Semantic Web 0 (0) 1 IOS Press

Reranking Answers of Large Language Models with Knowledge Graphs

Mikhail Salnikov^{* a,b}, Hai Le^{* b}, Olga Tsymboi^{c,d}, Ivan Lazichny^{a,d}, Dmitrii Iarosh^b,

Egor Cheremiskin^{c,e}, Andrey Savchenko^c, Dmitry Simakov^c and Alexander Panchenko^{a,b}

^a Artificial Intelligence Research Institute, Russia

^b Skolkovo Institute of Science and Technology, Russia

^d Moscow Institute of Physics and Technology, Russia

^e Moscow State University, Russia

Abstract. Answering natural language questions over knowledge graph data is challenging due to the vast number of facts, which can be difficult to process and navigate. One potential solution for this issue is to use mined subgraphs related to the query, although this process still requires extracting these subgraphs. This research presents a solution for extracting subgraphs related to entity candidates from a question-and-answer set, which can be obtained by inferring a large language model by calculating the shortest paths between entities. The proposed approaches detail various features that can be extracted from the subgraphs and reranking models to select the most probable answers from a list of candidates. Experiments were conducted on Wikidata to evaluate the effectiveness of the proposed approaches. This involved enumerating all the main feature types that can be extracted from mined subgraphs and a detailed analysis of the proposed features and reranking method combinations. In addition, a public web application that provides a useful web tool for studying the graph space between question and answer entities has been developed to work with subgraphs. This includes visualization of the extracted subgraph and automatic generation of natural language text to describe it.

Keywords: Large Language Model, Knowledge Graph, Shortest Path, Question Answering

1. Introduction

Answering factoid questions is a challenging task for any QA system, especially when the system does not have access to relevant knowledge from a Knowledge Graph (KG) or other external source. Despite the lack of a guaranteed correct answer, people still use Large Language Models (LLM) to address this challenge [26, 36], especially since the release of InstructGPT [16]. Given a natural language question and a knowledge base, this paper aims to generate an answer given the context of the KG. In addition to addressing the LLMs themselves, leveraging external structured knowledge bases such as Wikidata [32], DBPedia [1], and NELL [15] can potentially boost the suboptimal performance of these language models. These Knowledge Bases (Knowledge Graphs) are excellent examples of the implementation of Semantic Web principles. They demonstrate how ontologies, RDF (Resource Description Framework), and SPARQL can be used to create complex, interconnected, and machine-readable knowledge systems. While LLMs are great at generating human-like text, Semantic Web technologies are better at providing

^c Sber AI Lab, Russia

^{*}These authors contributed equally to this work.

structured, interpretable, and reusable knowledge. This makes them very useful for tasks that require precision and explainability.

As an alternative, unstructured data, such as plain text, can enhance the performance of LLM-powered factoid question answering through in-context learning. However, these techniques, known as Retrieval-Augmented Gener-ation (RAG) [13], may suffer from factual inaccuracies due to the use of untrusted general plain text sources. As a result, RAG platforms frequently require data source filtering and fine-tuning of the retrieval model for downstream applications to ensure accuracy or some other engineering hacks and tricks. Additionally, performance degradation may arise from either the controversially retrieved context or the model's inability to process context properly. LLM-based systems have been observed to suffer from hallucinations in their answer generation and reasoning pro-cesses, which has limited their utility for Knowledge Graph Question Answering (KGQA) tasks [30]. Furthermore, there is a challenge with long-tail facts, as dealing with information related to less frequent or obscure entities is difficult for all KGQA systems [8], and particularly for language models due to the limited knowledge available in training corpora about this type of facts. At the same time, state-of-the-art KGQA systems perform poorly on complex datasets [25].



Fig. 1. The motivation of our research: LLM's top-1 QA prediction accuracy is often incorrect. However, the correct answer is often in the top 10-40 alternative generations. Thus, a better reranking, such as the one presented in this article, is required. Dependence of Hit@N metric on the size of answer pool N on Mintaka [26] full dataset and Mistral [9], Mixtral [10], T5-Large-SSM and T5-XL-SSM [20] tuned models using Diverse Beam Search [31] as an inference strategy.

Despite all problems, LLMs can answer factoid questions correctly in some cases, especially when prompted to generate multiple answers using some variation of Beam Search [24, 31]. In other words, this shows that although the most likely option from the model's predictions can be incorrect, the answer pool can contain the correct answer. This highlights the importance of reranking answer pools and the use of accurate external data sources to improve the performance of a base language model after fine-tuning. As seen in Figure 1, where the Hit@N metric is evaluated for $1 \le N \le 40$, LLMs can generate accurate responses, at least for some beams. Therefore, if we have a method to select the correct responses, we can significantly improve the initial quality, bringing the performance closer to the ceiling. Our study will focus on the different reranking methods to enhance the base model's quality.

This paper presents a collection of methods to address the KGQA challenge by accurately reranking possible answers using subgraphs. Firstly, we utilize LLM-produced responses that have previously demonstrated satisfac-tory outcomes on a complex factoid question-answering dataset, further described in Section 3.1.1. In addition, it has been shown that incorporation of the KG information into LLMs significantly improves the results for various Natural Language Processing task [18]. Thus, this work focuses on extracting and utilizing the information about each question-answer pair by employing the information-dense Wikidata KG. To build upon the discussed novelty, we propose to look into this problem further in a reranking scope while still leveraging Wikidata as our external

source. Our hypothesis for this study is that there is important information in the space between entities in the KG

a web application that allows us to explore them and study their relevance to the question-and-answer pair.

- The contributions of our work are as follows: 1. We propose a novel approach to the KGQA that utilizes subgraphs generated from paths from entities mentioned in a question to answer candidate entities. This approach is based on the observation that, while the correct answer may not always be the most likely according to language model predictions, it often appears later within the sequence of generated predictions. Our study consistently improves the Hit@N metric by using various features and models in KG-based reranking.
- 2. Through our experiments, we have comprehensively compared the proposed method with well-known reranking techniques, ranging from classical algorithms based on graph features to more recent approaches using sentence transformers.
- 3. A publicly available web application¹ for KGQA, subgraph generation, visualization, and graph-to-text generation. This application demonstrates the relationships between the entities in the question and the answer, providing a better understanding of the graph's structure and motivation for one or another candidate to answer.

Our paper leverages extracted subgraphs and their features to rerank the language model (LM)'s generated answers. The main novelty compared to previous paper [25] is defined as follows:

- 1. The KGQA problem has been reformulated as a ranking problem, in contrast to previous work that focused on the top-1 answer, which does not provide sufficient information for a detailed understanding.
- 2. We conducted in-depth experiments on all possible features and combinations that can be extracted from subgraphs to address the lack of such detailed experiments in previous studies.
- 3. The system incorporates graph-to-natural text generation features, which enhances the interpretability and usability of generated answers.
- 4. A subgraph visualization and graph-to-text representation have been added to the demo to understand better the graph's structure and the relationships between its entities.

In this study, we have examined various aspects of working with subgraphs in KGQA and demonstrated that features derived from subgraphs can significantly enhance the ability of LLMs to answer questions. We release the publicly available source code on Github² for reprehensibility and transparent research.

2. Related Work

In this section, we provide an overview of the relevant research areas. Subsection 2.1 explores the KGs employed in this study. Next, in Subsection 2.2, we overview the key concepts of KGQA and discuss their advantages and limitations. In Subsection 2.3, we discuss factoid questions and methods to determine answer correctness.

2.1. Knowledge Graphs

KGs have emerged as powerful tools for organizing and representing structured knowledge and facilitating in-telligent applications. KGs consist of entities, relationships between these entities, and attributes that define each entity's characteristics. These components form a semantic framework that allows data from multiple sources to be integrated and analyzed. The concept of KGs traces back to earlier theories, such as semantic networks and ontolo-gies. A typical KG is a collection of triples, each consisting of a subject, a relationship, and an object. The subject entity can be linked to other entities through specific relationships or properties defined by the relationship type. There are several notable KGs:

¹https://kgga-nlp-zh.skoltech.ru/

²https://github.com/s-nlp/kbqa

1. Google KG: In 2012, Google introduced the KG, which significantly enhanced search capabilities by understanding the context and relationships between entities [28], marking a shift from keyword-based to entity-based search. This innovation improved the accuracy and relevance of search results, but unfortunately, we cannot use this proprietary internal resource. The existence of KGs suggests that other industries may adopt similar structures, indicating potential for further research and development in this area.

2. DBpedia: As a community-driven effort, DBpedia extracts structured information from Wikipedia and represents it as an RDF graph [2]. It has become a central hub in the Linked Open Data (LOD) cloud, facilitating the integration and interoperability of diverse datasets.

3. Wikidata: Launched by the Wikimedia Foundation, Wikidata is a collaborative knowledge base supporting Wikipedia and other projects within the Wikimedia community [32]. It is a centralized repository for structured data, facilitating more sophisticated querying and analysis capabilities. Wikidata is a fascinating example of a general KG, encompassing millions of entities and a wide range of attributes.

This paper focuses on the Wikidata KG, one of the most widely used general-purpose KGs. Their complexity makes the problem more challenging but also more realistic.

