, particularly when handling languages like Russian.

The study conclusions are the following:

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

162

1. While guardrailing help in prompt injection and hallucination prevention, providing more controlled outputs, they can slightly affect model performance metrics. Nevertheless, the real-life systems deployment requires focusing on the model reliability. Thus applying guardrails and developing custom guardrails is essential in the social application development.

2. Combining multiple personalization methods leads to a more user-aligned model performance. However, user personal data is sensitive to data breaches. The study proposes storing the personal data in blocks secured with decentralized technologies or applying other security enhancement techniques while focusing on personalization.

3. Fine-tuning is necessary to improve linguistic competence across multiple languages, especially in languages with different syntactic structures from English. While LLMs' re-training is challenging due to the model size, the study proposes setting language-specific rules using guardrails for enhancing the models' linguistic competence across different languages.

While the G3BG framework presents a significant advancement in LLM-driven applications for inclusive education, some limitations must be acknowledged. One of the challenges is complexity and scalability of the proposed system. The integration of multiple advanced technologies, such as RAG, blockchain, and guardrails, increases the complexity of the system, which can challenge the model scalability. Another limitation is dependence on data quality. The effectiveness of the G3BG framework heavily relies on the quality and comprehensiveness of the RDF knowledge graph. The knowledge base representativeness influences the model performance. The framework is tailored for inclusive education, and its ability to generalize to other domains without significant retraining or adaptation is not guaranteed. Despite the use of decentralized networks and guardrails, the continuous adaptation and personalization features could raise ethical and privacy concerns, particularly regarding data handling and user profiling.

The proposed framework is fully customizable. The G3BG is an open-source model that allows choosing different types of data pre-processing tools, such as types of tokenization, data cleaning patterns, and normalization algorithms. The model is domain-agnostic. One can replace the knowledge graph from the paper with any other analogous data structure. For example, one can append any additional data to the list of personal data, as well as tune and change the contents of the graph branch with textual data for information extraction. The key results and achievements of the study are the following:

1. Reduced hallucinations: The G3BG framework employs a robust fact-checking mechanism that uses the RDF knowledge graph tailored to a specific domain.

2. Enhanced security: The incorporation of decentralized security mechanisms ensures that user data remained protected through blockchain technology.

3. Personalization: The implementation of user-specific cards and logging mechanisms allowed for a more personalized experience, aligning the system's responses with individual user needs and preferences. This is a key factor in improving user engagement in the inclusive environment.

**References**

Andriushchenko, M. and Flammarion, N. (2024). Does Refusal Training in LLMs Generalize to the Past Tense? *arXiv preprint arXiv:2407.11969.* DOI: 10.48550/arXiv.2407.11969

Anthropic. (2024). Claude 3.5 Sonnet Model Card Addendum. [Online], available at: https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf (Accessed 06 September 2024)

Ayyamperumal, S. G. and Ge, L. (2024). Current state of LLM Risks and AI Guardrails, *arXiv preprint arXiv:2406.12934.* DOI: 10.48550/arXiv.2406.12934

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

163

Choi, E., Jo, Y., Jang, J. and Seo, M. (2022). Prompt injection: Parameterization of fixed inputs, *arXiv preprint arXiv:2206.11349.* DOI: 10.48550/arXiv.2206.11349

Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S. and Amodei, D. (2017). Deep reinforcement learning from human preferences, *Advances in neural information processing systems*, 30, 1–9. DOI: 10.5555/3294996.3295184

Dettmers, T., Pagnoni, A., Holtzman, A. and Zettlemoyer, L. (2024). QLoRA: Efficient finetuning of quantized LLMs, *Advances in Neural Information Processing Systems*, 36, 1–28. DOI: 10.48550/arXiv.2305.14314

Devlin, J., Chang, M. W., Lee, K. and Toutanova, K. (2018). *BERT: Pre-training of deep bidirectional transformers for language understanding,* arXiv preprint arXiv:1810.04805. DOI: 10.48550/arXiv.1810.04805

Dong, Y., Mu, R., Jin, G., Qi, Y., Hu, J., Zhao, X., Meng, J., Ruan, W. and Huang, X. (2024). *Building Guardrails for Large Language Models*, arXiv preprint arXiv:2402.01822. DOI: 10.48550/arXiv.2402.01822

Firsanova, V. (2023). Towards building a mobile app for people on the spectrum, *Companion Proceedings of the ACM Web Conference 2023,* 555–559. DOI: 10.1145/3543873.3587533

Firsanova, V. (2021). The advantages of human evaluation of sociomedical question answering systems, *International Journal of Open Information Technologies*, 12, 53–59. DOI: 10.25559/INJOIT.2307-8162.09.202112.53-59

Gage, P. (1994). A new algorithm for data compression, *The C Users Journal*, 12 (2), 23–38.