2.2. Knowledge Graph Question Answering

The KGQA task aims to generate accurate responses to a question based on facts extracted from KGs. KGQA systems utilize the structured data within KGs to provide accurate and context-aware responses to user inquiries. Leveraging language models has become crucial in various natural language processing applications, including KGQA. LLMs can perform zero-shot learning, producing answers for input prompts based on information stored in their pre-trained parameters, eliminating the need for additional training or labeled datasets.

Recent research proposes integrating knowledge from various sources, such as unstructured documents or tables (e.g., from Wikipedia) and factual information from KGs, into language models [11]. The rationale behind incorporating KGs into LMs stems from the fact that they represent a concise source of information.

Reranking candidates for Question Answering (QA) is an important and well-known problem that is used in various applications such as information KG completion [33, 34]. This is because the initial ranking of candidates often includes noise or less relevant options, and reranking helps refine the results to improve accuracy and relevance. The need for reranking in various applications further motivated us to apply this approach to KGQA with subgraphs, which is one of the novelties of this work.

2.3. Factoid Questions

Users commonly ask factoid questions, such as: (i) Which author wrote the most best-selling books in the decade that the internet was invented? (ii) Who is the grandmother of the eldest grandchild of the current British monarch?. Factoid questions usually have precise answers from a knowledge source such as KG or text corpus. We work with factoid questions and KGs as a source of knowledge in KGQA, a system for answering such questions. Answering factoid questions can be challenging without access to a KG, but KGQA systems that use KGs can achieve higher accuracy than traditional language models [6]. While language models can generate answers to factoid questions, they often produce incorrect responses, as illustrated in Figure 1. We hypothesize this is due to the lack of structured knowledge in KGs. Indeed, language models are trained on text, which may lack the same level of detail and accuracy as a regularly updated KG that can be edited.

3. Methods

In this study, we aim to explore the effectiveness of the subgraphs in reranking the LM-generated answer candidates. Our previous work highlighted the valuable information in the induced graph that includes all shortest paths between question and answer candidate entities. Despite the improvement in Hits@1 post-ranking, a few methods have been used to utilize the extracted subgraphs. Therefore, we aim to extend the reranking process while extracting as many useful features from the subgraphs as possible.



Fig. 2. The proposed method for reranking language model answers with KGs. The method includes subgraph extraction, features extraction, and various ranker approaches. The subgraphs consist of the shortest paths between question entities and answer candidates, as discussed in Section 3.1. Various features were extracted for the ranker approaches, including text, graph, and Graph2Text Sequence Features, as discussed in Section 3.2.

In this section, we examine each component of the process in detail. The overview of the proposed pipeline can be seen in figure 2. Firstly, we generate answer candidates using various LLMs and generate subgraphs, as discussed in Subsection 3.1. Next, in Subsection 3.2, we will look at which attributes are used to rank responses.

3.1. Subgraph Extraction

The main backbone of our approach is the procedure of subgraph extraction. We rely on the information conveyed in the relationships between question-answer pairs to improve the reranking of LLM generations. To further investigate how this relationship can improve performance, we employ a subgraph extraction algorithm that generates a KG's subgraph containing entities relevant to each question-answer pair and the shortest paths between them that contain relevant properties/relationships. We extract various features that can be used for reranking in addition to the subgraphs. This section presents the subgraph extraction algorithm, the features derived from the subgraphs, and our reranking approaches.

Referring to our previous research, we utilize Wikidata to extract subgraphs representing the relationship between each question-answer pair. For this paper, we deploy a similar subgraph extraction protocol, further discussed in the following section, with older (T5-large-ssm and T5-XL-ssm [23] as researched in the original paper) and more recent state-of-the-art LLMs (Mistral [9] and Mixtral [10]).

3.1.1. Answer Candidate Generation

As the subgraph extraction protocol requires answer candidates, we need a source of distinct answer candidates for each question. Most LLM approaches for QA, such as the one presented by [26], typically use Greed Search and evaluate the top-1 answer. However, it is important to note that the correct answer may not always be the top candidate. For example, the fine-tuned T5-XL-SSM [23] model achieved higher Mean Reciprocal Rank (MRR) scores for our task, indicating that reranking could improve the Hits@1 results. For an effective reranking pipeline, we require numerous unique answer candidates. However, even with Classical Beam Search, the output often consists of minor variations of a single sequence, which may not yield sufficient unique answer candidates for this reranking task.


Fig. 3. Extraction of subgraphs between question and answer in KG. We combine the extracted shortest paths into a single subgraph, adding links between intermediate nodes. Here, QE_n is the Question Entities, A is the current answer candidate entity, and the colored nodes represent the intermediate entities.

To solve the problem, we apply Diverse Beam Search [31], which produces a lot of candidates and generates them with higher variance. Diverse Beam Search is formulated as follows:

$$Y_{[t]}^{g} = \underset{y_{1}^{g}, \dots, y_{br}^{g} \in Y_{t}^{g}}{\operatorname{argmax}} \underbrace{\sum_{b \in [Bt]} \Theta(y_{b,[t]}^{g})}_{\text{diversity penalty}} + \underbrace{\sum_{b=1}^{g-1} \lambda_{g} \Delta(y_{b,[t]}^{g}, Y_{[t]}^{h})}_{\text{dissimilarity term}},$$
(1)

The formula involves splitting the set of beams at time t into g disjointed subsets $Y_{[t]}^g$, and then selecting the candidate with the highest diversity penalty, which is calculated as the sum of a diversity penalty function $\Theta(y_{b,[t]}^g)$ over all candidates in the subset. Additionally, a dissimilarity term is included, which is calculated as the sum of a dissimilarity function $\Delta(y_{b,[t]}^g, Y_{[t]}^h)$ over all previous subsets $Y_{[t]}^h$ up to time g-1. The dissimilarity term is weighted by a parameter λ_g . This formula is used to optimize the selection of answer candidates computationally efficiently. We apply Diverse Beam Search to the following LLMs: T5-large-ssm, T5-XL-ssm, Mistral, and Mixtral with 200 beams, 20 beam groups, and a 0.1 diversity penalty. We extend our previous research [25] by fine-tuning the proposed T5-like models and comparing them to more recent state-of-the-art models like Mistral and Mixtral, which should make our research more applicable to real-world use cases. T5-large-SSM and T5-XL-SSM were reported to be state-of-the-art both in the original Mintaka paper and our previous work, serving as a good baseline for comparison in this study.

To finetune the T5-like models, we first train them on English questions for 10000 steps, following the protocols outlined in the original Mintaka paper [26]. For the more state-of-the-art Mistral and Mixtral, we finetune with LoRA and train on English questions by generating the answer candidates with "*Answer as briefly as possible*

 $shortest_paths \leftarrow get_shortest_path_from_entity_to_candidate(entity, candidate)$ unique_node_neighbor
<--- get_neighboring_nodes(unique_node)

without additional information. [Question]". However, for the T5-like models, despite adhering to these protocols, we could not achieve the reported Hits@1 accuracy in the original paper. Despite this challenge, the main focus of the study is on the reranking aspect of the pipeline. Therefore, this paper's primary contribution is improving our fine-tuned models. 3.1.2. Question-Answer Subgraph Construction With our LM's produced answer candidates (Section 3.1.1) and the question entities, we seek to combine the shortest paths between each question entity and the current answer candidate. Thus, for each question-answer can-didate pair, the desired subgraph G is mathematically defined as an induced subgraph of the Wikidata KG. Thus, given our shortest paths from $e_i \rightarrow A$, where e_i — entity extracted from the question and A — Answer. We can use the following Algorithm 1 to extract G. Let us define H as the set of all distinct nodes within our shortest paths P_i . We want to preserve all edges between the nodes within H. We aim to retain the relationship between our question entities E and answer candidate entity A_i for all question-answer pairs. The process is schematically depicted in Figure 3. Given the computational limitations of Wikidata Ouery Services, we cannot extract the shortest paths to construct each respective subgraph using the subgraph extraction algorithm. A time-out protocol exists for each SPARQL shortest path query to Wikidata Query Services. As a solution to this limitation, we utilize *igraph*³, a library con-sisting of network analysis tools with an emphasis on efficiency and portability. With *igraph*, we parse the entire

Wikidata KG via Wikidata's online RDF dumps⁴. By building the parsed local Wikidata KG via *igraph*, we can efficiently construct the subgraphs dataset with the extraction algorithm discussed in 1.

3.2. Features based on Extracted Subgraphs

Algorithm 1 Subgraph Extraction

 $H \leftarrow get_unique_nodes_shortest_paths_flattened(shortest_paths)$

G.add_edge_between(unique_node, neighbour_node)

for neighbour_node in unique_node_neighbor do

if neighbour_node in H then

Require: entities, candidate

for entity in entities do

 $G \leftarrow DirectedGraph()$

end if

end for

end for

return G

for unique node in H do

Ensure: subgraph G

end for

After extracting subgraphs for all answer candidates of our LMs, we use all possible useful features for reranking. Referring to our previous study, we mainly focused on a simple text representation of the extracted subgraphs to rank our answer candidates. Thus, in this study, we propose extracting as many useful features as possible and analyzing each feature's importance in this reranking problem. We have divided the features into the following main categories: graph, text, and Graph2Text sequence features.

³https://jgraph.org

⁴https://dumps.wikimedia.org/wikidatawiki/entities/

3.2.1. Graph Features

With our extracted subgraphs and their corresponding answer candidate, we seek to use the relationship from the subgraphs to classify the correct answer candidate. As the first simple baseline, we utilize graph features consisting of simple numerical subgraph statistics. We hypothesize that subgraphs with the correct answer will be less "complex" than subgraphs with the incorrect answer candidate. Therefore, we would want the graph features to convey the complexity of the respective subgraph. With a clear objective in mind, we experiment with the following graph features:

- Number of nodes and edges: basic statistics of the nodes and edges of graph G.
- Number of cycles: a cycle of graph G is a non-empty path that starts from a given node and ends at the same node.
- Number of bridges: a bridge of graph G is an edge, where its deletion increases the number of connection components.
 - Average shortest path: the average of each shortest path between the question entity and the answer entity.
- **Density**: measurement of the density of a graph, where the number of edges in a dense graph is close to the maximal number of edges (each pair of nodes is connected by an edge). The density d for the graph G is formulated as $d = \frac{m}{n(n-1)}$, where n is the number of nodes and m is the number of edges in G.
- **Katz centrality** [12]: measurement of the importance (or "centrality" how "central" a node is in the graph) of a specific node *i* in a graph *G*. The Katz centrality for node *i* of graph *G* is formulated as $x_i = \alpha \sum_j A_{ij} x_j + \beta$, where *A* is the adjacency matrix of graph *G* with eigenvalues λ , β is the parameter that controls the initial centrality, and $\alpha < \frac{1}{\lambda_{\text{max}}}$.
 - PageRank [17]: a popular algorithm used by Google to rank web pages in the search query by counting the number and quality of links to a page to determine an estimate of its importance. In graph theory, the "web pages" and "links" are synonymous with nodes and edges.

We hypothesize that these features may provide ranker models with insights into the complexity of the respective subgraphs.

3.2.2. Text Features

As researched in our original paper [25], the ablation study showcased the importance of including the question within the text representation of the subgraph. Therefore, besides the simple graph features, we want to emphasize each question/answer pair without using extracted subgraphs. Thus, the text features represent the concatenation between the question and answer, separated by a semicolon — ";". To use this simple concatenation for all ranker approaches, we encode the string using the MPNet⁵ embedding model [29], discussed more in A.

3.2.3. Graph2Text Sequence Features

Given the vast amount of data contained in KGs, it is essential to convert this information into natural language to facilitate understanding and accessibility. Converting a KG into text, known as KG-to-text or Graph2Text, has demonstrated notable success in various applications [35]. Therefore, when generating text from a KG, it is crucial to analyze the underlying graph structure carefully to ensure accurate translation.

Without an obvious way of incorporating the question within the subgraphs, relying purely on the subgraphs to rerank is ineffective [25]. Therefore, we address this issue by further exploration of different KG-to-text methods. The main objective is experimenting with various techniques to represent the extracted subgraphs more explicitly. For this type of textual feature, we researched and developed three methods for representing subgraphs as a text, including **Graph2Text Deterministic, Graph2Text T5, and Graph2Text GAP**.

Firstly, we employ the **Graph2Text Deterministic** approach, the most straightforward text linearization approach. In simple terms, the subgraphs are unraveled by their matrix representation. Firstly, to linearize, we convert the subgraph into its binary adjacency matrix representation, A. Given n nodes in the subgraph, the resulting matrix's dimension will be $n \times n$. The matrix's element [i, j] represents the existence of an edge between a node with index *i* and a node with index *j*. Then, we replace the edges in the matrix with the edge label and call it A'. Lastly,

	Entities	2,730		_	Total	623,902	
	Relations Triples	354 81 927			Unique Entity	8,075 60%	
(a) Knowledge Gra	ph statistics	. Total number	r of KG compo-	(b) Texts statistics. T	he percer	ntage of text er	ntities represents
nents, number of to	okens in the i	narratives.		the portion of the text	t that incl	udes entity lab	els.
we unravel A' ro	w by row to	produce ou	r final sequenc	e and add the triple (no	de_from	, edge, node_	to) to our final
sequence. Algori	thm 2 sumr	narizes the a	forementioned	steps.			
Algorithm 2 Sut	graphs to S	Sequence					
Require: Subgra	aph G						
Ensure: Text rep	presentation	n of subgraph	1 Seq				
$adj_matrix \leftarrow$	get_adjace	ncy_matrix(C	3)				
sey ← ^m	atrix do						
for i in i d	0						
if j not	0 then						
edg	ge_info = ge	et_edge_betw	veen_nodes(G	<i>i</i> , <i>j</i>)			
Sec	$l \leftarrow Node(i$).label + edg	e_info + Node	(j).label			
end if							
end for							
end for							
Teturn Seq							
F 4	• , ,	. 1					
For the remain	ang two tex	t linearization	on approaches,	Graph2 lext 15 and \mathbf{G}	rapn21	ext GAP, We	stances where
each includes a k	G from D	Bredia [1] ai	nd a target text	$\frac{1}{2.0}$ comprising one or more	e senten	ces that desc	ribe the graph.
The test set is div	vided into p	artitions of s	seen (DBpedia	categories present in the	e trainin	g set) and un	seen (DBpedia
categories not pre	esent in the	training set).	The statistics (of this hand-crafted and	human-v	verified datase	et are described
in detail in Table	1.	U ,					
The idea behin	d the Grap	h2Text T5 a	approach is to	extract informative and	useful fe	atures from l	KGs using pre-
trained text-to-te	xt LMs. Wi	th the impre	ssive capabilit	es of pretrained LMs in	the text	t-to-text gene	ration task, we
	such result	s in the grap	h-to-text scope	e. Our idea is built upor	the ana	logous algor	ithm discussed
seek to replicate	ore tackle f	ne granh_to_	Tavt ganarotiot		two por	ouiar text-to-t	ext pre-trained
seek to replicate in [22]. The auth	T5 Those	models have	an ancodor d	task in this work with	oh moleo	a thom wall	cuited for acr
seek to replicate in [22]. The auth LMs, BART and ditional text gene	T5. These	models have	an encoder-de	task in this work with ecoder architecture, whi	ch make	es them well-	suited for con-
seek to replicate in [22]. The auth LMs, BART and ditional text gene BART and T5 us	T5. These eration task	models have s. To adapt to owing approx	an encoder-de these models f aches:	a task in this work with ecoder architecture, whi for the graph-to-text tas	ch make k, the au	es them well- uthors continu	suited for con- ue pre-training
seek to replicate in [22]. The auth LMs, BART and ditional text gene BART and T5 us	T5. These eration task	models have s. To adapt	an encoder-det these models f aches:	or the graph-to-text tas	ch make k, the au	es them well- uthors continu	suited for con- ue pre-training
seek to replicate in [22]. The auth LMs, BART and ditional text gend BART and T5 us 1. Language M	T5. These eration task ing the follo Aodel Adap	models have s. To adapt owing approa	an encoder-de these models f aches: (): the models a	a task in this work with ecoder architecture, whi for the graph-to-text tas are trained on reference	ch make k, the au texts tha	es them well- athors contine t describe gra	suited for con- ue pre-training uphs, following
seek to replicate in [22]. The auth LMs, BART and ditional text gene BART and T5 us 1. Language N the BART a	T5. These eration task ing the follo Aodel Adap and T5 pre-	models have s. To adapt to owing approximation (LMA training strat	an encoder-de these models f aches: (): the models a egies.	a task in this work with ecoder architecture, whi for the graph-to-text tas are trained on reference	texts tha	es them well- uthors continu t describe gra	suited for con- ue pre-training uphs, following
seek to replicate in [22]. The auth LMs, BART and ditional text gene BART and T5 us 1. Language N the BART a 2. Supervised collected fr	T5. These eration task ing the follo Model Adap and T5 pre- Task Adap	models have s. To adapt to owing approa- potation (LMA training strat tation (STA)	aches: (): the models a egies. (): the models a egies. (): the models a r domain as the	a task in this work with ecoder architecture, whi for the graph-to-text tas are trained on reference are trained on pairs of g	texts tha	es them well- athors contine t describe gra ad their corre	suited for con- ue pre-training uphs, following sponding texts
seek to replicate in [22]. The auth LMs, BART and ditional text gene BART and T5 us 1. Language M the BART a 2. Supervised collected fr	T5. These eration task ing the follo Aodel Adap and T5 pre- Task Adap om the sam	models have s. To adapt owing approa ptation (LMA training strat tation (STA) e or a simila	 an encoder-dethese models faches: aches: aches: b): the models a egies. b): the models a r domain as the 	a task in this work with ecoder architecture, whi for the graph-to-text tas are trained on reference are trained on pairs of g e target task — graph-to	texts that raphs ar	es them well- uthors continu t describe gra nd their corre this case.	suited for con- ue pre-training uphs, following sponding texts
seek to replicate in [22]. The auth LMs, BART and ditional text gene BART and T5 us 1. Language N the BART a 2. Supervised collected fr Building on the	T5. These eration task ing the follo Aodel Adap and T5 pre- Task Adap om the sam STA appro	models have s. To adapt to owing approxi- otation (LMA training strat tation (STA) e or a simila ach via T5 a	an encoder-de these models f aches: (a): the models a egies. (b): the models a r domain as the and WebNLG	are trained on reference re trained on pairs of g target task — graph-to-to- 2.0, we obtain graph-to	texts tha raphs ar -text se	es them well- uthors continu t describe gra nd their corre this case. quences by f	suited for con- ue pre-training uphs, following sponding texts first converting
seek to replicate in [22]. The auth LMs, BART and ditional text gene BART and T5 us 1. Language N the BART a 2. Supervised collected fr Building on the the graph into a	T5. These eration task ing the follo Aodel Adap and T5 pre- Task Adap om the sam STA appro sequence of	models have s. To adapt owing approa- potation (LMA training strat tation (STA) e or a simila ach via T5 a of tokens thr	an encoder-de these models t aches: (a): the models a egies. (b): the models a r domain as the and WebNLG rough lineariza	are trained on reference trained on pairs of g target task — graph-to- to, we obtain graph-to tion. We use the string	texts tha texts tha raphs ar -text in to- text se	es them well- athors continue t describe gra- nd their corre- this case. quences by f rt the [graph	suited for con- ue pre-training uphs, following sponding texts irst converting] to [text]:" to
seek to replicate in [22]. The auth LMs, BART and ditional text gene BART and T5 us 1. Language N the BART a 2. Supervised collected fr Building on the the graph into a acquire this linea	T5. These eration task ing the follo Model Adap and T5 pre- Task Adap om the sam STA appro sequence of rised seque	models have s. To adapt to owing approa- otation (LMA training strat tation (STA) e or a simila ach via T5 a of tokens thr ence. This ou	aches: (a): the models is aches: (a): the models is egies. (b): the models is r domain as the and WebNLG rough lineariza (tput sequence (T5 approximate)	are trained on pairs of g e target task — graph-to- 2.0, we obtain graph-to tion. We use the string is then fed into the inp	ch make k, the au texts tha raphs ar o-text in co-text se g "conve ut seque	es them well- athors continue t describe gra- ad their correct this case. quences by f rt the [graph nce for the T	suited for con- ue pre-training uphs, following sponding texts first converting to [text]:" to 5 model tuned