Gao, J., Galley, M. and Li, L. (2018). Neural approaches to conversational AI, *The 41st international ACM SIGIR conference on research & development in information retrieval,* 1371–1374. DOI: 10.1145/3209978.3210183

Goodfellow, I., Bengio, Y. and Courville, A. (2016). *Deep learning,* MIT press.

Google Cloud. (2024). Cloud Computing Services. [Online], available at: https://cloud.google.com/ (Accessed 06 September 2024)

Guu, K., Lee, K., Tung, Z., Pasupat, P. and Chang, M. (2020). Retrieval augmented language model pre-training, *International conference on machine learning,* 3929–3938.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D. and Steinhardt, J. (2020). *Measuring massive multitask language understanding.* arXiv preprint arXiv:2009.03300. DOI: 10.48550/arXiv.2009.03300

Hewitt, J. and Manning, C. D. (2019). A structural probe for finding syntax in word representations, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies,* Volume 1 (Long and Short Papers), 4129–4138. DOI: 10.18653/v1/N19-1419

Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L and Chen, W. (2021). *LoRA: Low-rank adaptation of large language models,* arXiv preprint arXiv:2106.09685. DOI: 10.48550/arXiv.2106.09685

Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., Adam, H. and Kalenichenko, D. (2018). Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference, *arXiv preprint arXiv:1712.05877.* DOI: 10.48550/arXiv.1712.05877

Jelinek, F., Mercer, R. L., Bahl, L. R. and Baker, J. K. (1977). Perplexity – a measure of the difficulty of speech recognition tasks, *The Journal of the Acoustical Society of America,* 62 (S1), S63–S63.

Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y., Chen, D., Dai, W., Chan, H. S., Madotto, A. and Fung, P. (2023). Survey of hallucination in natural language generation, *ACM Computing Surveys,* 55 (12), 1–38.

Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. D. L., Bressand, F., Lengyel, G., Lample, G., Saulnier, L. and Lavaud, L.R. (2023). Mistral 7B, *arXiv preprint arXiv:2310.06825*. DOI: 10.48550/arXiv.2310.06825

Jurafsky, D. and Martin, J. H. (2023). *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition,* Stanford University, University of Colorado at Boulder.

LM Studio. (2024). LM Studio Documentation. [Online], available at: https://lmstudio.ai/docs/welcome (Accessed 06 September 2024).

Luo, H., Luo, J. and Vasilakos, A. V. (2023). *BC4LLM: Trusted artificial intelligence when blockchain meets large language models,*

*Фирсанова В. И. Новый графовый подход к генерации текстов ...*
*Firsanova V. I. A graph-based approach to closed-domain natural language generation*

164

arXiv preprint arXiv:2310.06278. DOI: 10.48550/arXiv.2310.06278

McCarthy, J. (1987). Generality in artificial intelligence, *Communications of the ACM*, 30 (12), 1030–1035.

Meister, C., Cotterell, R. (2021). Language Model Evaluation Beyond Perplexity, *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)*, 5328–5339.

Mikolov, T., Chen, K., Corrado, G. and Dean, J. (2013). *Efficient estimation of word representations in vector space*, arXiv preprint arXiv:1301.3781. DOI: 10.48550/arXiv.1301.3781

Mistral. (2024). Mistral Large 2. [Online], available at: https://mistral.ai/news/mistral-large-2407/ (Accessed 06 September 2024)

Morris, J., Hirst, G. (1991). Lexical Cohesion Computed by Thesaural relations as an indicator of the structure of text, *Computational Linguistics*, 17 (1), 21–48.

OpenAI. (2024). GPT-4o mini: advancing cost-efficient intelligence. [Online], available at: https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/ (Accessed 06 September 2024)

OpenAI API. (2024). Open AI API. [Online], available at: https://openai.com/index/openai-api (Accessed 06 September 2024)

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A. and Schulman, J. (2022). Training language models to follow instructions with human feedback, *Advances in neural information processing systems*, 6 (35), 27730–27744. DOI: 10.48550/arXiv.2203.02155

Polyzotis, N. and Zaharia, M. (2021). *What can data-centric AI learn from data and ML engineering?* arXiv preprint arXiv:2112.06439. DOI: 10.48550/arXiv.2112.06439

Priest, G. (2000). *Logic: A Very Short Introduction*, Oxford University Press, Oxford, UK.

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W. and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer, *Journal of machine learning research*, 21 (140), 1–67.

Rajpurkar, P., Zhang, J., Lopyrev, K. and Liang, P (2016). *SQuAD: 100,000+ questions for machine comprehension of text*, arXiv preprint arXiv:1606.05250. DOI: 10.48550/arXiv.1606.05250

Rajpurkar, P., Jia, R. and Liang, P. (2018). *Know what you don't know: Unanswerable questions for SQuAD*, arXiv preprint arXiv:1806.03822. DOI: 10.48550/arXiv.1806.03822

Ruder, S. (2019). *Neural transfer learning for natural language processing*, NUI Galway.

Schmidhuber, J. (1987). *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*, Technische Universität München.

Talmor, A, Herzig, J, Lourie, N. and Berant, J. (2018). *Commonsenseqa: A question answering challenge targeting commonsense knowledge*, arXiv preprint arXiv:1811.00937. DOI: 10.48550/arXiv.1811.00937

Thakur, N., Reimers, N., Rücklé, A., Srivastava, A. and Gurevych, I. (2021). *Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models*. arXiv preprint arXiv:2104.08663. DOI: 10.48550/arXiv.2104.08663bs/2104.08663

Van Rijsbergen, P. J. (1979). *Information Retrieval*, London: Butterworths.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł. and Polosukhin, I. (2017). Attention is all you need, *Advances in neural information processing systems*, 30, 261–272. DOI: 10.48550/arXiv.1706.03762

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M. and Davison, J. (2019). *HuggingFace's Transformers: State-of-the-art natural language processing*, arXiv preprint arXiv:1910.03771. DOI: 10.48550/arXiv.1910.03771

Zhang, P., Xiao, S., Liu, Z., Dou, Z. and Nie, J. Y. (2023). *Retrieve anything to augment large language models*, arXiv preprint arXiv:2310.07554. DOI: 10.48550/arXiv.2310.07554

Zhong, W., Cui, R., Guo, Y., Liang, Y., Lu, S., Wang, Y., Saied, A., Chen, W. and Duan, N. (2023). *AGIEval: A human-centric benchmark for evaluating foundation models*, arXiv preprint arXiv:2304.06364. DOI: 10.48550/arXiv.2304.06364

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

165

**Список литературы**

Andriushchenko M., Flammarion N. Does Refusal Training in LLMs Generalize to the Past Tense? arXiv preprint arXiv:2407.11969. 2024. P. 16. DOI: 10.48550/arXiv.2407.11969

Anthropic. Claude 3.5 Sonnet Model Card Addendum, 2024. URL: https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf (дата обращения: 06.09.2024).

Ayyamperumal S. G., Ge L. Current state of LLM Risks and AI Guardrails. arXiv preprint arXiv:2406.12934. 2024. P. 9. DOI: 10.48550/arXiv.2406.12934

Choi E. Prompt injection: Parameterization of fixed inputs / Choi E., Jo Y., Jang J., Seo M. arXiv preprint arXiv:2206.11349. 2022. DOI: 10.48550/arXiv.2206.11349

Christiano P. F. Deep reinforcement learning from human preferences / P. F. Christiano, J. Leike, T. B. Brown, M. Martic, S. Legg, D. Amodei // Advances in neural information processing systems. 2017. V. 30. Pp. 1–9. DOI: 10.5555/3294996.3295184

Dettmers T. QLoRA: Efficient finetuning of quantized LLMs / Dettmers T., Pagnoni A., Holtzman A., Zettlemoyer L. // Advances in Neural Information Processing Systems. 2024. V. 36. Pp. 1–28. DOI: 10.48550/arXiv.2305.14314

Devlin J. BERT: Pre-training of deep bidirectional transformers for language understanding / Devlin J., Chang M. W., Lee K., Toutanova K. // Proceedings of NAACL-HLT. 2019. Pp. 4171–4186. DOI: 10.48550/arXiv.1810.04805

Dong Y. Building Guardrails for Large Language Models / Dong Y., Mu R., Jin G., Qi Y., Hu J., Zhao X., Meng J., Ruan W. and Huang X. // arXiv preprint arXiv:2402.01822. 2024. DOI: 10.48550/arXiv.2402.01822

Firsanova V. Towards building a mobile app for people on the spectrum // Companion Proceedings of the ACM Web Conference 2023. 2023. Pp. 555–559. DOI: 10.1145/3543873.3587533

Firsanova V. The advantages of human evaluation of sociomedical question answering systems // International Journal of Open Information Technologies. 2021. V. 9. № 12. Pp. 53–59. DOI: 10.25559/INJOIT.2307-8162.09.202112.53-59

Gage P. A new algorithm for data compression // The C Users Journal. 1994. V. 12. №. 2. Pp. 23–38.