Lastly, the **Graph2Text GAP** approach is based on the current state-of-the-art graph-to-text task, GAP, built on BART [3]. The main idea of GAP is a fully graph-aware encoding combined with the coverage of pre-trained LMs. The GAP KG-to-text framework fuses graph-aware elements into existing pre-trained LMs, capturing the advantages brought forth by both model types. The architecture of this solution consists of two main components:

- 1. **Global Attention**: to capture the graph's global semantic information, the graph's components are first encoded using an LM. This allows the model to leverage the lexical coverage of pre-trained LMs.
- 2. **Graph-aware Attention**: to attend to and update the representations of entities, relations, or both, a topological-aware graph attention mechanism was introduced, which includes entity and relation type encoding.

Applying the work of GAP, we first linearize the input graph into a text string by creating a sequence of all triples in the KG, interleaved with tokens that separate each triple and the triple's components (head, relation, and tail). Then, we use a transformer encoder to obtain vector representations. The first module in each transformer layer acts as a Global Attention and captures the semantic relationships between all tokens. Moreover, we use a Graph-aware Attention module to capture the sparse nature of adjustment in a graph and apply it to entity and relation vectors from word vectors. By proposing this flexible framework, where graph-aware components can be interchanged, the current architecture aims to generate coherent and representative text descriptions of the KG. Like the Graph2Text T5 approach, we pretrain the model on the WebNLG 2.0 dataset and get the final predictions through the fine-tuned model. For finetuning the Graph2text GAP approach, we use the following hyperparameters: *learning rate*: $2e^{-5}$, batch size: 16; beam size: 5, Adam Optimizer, 50 nodes, and 60 relations.

In this research, we introduce the more complex neural-based graph-to-text approach to explore further the reranking capabilities of the textual representation of our extracted subgraphs. The initial rudimentary text linearization approach has already achieved state-of-the-art Hits@1. We look for a more complete case study on reranking the text linearization with these two neural-based linearization approaches. To further digest the two methods, a comparison between the Graph2Text T5 and Graph2Text GAP sequences can be seen in the table 2. Additionally, for better visualization, we implement a web application that automatically applies the T5 and GAP approaches to the desired subgraph, discussed in detail in Section 6.

All three variation of Graph2Text Sequence features are further encoded with MPNet embedding model [29], discussed more in A. Moreover, motivated by our previous research, we employ context and highlight these Graph2Text sequences, discussed further in 4.4.

3.3. Rankers

With the subgraphs and their extracted features discussed above, we devise several reranking approaches to maximize the performance of the base models. As the focal point of the research is the reranking scope, we employ reranking methodologies from least to most complex. The hypothesis is a positive trend in performance as we apply more complex models and features.

As a starting point, we employ semantic reranking. This is a popular solution in information retrieval [5, 7], implemented differently under the same name. Building on this foundation, our semantic ranker utilizes the MPNET [29] embeddings of the answer candidates, further justified in 4.3. We then rank the answer candidates by the cosine similarity between the embedding vectors.

In the next layer of complexity, we utilize regression-based models, namely, linear and logistic regression. For the features set, we apply all features discussed in 3.2 (for features in text/string format, we apply MPNET embeddings, discussed in 4.3). In the case of linear regression, we employ ordinary least squares linear regression to predict either 1 or 0, corresponding to correct and incorrect responses, respectively. The predicted score was then used to rank the potential answers by sorting the values from highest to lowest. Although we employ logistic regression for the same reranking task, we reformat the problem to a classic classification problem. We sort the answers with the highest classification confidence to rank the candidates. We use a standard logistic regression model with L2 regularisation. In addition to regression-based models, we want to utilize the same features with a more complex ranker. Thus, we explore and experiment with the gradient boosting model, specifically the CatBoost regression model [19]. Before training, we use grid search to finetune learning_rate, depth, and iteration. We use root-mean-square deviation



(RMSE) to evaluate. Like linear regression, we use the predicted score to rank the answer candidates by sorting the Last but not least, the most effective and complex approach for reranking answer candidates is a neural-based

ranker with textual features as input. We experiment with a transformer-based model with an additional regression head layer, which is fine-tuned using mean-square loss and AdamW optimisation [14]. To keep the experiments clear and transparent, we keep the same MPNet model for this ranker. We employ this variation of sentence trans-former throughout this research for various sentence/text embeddings, as mentioned in 4.3. Due to the lack of a straightforward method for utilizing numeric or table-like features with this ranker, we choose not to apply graph

values.

4. Experiments

In this section, we describe an experimental setup to test the usability of our proposed language model reranking methods based on KGs. We explore the impact of different combinations of features on reranking performance and examine the effects of various reranking methods, ranging from simple to more sophisticated, on reranking accuracy.

4.1. Dataset

To further enhance the results gathered in the original paper, we also conduct our research on Mintaka [26] dataset, which is a large-scale, complex and natural dataset that can be used for end-to-end question-answering models, composed of 20,000 question-answer pairs. This dataset is annotated with Wikidata entities and comprises 8 types of complex questions. These types include:

- Count, e.g., Q: How many astronauts have been elected to Congress? A: 4.
- Comparative, e.g., Q: Is Mont Blanc taller than Mount Rainier? A: Yes.
- Superlative, e.g., Q: Who was the youngest tribute in the Hunger Games? A: Rue.

- Ordinal, e.g., Q: Who was the last Ptolemaic ruler of Egypt? A: Cleopatra.

- Multi-hop, e.g., Q: Who was the quarterback of the team that won Super Bowl 50? A: Peyton Manning.
- Intersection, e.g., Q: Which movie was directed by Denis Villeneuve and stars Timothee Chalamet? A: Dune.
- Difference, e.g., Q: Which Mario Kart game did Yoshi not appear in? A: Mario Kart Live: Home Circuit.
- Yes/No, e.g., Q: Has Lady Gaga ever made a song with Ariana Grande? A: Yes.
- Generic, e.g., Q: Where was Michael Phelps born? A: Baltimore, Maryland.

⁶https://github.com/s-nlp/subgraph_kgqas-nlp/KGQA_Subgraphs_Ranking

This research centers around the reranking aspect of the pipeline, discussed in our previous research [25]. Thus, with the question types listed above, we exclude Yes/No and Count questions. These question types offer no value information, as numbers and "yes/no" have no respective Wikidata entities, leading to a non-existent/impractical relationship between the question/answer pair. Thus, we deem Yes/No and Count pointless in the scope of our research. However, one could still utilize the entire pipeline to compute and evaluate the results based on the whole Mintaka dataset, including Yes/No and Count. Referring to the question type classifier introduced in our original research [25], this additional component allows for Yes/No and Count questions to receive special treatment.

We also compile and publish⁶ the dataset of subgraphs for the whole Mintaka dataset (for train, validation, and test splits separately). Additionally, we also publish the answer candidates generated by the base LMs. Subgraphs are collected using the process discussed in Section 3.1: we generate candidate answers, we take the true answer and the entities from the question entity neighbors as candidates, and construct subgraphs with Algorithm 1. As a result, we construct a "correct" subgraph containing the correct highlighted answer and several "incorrect" subgraphs from the incorrect candidate answers generated by the model. We present four versions of the dataset with subgraphs: with candidates generated by T5-Large-SSM, T5-XL-SSM, Mistral, and Mixtral models.

4.2. Question Entities

Referencing the subgraph extraction protocol discussed in 3.1, we require question entities and the answer entity for each question-answer pair. Regarding the question entities, any entity linker such as mGENRE [4] could be applied. However, with the objectives outlined above, utilizing an entity linker would derail the main focus of evaluating this reranking scope. As a result, we leverage the golden truth question entities provided by the Mintaka dataset.

4.3. Text Embeddings

Some features derived from the extracted subgraphs are in their natural language form (answer candidates for semantic ranker, Graph2Text sequences, and text features). Therefore, we use these features for our reranking objectives based on the MPNet embedding model [29]. In Performance Sentence Embeddings (evaluation of the quality of the embedded sentence) and Performance Semantic Search (evaluation of the quality of the embedded search queries & paragraph), the MPNet Embedding model outperforms 37 other sentence transformer models [21]. These models were compared by averaging the Performance Sentence Embeddings and Performance Semantic Search while considering the speed and the model size. In addition to the highest embedding performance, MPNet is relatively small while fast in training time. Moreover, this model is very popular within the HuggingFace community ⁷. Utilizing a well-known and widespread embedding model would enhance the aim of formulating our approach as a reranking problem motivated by end-user requirements.

4.4. Graph2Text with Highlight & Context

As mentioned in 3.2.3, we focus on different text representations of the subgraphs to rank the respective answer candidates. We employ context for linearised text representation of the subgraph. This addition is a simple concatenation between the question and the linearised sequence, separated by a special token </s> to emphasize the question in the question-answer pairing. Moreover, the study showcased the effectiveness of highlighting (HL) the answer candidate within the linearised sequence of the concatenation. The context is motivated by the assumption that the subgraph alone does not provide the necessary information to answer the question. In other words, the model cannot answer the question without it. Similarly, it is difficult to rank the answers if the model has no idea which entities in the subgraph are potential answer candidates. An example of such concatenation, the rank by MPNet to achieve the current state-of-the-art, is shown below:

Question: Which actor was the star of Titanic and was born in Los Angeles, California? **Answer**: Leonardo DiCaprio

Original Deterministic Sequence with HL and Context[25]:

Which actor was the star of Titanic and was born in Los Angeles, California? </s> [unused1]Leonardo Di-Caprio[unused2], place of birth, Los Angeles, Titanic, cast member, [unused1]Leonardo DiCaprio[unused2] We employ a similar **HL** and **context** protocol to our three Graph2Text sequences. Similar to the original approach [25], the objective is to assist the model in understanding the question-answer pair and the Graph2Text sequence. An example of such processing for the three sequences can be seen below:

Question: Which actor was the star of Titanic and was born in Los Angeles, California? **Answer**: Leonardo DiCaprio

Graph2Text Deterministic: Which actor was the star of Titanic and was born in Los Angeles, California? </s> [unused1]Leonardo DiCaprio[unused2], place of birth, Los Angeles, Titanic, cast member, [unused1]Leonardo Di-Caprio[unused2]

Graph2Text T5: Which actor was the star of Titanic and was born in Los Angeles, California?</s> Los Angeles born [unused1]Leonardo DiCaprio [unused2], who played the role of Jack Sparrow in the film Titanic, was born in the United States.

Graph2Text GAP: Which actor was the star of Titanic and was born in Los Angeles, California?</s>Born in Los Angeles, the actor, [unused1]Leonardo DiCaprio [unused2], was a member of the crew of the Titanic.

It is important to note that we employ the **context** and **HL** approaches only for the MPNet approach, discussed in 3.3. Sensibly, the transformer-based model trained on sentences and paragraphs is adept at extracting useful information from natural texts. Thus, unlike our other proposed approaches, the MPNet approach classically only handles

⁷https://huggingface.co/sentence-transformers/all-mpnet-base-v2

text embedding as input. Therefore, we utilize context and HL to give the model the best chance at extracting useful information to rank answer candidates. Other approaches, such as regression-based and gradient boosting, have various other features (i.e., graph, text, and sequence). Thus, we choose not to employ context and HL transformation for Graph2Text sequences for regression-based and gradient-boosting rankers. We discuss the pipeline for each ranker in more detail in 4.5.

4.5. Experimental Pipeline

With our two objectives of observing the effects of 1) different combinations of feature sets and 2) different reranking approaches in varying complexity, we devise our experiments for each answer candidate source (from either T5-Large-SSM, T5-XL-SSM, Mistral, or Mixtral) as the following:

- apply each feature set extracted from the subgraphs (from the least to most complex) to each proposed ranker (from the least to most complex). The feature sets A, B, C are ranked from least to most complex. We feed A, then B, then C to each ranker. The goal is to observe the performance of varying complexity feature sets with varying complexity rankers.
- add/combine different feature sets (from the least to most complex) to each proposed ranker (from the least to most complex). The feature sets A, B, C are ranked from least to most complex. We feed A, then A + B, then A + B + C to each ranker. The goal is to observe how adding more complex feature sets affects the performance of each ranker.

With the above experimental pipeline, we seek to provide an exhaustive case study on reranking LLM's answer candidates with rankers and feature sets of varying complexity.

4.6. Evaluation

As outlined in 4.5, we experiment with numerous feature sets and rankers. Leveraging the limited amount of answer candidates, the primary objective is to determine whether selecting a final answer from this set is appropriate for our question-answering task. We evaluate using the Hits@N metric for all experiments discussed. Even if the top answer (Hits@1) is incorrect, this metric Hits@N allows us to understand the potential effectiveness of the respective ranker. For example, let us define two QA systems that can provide dozens of possible answers. These two systems generate the correct answer at the second and tenth positions. From an end-user point of view, the system that provides the correct answer at the second position or earlier in the sequence of answers will be much more beneficial. With that being said, as our research's main focus is the reranking task, Hits@N will provide us with valuable insights into each feature set and ranker.

5. Results & Discussion

This section provides comprehensive results and analysis showcasing the effectiveness of 1) different combina-tions of feature sets and 2) different ranking methods. Table 4 shows the final Hits@1 results for all rankers with different answer candidates sources and various feature sets as input. We can observe several trends within this ta-ble. Firstly, for each feature set, the Hits@1 increases as the answer candidate source model gets more complex. For instance, for all rankers with text features as input, the Hits@1 increases gradually as we move from T5-Large-SSM \rightarrow T5-XL-SSM \rightarrow Mistral \rightarrow Mixtral. This behavior is observable in any feature set. This is, of course, a sensi-ble trend as the model increases in tunable parameters and size. Furthermore, the experiment results demonstrate the importance of text features, including initial questions, which is a logical conclusion for all answer candidates' sources.

Additionally, we can observe an interesting performance difference between Graph2Text T5 and GAP sequences. The counterintuitive high-quality results from these Graph2Text sequences indicate that the T5 performs better than the GAP on the downstream task. However, as discussed in Section 3.2.3, the GAP model shows better results on the

WebNLG 2.0 dataset. This finding motivates us to explore further the quality of Graph2Text models on the Mintaka

M. Salnikov et al. / Reranking Answers of Large Language Models with Knowledge Graphs Table 3

Answers Source		Question Entities Accuracy	Answer Entities Accuracy
	Gra	ph2Text (T5)	
TE Laws COM	TRAIN	0.995	0.833
15-Large-SSM	TEST	0.954	0.835
T5-XL-SSM	TRAIN	0.955	0.827
	TEST	0.955	0.829
	Grap	bh2Text (GAP)	
T5 Laws COM	TRAIN	0.922	0.773
13-Large-SSM	TEST	0.923	0.776
T5 VI SSM	TRAIN	0.924	0.767
13-AL-55M	TEST	0.926	0.77

dataset. To do so, we calculate the accuracy of the entity label represented from the subgraph in the generated text for our Mintaka subgraph dataset. This simple exploration reveals that GAP tends to omit some entities from the provided subgraph, as shown in Table 3. Moreover, based on our subjective human assessment, GAP generates more hallucinations in the Graph2Text task, shown in Table 2.

Table 4

The impact of features on the Hit@1 performance of different reranking models. We carefully tune each answer source model in the same manner on the Mintaka train set.