Gao J., Galley M., Li L. Neural approaches to conversational AI // The 41st international ACM SIGIR conference on research & development in information retrieval. 2018. Pp. 1371–1374. DOI: 10.1145/3209978.3210183

Goodfellow I., Bengio Y., Courville A. Deep learning. MIT press, 2016. P. 781.

Google Cloud. Cloud Computing Services, 2024. URL: https://cloud.google.com/ (дата обращения: 06.09.2024).

Guu K. Retrieval augmented language model pre-training / Guu L., Lee K, Tung Z, Pasupat P, Chang M. // InInternational conference on machine learning. Pp. 3929–3938.

Hendrycks D. Measuring massive multitask language understanding / Hendrycks D., Burns C., Basart S., Zou A., Mazeika M., Song D., Steinhardt J. arXiv preprint arXiv:2009.03300. 2020. P. 27. DOI: 10.48550/arXiv.2009.03300

Hewitt J., Manning P. D. A structural probe for finding syntax in word representations // Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). 2019. Pp. 4129–4138. DOI: 10.18653/v1/N19-1419

Hu E. J. Lora: Low-rank adaptation of large language models / Hu E. J., Shen Y., Wallis P., Allen-Zhu Z., Li Y., Wang S., Wang L., Chen W. arXiv preprint arXiv:2106.09685. 2021. P. 26. DOI: 10.48550/arXiv.2106.09685

Jacob B. Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference / Jacob B., Kligys S., Chen B., Zhu M., Tang M., Howard A., Adam H., Kalenichenko D. arXiv preprint arXiv:1712.05877. 2018. P. 14. DOI: 10.48550/arXiv.1712.05877

Ji, Z. Survey of hallucination in natural language generation / Ji Z., Lee N., Frieske R., Yu T., Su D., Xu Y., Ishii E., Bang Y., Chen D., Dai W., Chan H. S., Madotto A., Fung P. // ACM Computing Surveys. 2023. V. 55. № 12. Pp. 1–38.

Jiang A. Q. Mistral 7B / Jiang A. Q., Sablayrolles A., Mensch A., Bamford C., Chaplot D. S., Casas D. D., Bressand F., Lengyel G., Lample G., Saulnier L., Lavaud L. R. arXiv preprint arXiv:2310.06825. 2023. P. 9. DOI: 10.48550/arXiv.2310.06825

Фирсанова В. И. Новый графовый подход к генерации текстов ...
Firsanova V. I. A graph-based approach to closed-domain natural language generation

166

Jelinek F. Perplexity – a measure of the difficulty of speech recognition tasks / Jelinek F., Mercer R. L., Bahl L. R., Baker J. K. // The Journal of the Acoustical Society of America. 1977. V. 62. №. S1. Pp. S63–S63.

Jurafsky D., Martin J. H. Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition. Stanford University, University of Colorado at Boulder. 2023. P. 577.

LM Studio. LM Studio Documentation, 2024. URL: https://lmstudio.ai/docs/welcome (дата обращения: 06.09.2024).

Luo H., Luo J., Vasilakos A. V. BC4LLM: Trusted artificial intelligence when blockchain meets large language model. arXiv preprint arXiv:2310.06278. 2023. P. 42. DOI: 10.48550/arXiv.2310.06278

McCarthy J. Generality in artificial intelligence // Communications of the ACM. 1987. V. 30. № 12. Pp. 1030–1035.

Meister C., Cotterell R. Language Model Evaluation Beyond Perplexity // 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). 2021. Pp. 5328–5339.

Mikolov T. Efficient estimation of word representations in vector space / T. Mikolov, Chen K., Corrado G., Dean J. arXiv preprint arXiv:1301.3781. 2013. P. 12. DOI: 10.48550/arXiv.1301.3781

Mistral. Mistral Large 2, 2024. URL: https://mistral.ai/news/mistral-large-2407/ (дата обращения: 06.09.2024).