Answers Source	Features	Linear Regression	Logistic Regression	CatBoost	MPNet
	Text	0.2695	0.2605	0.2458	0.2620
10°	Graph	0.2338	0.2335	0.1935	-
s5-Larn	Graph2Text (Determ Lin)	0.2550	0.2440	0.2405	0.3398
1. 2	Graph2Text (T5)	0.2505	0.2313	0.2398	0.3493
	Graph2Text (GAP)	0.2393	0.2925	0.2395	0.3395
	Text	0.2955	0.2850	0.2593	0.3418
Ň	Graph	0.2550	0.2613	0.2760	-
15. M	Graph2Text (Determ Lin)	0.2640	0.2598	0.2580	0.3923
್ರೆ	Graph2Text (T5)	0.2593	0.2538	0.2485	0.3905
	Graph2Text (GAP)	0.2563	0.2503	0.2590	0.3573
	Text	0.4760	0.4730	0.3917	0.5115
\$	Graph	0.3575	0.3558	0.3632	-
Mistra	Graph2Text (Determ Lin)	0.3960	0.3970	0.3862	0.5007
Ż,	Graph2Text (T5)	0.4013	0.3985	0.4012	0.4965
	Graph2Text (GAP)	0.3885	0.3850	0.3787	0.4917
	Text	0.4883	0.4853	0.4040	0.5237
\$	Graph	0.3698	0.3680	0.3755	-
Nixtra	Graph2Text (Determ Lin)	0.4083	0.4093	0.3985	0.5130
Ź,	Graph2Text (T5)	0.4135	0.4108	0.3940	0.5087
	Graph2Text (GAP)	0.4008	0.3973	0.3910	0.5040

5.1. Features Importance

In addition to the peculiar performance of neural G2T sequences and GAP positive correlation between Hits@1 and the complexity of the answer candidates LLM, we are interested in the effectiveness of the extracted subgraphs features. Thus, we calculate the importance of the permutation feature for the regression-based and gradientboosting rankers. The most effective approach, MPNet & G2T sequence ranker, does not utilize various feature sets. Thus, we sensibly do not calculate the importance of the permutation feature for this approach.

Firstly, we prepare graph features and MPNet embeddings of the text and G2T sequence features for the regression-based rankers as input. Without the embedding-like features, these regression-based rankers fail to achieve sufficient performance while solely on graph features, seen in table 4. Due to the nature of Logistic and Linear Regression, we split up the embedding-like features into individual values of separate columns. Thus, there are no trivial ways to assess the importance of the overall embeddings. Therefore, we hone in on the importance of the graph features for each answer candidate source.

We calculate the importance of the permutation feature using 10 repetitions to estimate the graph features' effectiveness for regression-based rankers. The results of our feature importance analysis for T5-Large-SSM/T5-XL-SSM and Mistral/Mixtral can be seen in figure 4 and 5 respectively. These results indicate the overwhelming influence of PageRank in the ranking of Logistic and Linear Regression, thus motivating the future usage of this popular graph feature. Furthermore, other features appear important, particularly the number of bridges, nodes/edges, and the average shortest path between the question entities and answer candidates. Sensibly, for T5-Large-SSM and T5-XL-SSM, several features differ in importance compared to Mistral and Mixtral.

Lastly, we calculate the importance of the feature for our gradient-boosting rankers. Catboost could support various feature types, including numerical, categorical, text, and embedding features. In other words, in addition to table-like features, Catboost natively supports embedding features without splitting them up individually. Therefore, unlike regression-based rankers, it is sensible to calculate the importance of both graph and embedding-like features. Furthermore, as mentioned, Catboost also supports categorical or string-like features. However, Catboost's native embeddings support order target encodings, which is ambiguous. Furthermore, this embedding algorithm is less well-known and utilized in natural language processing. As we sought to format the problem statement from the end-user viewpoint, we decided not to utilize this component. With that being said, for all answer candidates models, the feature importance of all graph, text, and each respective Graph2Text sequence features can be seen in figure 6 and 7.

Looking closer, we can see that the Graph2Text sequence and text features overwhelmingly dominate other im-portant features. This further fortifies the assertion of the effectiveness of text features. In other words, including the initial questions strongly affects the ranker's performance. Furthermore, this analysis elaborates on the effec-tiveness of each Graph2Text sequence, which can also be seen in table 4. For both T5-like and our state-of-the-art Mistral/Mixtral, the PageRank is deemed quite important in contributing toward the final performance of Catboost. This further bolsters the claim of future usage of PageRank, both in classical graph classification and KGQA scope. Other features such as the density, Katz centrality, and average shortest path between the question entities and an-swer candidates are also deemed important. Lastly, unlike regressor-based rankers, the number of bridges does not make a huge contribution toward the performance of Catboost.

5.2. Hits@N Evaluation

Lastly, solely relying on the Hits@1 or the highest-ranking answer will not give a complete look at the effec-tiveness of the ranking approaches. As mentioned in 4.6, we choose the Hits@N metrics to showcase an exhaustive case study of all rankers and feature sets. That said, we analyze the complete sequence of answers before and after the reranking in this subsection. All proposed methods have been tested on different features, and Hit@N has been calculated for various N values to ensure that the suggested approaches have had an impact. Figure 8 shows the most representative and significant results. The reranking produced by MPNet with text and Graph2Text features demonstrates a clear improvement in quality compared to both initial predictions and baseline models. Additionally, we can observe a clear increase in the quality of the Hits@1 answers, regardless of the size of the LLMs.



Fig. 4. Permutation importance of graph features for Linear and Logistic Regression rankers on answer candidates generated by T5-Large-SSM and by T5-XL-SSM

Furthermore, the results for all possible combinations of answer candidates models, reranking models, feature sets, and answer candidate sources are presented in Appendix A. The results demonstrate the potential of utilizing subgraphs to rerank LLM's responses. Moreover, it is essential to accurately transfer all relevant information from the subgraph to the reranking model. Therefore, the outcomes of this process significantly depend on the quality of the Graph2Text process.

6. Tool for KGQA and Subgraphs Exploration

To explore the space of KGs between arbitrary questions and corresponding answer entities, we have developed a web tool that visualizes subgraphs and automatically applies the Graph2Text (T5 and GAP) methods to these subgraphs. These methods were discussed in Section 3.2.3. In addition, this web application has features for predicting answers to questions and functions as a KGQA system. This web application, available at https://kgqa-nlp-zh.skoltech.ru/graph, can be used to systematically study subgraphs in detail and evaluate the motivation behind provided answers, as well as explore the space of KGs between questions and answers, which often contains interesting information about knowledge structure.

Figure 9 presents two subgraph visualizations that showcase the versatility of our tool in exploring diverse knowl edge domains. The Figure is divided into two parts, each demonstrating a unique set of entity relationships. In the
 first part of Figure 9, we see a subgraph connecting the countries of the United States, China, and the historical figure



Fig. 5. Permutation importance of graph features for Linear and Logistic Regression rankers on answer candidates generated by Mistral and by Mixtral

Yuri Gagarin. This visualization allows users to explore the complex relationships between countries and important historical figures in the context of space exploration. Interestingly, the shortest path between Yuri Gagarin and China went through Japan. This specific aspect of the KG can be confusing for users and researchers, but techniques like those we provided can help shed light on such matters. The second part of Figure 9 shows a subgraph connecting James Bond and the United Kingdom. This visualization section makes users more interested if a question about James Bond is asked. By presenting these contrasting examples within a single figure, we demonstrate the tool's capability to visualize and analyze relationships across a wide spectrum of knowledge, from historical events to popular culture, because of general KG - Wikidata [32].

We use the FastAPI⁸ framework to develop web applications that provide a public API⁹ in addition to web pages. This choice of framework allows for high-performance, easy-to-maintain code with built-in sup-port for asynchronous operations and automatic API documentation generation. To speed up backend pro-cesses, we implement a caching mechanism that stores the most popular entities our users interact with. This cache significantly reduces response times for frequently accessed data, enhancing the overall user experi-ence. The caching strategy is dynamically adjusted based on usage patterns and entity popularity, ensuring op-timal performance. One of the application's key features is a subgraph visualization in SVG format. Third parties can easily integrate this scalable vector graphics output into their applications or research projects.

- ⁸https://fastapi.tiangolo.com
- ⁹https://kgqa-nlp-zh.skoltech.ru/docs



Fig. 6. Permutation importance of graph, text, and G2T features for Catboost rankers on answer candidates generated by T5-Large-SSM and T5-XL-SSM

The web application also provides a KGQA component with various backends, including the T5 sequence-tosequence model, trained on the Mintaka [26] dataset. To illustrate the capabilities of our system, Figure 10 presents an example query: "Which movie is part of the Marvel Cinematic Universe and has Chadwick Boseman in the titular role?" The Figure showcases the answers generated using beam search, a technique that allows exploring multiple potential answer paths. The correct answer for this one question is "Black Panther", and the subgraph for this one answer looks simpler and clearer (Figure 11) can help the user to be sure about the answer.



7. Conclusion

In this paper, we propose a novel methodology for reranking answers in KGQA based on the use of subgraphs. We demonstrate the effectiveness of this method by combining it with various feature sets and reranking models and show that incorporating subgraph information can significantly improve accuracy in answer selection. Our work also emphasizes the importance of Semantic Web technologies, which form the basis for Knowledge Graphs. Despite the rise of modern Large Language Models (LLMs), Semantic Web standards remain crucial for creating structured, interpretable, and semantically rich knowledge representations. These technologies enable the creation of Knowledge Graphs, which are essential for tasks such as Question Answering, where precise and reliable information retrieval is essential.



Fig. 8. Hit@N Metrics for Different N Values on Mintaka [26] full dataset for generated by Diverse Beam Search [31] answers by LLM tuned models with different reranking approaches.

Our work has several implications for further research in the field of KGQA. First, it emphasizes the importance of considering structural relationships between entities to answer complex questions. Second, it demonstrates that integrating graph-based techniques with traditional language models can be beneficial. Finally, our approach suggests that subgraphs may be a valuable tool for enhancing the transparency and interpretability of KGQA systems.

The proposed methods for reranking answers to factual questions could be implemented in real-world applications, such as zero-click search, a common feature of modern search engines. A web application that visually represents subgraphs of questions and candidate answers, along with a graph-to-text representation, could serve as an additional tool for assisting humans in verifying the accuracy of provided answers.

Our study offers a novel perspective on answer selection in KGQA, emphasizing the potential advantages of incorporating subgraph data into the reranking procedure. We anticipate that this study stimulate further investigation and advancement in this field, ultimately resulting in more precise and efficient question-answering systems.







2	
3	

4
5

		Table 5			
Hit@1-3 for all propo	sed reranking appro	baches and features after rerankir	ng the T5-Lar	ge-SSM	generated answer candidates.
1 1		T (
	Reranking Model	Features	Hit@1	Hit@2	Hit@3
	Without reranking		0.2395	0.3580	0.4020
	Full random		0.0278	0.0562	0.0828
		—	0.0200	0.0647	0.0005
	Semantic reranking	Text	0.0298	0.0647	0.0965
		Text	0.2695	0.4502	0.5400
		Graph	0.2338	0.4007	0.5000
		Text + Graph	0.2953	0.4697	0.5558
		G2T (Det. Lin.)	0.2550	0.4262	0.5230
	Linear Regression	G2T (T5)	0.2505	0.4230	0.5193
		G2T (GAP)	0.2393	0.4133	0.5123
		Text + Graph + G2T (Det. Lin.)	0.3065	0.4873	0.5755
		Text + Graph + $G2T(T5)$	0.3070	0.4835	0.5710
		Text + Graph + G2T (GAP)	0.3033	0.4870	0.5728
		Text	0.2605	0.4353	0.5360
		Graph	0.2335	0.4040	0.4980
		Text + Graph	0.2470	0.4240	0.5218
		G2T (Det. Lin.)	0.2440	0.4203	0.5163
		G2T (T5)	0.2313	0.4068	0.5028
	Logistic Regression	G2T (GAP)	0.2925	0.4688	0.5625
	0 0	Text + Graph + G2T (Det. Lin.)	0.3060	0.4873	0.5745
		Text + Graph + $G2T(T5)$	0.2794	0.4780	0.5668
		Text + Graph + G2T (GAP)	0.2925	0.4688	0.5625
		Text	0.2458	0.3753	0.4525
		Graph	0.1935	0.3353	0.4238
		Text + Graph	0.1975	0.3463	0.4510
		G2T (Det. Lin.)	0.2405	0.4178	0.5068
		G2T (T5)	0.2398	0.4025	0.4868
	CatBoost	G2T (GAP)	0.2395	0.3975	0.4873
		Text + Graph + G2T (Det. Lin.)	0.1808	0.3208	0.4120
		Text + Graph + $G2T(T5)$	0.1888	0.3318	0.4233
		Text + Graph + G2T (GAP)	0.1915	0.3350	0.4268
		Text	0.2620	0.4380	0.5345
		Text + G2T (Det. Lin.)	0.3399	0.5095	0.5900
		Text + G2T (Det. Lin.) + HL	0.3105	0.4853	0.5735
		Text + G2T (T5)	0.3493	0.5203	0.5945
	MPNet	Text + G2T $(T5)$ + HL	0.3468	0.5178	0.5938
	1911 1901	Text + G2T (GAP)	0.3395	0.5080	0.5890
		Text \pm G2T (GAP) \pm HI	0 3435	0 5170	0 5973

		Table 6			
1-3 for all r	proposed reranking app	roaches and features after reranki	ng the T5-XI	-SSM	generated answer candid
, i b ioi un p			ing the re ru	5 5511	
	Reranking Model	Features	Hit@1	Hit@2	Hit@3
	Without reranking		0.3042	0.4298	0.4688
	Full random		0.0235	0.0473	0.0723
	Semantic reranking	Text	0.0245	0.0545	0.0830
		Text	0.2955	0.4763	0.4763
		Graph	0.2550	0.4353	0.5285
		Text + Graph	0.3003	0.4825	0.5693
		G2T (Det. Lin.)	0.2640	0.4480	0.5375
	Linear Regression	G2T (T5)	0.2593	0.4438	0.5333
		G2T (GAP)	0.2563	0.4320	0.5215
		Text + Graph + $G2T$ (Det. Lin.)	0.3220	0.5073	0.5915
		Text + Graph + G2T (15) Text + Graph + G2T (GAP)	0.3170	0.5035	0.5818
		Text	0.2850	0.4683	0 5603
		Graph	0.2613	0.4385	0.5283
		Text + Graph	0.3003	0.4825	0.5693
		G2T (Det. Lin.)	0.2598	0.4420	0.5423
		G2T (T5)	0.2538	0.4393	0.5315
	Logistic Pagrossion	G2T (GAP)	0.2503	0.4310	0.5218
	Logistic Regression	Text + Graph + G2T (Det. Lin.)	0.3210	0.5025	0.5885
		Text + Graph + G2T (T5)	0.3080	0.4955	0.5778
		Text + Graph + G2T (GAP)	0.3043	0.4908	0.5815
		Text	0.2593	0.3985	0.4720
		Graph	0.2760	0.4573	0.5433
		Text + Graph	0.1950	0.3370	0.4298
		G2T (Det. Lin.)	0.2580	0.4218	0.5130
		G2T (T5)	0.2485	0.4158	0.5025
	CatBoost	G2T (GAP)	0.2590	0.4390	0.5258
		Text + Graph + G2T (Det. Lin.) Text + Graph + G2T (T5)	0.1738	0.3160	0.4065
		Text + Graph + $G2T$ (GAP)	0.1670	0.3108	0.4048
			0.1000	0.5075	0.1010
		Text	0.3418	0.5185	0.5960
		Text + G2T (Det. Lin.)	0.3923	0.5648	0.6260
		1ext + G2T (Det. Lin.) + HL Text + G2T (T5)	0.3628	0.5453	0.6175
		1ext + G2T (T5) $Text + G2T (T5) + HI$	0.3905	0.5610	0.6160
	MPNet	$T_{ext} + G_{21} (T_{3}) + HL$ $T_{ext} + G_{2T} (G \Delta P)$	0.3700	0.5490	0.6040
		Text + $G2T (GAP) \pm HI$	0.3575	0.5578	0.6165
		10AT (0/11) T IIL	0.3073	0.5520	0.0105

2	
3	

4
5

6		Table 7				
7	Hit@1-3 for all proposed reranking a	approaches and features after r	eranking the Mi	stral ger	nerated ansv	ver candidate
8						
9	Reranking Model	Features	Hit@1	Hit@2	Hit@3	
10	Without reranking		0.4655	0.5313	0.5590	
11	Full random		0.0920	0.1585	0.2165	
12		m	0.0050	0.1000	0.1700	
13	Semantic reranking	Text	0.0858	0.1393	0.1790	
14		Text	0.4760	0.5818	0.6173	
15		Graph	0.3575	0.4918	0.5653	
16		Text + Graph	0.4810	0.5848	0.6193	
17		G2T (Det. Lin.)	0.3960	0.5298	0.5898	
18	Linear Regression	G2T (T5)	0.4013	0.5295	0.5910	
19		G2T (GAP)	0.3885	0.5220	0.5845	
20		Text + $G2T$ (Det. Lin.)	0.4860	0.5875	0.6198	
21		Text + $G2T(T5)$	0.4825	0.5845	0.6200	
22		Text + $G2T$ (GAP)	0.4770	0.5798	0.6180	
22		Text	0.4730	0.5795	0.6148	
23		Graph	0.3558	0.4925	0.5645	
24		Text + Graph	0.4745	0.5833	0.6200	
25		G2T (Det. Lin.)	0.3970	0.5260	0.5960	
26		G2T (T5)	0.3985	0.5305	0.5933	
27	Logistic Regression	G2T (GAP)	0.3850	0.5195	0.5878	
28		Text + G2T (Det. Lin.)	0.4790	0.5833	0.6210	
29		Text + $G2T(T5)$	0.4745	0.5838	0.6205	
30		Text + G2T (GAP)	0.4728	0.5785	0.6183	
31		Text	0.3918	0.5070	0.5708	
32		Graph	0.3633	0.4995	0.5715	
33		Text + Graph	0.3128	0.4443	0.5295	
34		G2T (Det. Lin.)	0.3863	0.5078	0.5770	
35		G2T (T5)	0.4013	0.5295	0.5910	
36	CatBoost	G2T (GAP)	0.3788	0.5098	0.5775	
37	Carboost	Text + G2T (Det. Lin.)	0.2883	0.4208	0.5163	
37		Text + $G2T(T5)$	0.2568	0.3843	0.4958	
38		Text + G2T (GAP)	0.3013	0.4393	0.5265	
39		Text	0.5115	0.6035	0.6273	
40		Text + G2T (Det. Lin.)	0.5008	0.5913	0.6225	
41		Text + G2T (Det. Lin.) + HL	0.5140	0.6013	0.6288	
42		Text + G2T (T5)	0.4965	0.5913	0.6230	
43	MDNIat	Text + G2T (T5) + HL	0.5148	0.6003	0.6253	
44	MPINet	Text + G2T (GAP)	0.4918	0.5880	0.6205	
45		Text + G2T (GAP) + HL	0.5073	0.5988	0.6250	
46						