Morris, J., Hirst, G. Lexical Cohesion Computed by Thesaural relations as an indicator of the structure of text // Computational Linguistics. 1991. V. 17. № 1. Pp. 21–48.

Ouyang L. Training language models to follow instructions with human feedback / Ouyang L., Wu J., Jiang X., Almeida D., Wainwright C., Mishkin P., Zhang C., Agarwal S., Slama K., Ray A., Schulman J. // Advances in Neural Information Processing Systems. 2022. V. 35. Pp. 27730-27744. DOI: 10.48550/arXiv.2203.02155

OpenAI. GPT-4o mini: advancing cost-efficient intelligence, 2024. URL: https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/ (дата обращения: 06.09.2024).

OpenAI API. Open AI API, 2024. URL: https://openai.com/index/openai-api/ (дата обращения: 06.09.2024).

Polyzotis N., Zaharia M. What can data-centric AI learn from data and ML engineering? arXiv preprint arXiv:2112.06439. 2021. P. 5.

Priest, G. Logic: A Very Short Introduction. Oxford University Press. 2000. P. 160.

Raffel C. Exploring the limits of transfer learning with a unified text-to-text transformer / Raffel C., Shazeer N., Roberts A., Lee K., Narang S., Matena M., Zhou Y., Li W., Liu P. J. // Journal of machine learning research. 2020. V. 21. №. 140. Pp. 1–67.

Rajpurkar P. SQuAD: 100,000+ questions for machine comprehension of text / Rajpurkar P., Zhang J., Lopyrev K., Liang P. arXiv preprint arXiv:1606.05250. 2016. P. 10. DOI: 10.48550/arXiv.1606.05250

Rajpurkar P., Jia R., Liang P. Know what you don't know: Unanswerable questions for SQuAD. arXiv preprint arXiv:1806.03822. 2018. P. 9. DOI: 10.48550/arXiv.1806.03822

Ruder S. Neural transfer learning for natural language processing. NUI Galway. 2019. P. 330.

Schmidhuber J. Evolutionary principles in self-referential learning, or on learning how to learn. Technische Universität München. 1987. P. 64.

Talmor A. Commonsenseqa: A question answering challenge targeting commonsense knowledge / Talmor A., Herzig J., Lourie N., Berant J. arXiv preprint arXiv:1811.00937. 2018. P. 10. DOI: 10.48550/arXiv.1811.00937

Thakur N. Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models / Thakur N., Reimers N., Rücklé A., Srivastava A., Gurevych I. arXiv preprint arXiv:2104.08663. P. 24. DOI: 10.48550/arXiv.2104.08663bs/2104.08663

Van Rijsbergen P. J. Information Retrieval. London: Butterworths. 1979. P. 147.

Vaswani A. Attention is all you need / Vaswani A., Shazeer N., Parmar N., Uszkoreit J., Jones L., Gomez A. N. // Advances in neural information processing systems. 2017. T. 30. Pp. 261–272. DOI: 10.48550/arXiv.1706.03762

Wolf T. HuggingFace's Transformers: State-of-the-art natural language processing / Wolf T.,

*Научный результат. Вопросы теоретической и прикладной лингвистики. Т. 10, №3. 2024*
*Research result. Theoretical and Applied Linguistics, 10 (3). 2024*

167

Debut L., Sanh V., Chaumond J., Delangue C., Moi A., Cistac P., Rault T., Louf R., Funtowicz M., Davison J. arXiv preprint arXiv:1910.03771. 2019. P. 8. DOI: 10.48550/arXiv.1910.03771

Zhang P. Retrieve anything to augment large language models. / Zhang P., Xiao S., Liu Z., Dou Z., Nie J. Y. arXiv preprint arXiv:2310.07554. 2023. P. 16. DOI: 10.48550/arXiv.2310.07554

Zhong W. AGIEval: A human-centric benchmark for evaluating foundation models / Zhong W., Cui R., Guo Y., Liang Y., Lu S., Wang Y., Saied A., Chen W., Duan N. arXiv preprint arXiv:2304.06364. 2023. P. 22. DOI: 10.48550/arXiv.2304.06364

**Victoria I. Firsanova**, PhD Student, St. Petersburg State University, St. Petersburg, Russia.

**Фирсанова Виктория Игоревна**, аспирант, Санкт-Петербургский государственный университет, Санкт-Петербург, Россия.