					27
1					1
2					2
3					3
4					4
5					5
6		Table 8			6
7	Hit@1-3 for all proposed reranking	approaches and features after rer	anking the Mixtral ge	enerated answer candidates.	7
8			6		8
9	Reranking Model	Features	Hit@1 Hit@2	Hit@3	9
10	Without reranking		0.5173 0.5628	0.5735	10
11	Entland down		0.0590 0.1000	0.1455	11
12	Full random		0.0580 0.1090	0.1455	12
13	Semantic re-anking	g Text	0.0303 0.0550	0.0715	13
14		Text	0.4883 0.5883	0.6188	14
15		Graph	0.3698 0.4983	0.5668	15
16		Text + Graph	0.4933 0.5913	0.6208	16
17		G2T (Det. Lin.)	0.4083 0.5363	0.5913	17
10	Linear Regression	G2T (T5)	0.4135 0.5360	0.5925	17
18		G2T (GAP)	0.4008 0.5285	0.5860	18
19		Text + G2T (Det. Lin.)	0.4983 0.5940	0.6213	19
20		Text + G2T (T5)	0.4948 0.5910	0.6215	20
21		Text + G2T (GAP)	0.4850 0.5850	0.6198	21
22		Text	0.4853 0.5860	0.6163	22
23		Graph	0.3680 0.4990	0.5660	23
24		Text + Graph	0.4868 0.5898	0.6215	24
25		G2T (Det. Lin.)	0.4093 0.5325	0.5975	25
26		G2T (T5)	0.4108 0.5370	0.5948	26
27	Logistic Regression	G2T (GAP)	0.3973 0.5260	0.5893	27
28	Logistic Regressio.	Text + G2T (Det. Lin.)	0.4913 0.5898	0.6225	28
29		Text + G2T (T5)	0.4868 0.5903	0.6220	29
30		Text + G2T (GAP)	0.4850 0.5850	0.6198	30
31		Text	0.4040 0.5135	0.5723	31
32		Graph	0.3755 0.5060	0.5730	32
33		Text + Graph	0.3250 0.4508	0.5310	33
34		G2T (Det. Lin.)	0.3985 0.5143	0.5785	34
35		G2T (T5)	0.3940 0.5143	0.5793	35
36	CatBoost	G2T (GAP)	0.3910 0.5163	0.5790	36
37		Text + G2T (Det. Lin.)	0.3005 0.4273	0.5178	37
3.8		Text + G2T $(T5)$	0.2690 0.3908	0.4973	39
30		Text + G2T (GAP)	0.3135 0.4458	0.5280	30
10		Text	0.5238 0.6100	0.6288	39
4.1		Text + G2T (Det. Lin.)	0.5130 0.5978	0.6240	40
41		Text + G2T (Det. Lin.) + HL	0.5263 0.6078	0.6303	41
42		Text + G2T (T5)	0.5088 0.5978	0.6245	42
43	MPNet	Text + $G2T(T5)$ + HL	0.5270 0.6068	0.6268	43
44		Text + G2T (GAP)	0.5040 0.5945	0.6220	44
45		Text + G2T (GAP) + HL	0.5195 0.6053	0.6265	45
46					46
47					47
48					48

References

S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak and Z.G. Ives, DBpedia: A Nucleus for a Web of Open Data, in: *The Semantic Web, 6th International Semantic Web Conference, 2nd Asian Semantic Web Conference, ISWC 2007 + ASWC 2007, Busan, Korea, November 11-15, 2007, K. Aberer, K. Choi, N.F. Noy, D. Allemang, K. Lee, L.J.B. Nixon, J. Golbeck, P. Mika, D. Maynard, R. Mizoguchi, G. Schreiber and P. Cudré-Mauroux, eds, Lecture Notes in Computer Science, Vol. 4825, Springer, 2007, pp. 722–735. doi:10.1007/978-3-540-76298-0_52.*

- [2] C. Bizer, J. Lehmann, G. Kobilarov, S. Auer, C. Becker, R. Cyganiak and S. Hellmann, DBpedia A crystallization point for the Web of Data, J. Web Semant. 7(3) (2009), 154–165. doi:10.1016/J.WEBSEM.2009.07.002. https://doi.org/10.1016/j.websem.2009.07.002.
- [3] A.M. Colas, M. Alvandipour and D.Z. Wang, GAP: A Graph-aware Language Model Framework for Knowledge Graph-to-Text Generation, in: *Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022, N. Calzolari, C. Huang, H. Kim, J. Pustejovsky, L. Wanner, K. Choi, P. Ryu, H. Chen, L. Donatelli, H. Ji, S. Kurohashi,* P. Paggio, N. Xue, S. Kim, Y. Hahm, Z. He, T.K. Lee, E. Santus, F. Bond and S. Na, eds, International Committee on Computational Linguistics, 2022, pp. 5755–5769. https://aclanthology.org/2022.coling-1.506.
- [4] N. De Cao, L. Wu, K. Popat, M. Artetxe, N. Goyal, M. Plekhanov, L. Zettlemoyer, N. Cancedda, S. Riedel and F. Petroni, Multilingual autoregressive entity linking, *Transactions of the Association for Computational Linguistics* 10 (2022), 274–290.
- [5] J.H. Figueroa, F. Pérez-Téllez and D. Pinto, Measuring semantic similarity of documents with weighted cosine and fuzzy logic, *J. Intell. Fuzzy Syst.* **39**(2) (2020), 2263–2278. doi:10.3233/JIFS-179889.
- [6] V. Gupta, M.K. Chinnakotla and M. Shrivastava, Retrieve and Re-rank: A Simple and Effective IR Approach to Simple Question Answering over Knowledge Graphs, in: *Proceedings of the First Workshop on Fact Extraction and VERification, FEVER@EMNLP 2018, Brussels, Belgium, November 1, 2018, J. Thorne, A. Vlachos, O. Cocarascu, C. Christodoulopoulos and A. Mittal, eds, Association for Computational Linguistics, 2018, pp. 22–27. doi:10.18653/V1/W18-5504. https://doi.org/10.18653/v1/W18-5504.*
 - [7] M.L. Henderson, R. Al-Rfou, B. Strope, Y. Sung, L. Lukács, R. Guo, S. Kumar, B. Miklos and R. Kurzweil, Efficient Natural Language Response Suggestion for Smart Reply, *CoRR* abs/1705.00652 (2017). http://arxiv.org/abs/1705.00652.
- [8] W. Huang, G. Zhou, M. Lapata, P. Vougiouklis, S. Montella and J.Z. Pan, Prompting Large Language Models with Knowledge Graphs for Question Answering Involving Long-tail Facts, *CoRR* abs/2405.06524 (2024). doi:10.48550/ARXIV.2405.06524. https://doi.org/10. 48550/arXiv.2405.06524.
- [9] A.Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D.S. Chaplot, D. de Las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L.R. Lavaud, M. Lachaux, P. Stock, T.L. Scao, T. Lavril, T. Wang, T. Lacroix and W.E. Sayed, Mistral 7B, *CoRR* abs/2310.06825 (2023). doi:10.48550/ARXIV.2310.06825. https://doi.org/10.48550/arXiv.2310.06825.
- [10] A.Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D.S. Chaplot, D. de Las Casas, E.B. Hanna, F. Bressand, G. Lengyel, G. Bour, G. Lample, L.R. Lavaud, L. Saulnier, M. Lachaux, P. Stock, S. Subramanian, S. Yang, S. Antoniak, T.L. Scao, T. Gervet, T. Lavril, T. Wang, T. Lacroix and W.E. Sayed, Mixtral of Experts, *CoRR* abs/2401.04088 (2024). doi:10.48550/ARXIV.2401.04088. https://doi.org/10.48550/arXiv.2401.04088.
- [11] A.S. Jinheon Baek Alham Fikri Aji, Knowledge-Augmented Language Model Prompting for Zero-Shot Knowledge Graph Question Answering Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE) (2023), 78–106–.
- [12] L. Katz, A new status index derived from sociometric analysis, *Psychometrika* **18**(1) (1953), 39–43.
- [13] P.S.H. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, S. Riedel and D. Kiela, Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks, in: Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan and H. Lin, eds, 2020. https://proceedings.neurips.cc/paper/2020/hash/ 6b493230205f780e1bc26945df7481e5-Abstract.html.
- [14] I. Loshchilov and F. Hutter, Decoupled Weight Decay Regularization, in: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019, OpenReview.net, 2019. https://openreview.net/forum?id=Bkg6RiCqY7.
- [15] T. Mitchell, W. Cohen, E. Hruschka, P. Talukdar, B. Yang, J. Betteridge, A. Carlson, B. Dalvi, M. Gardner, B. Kisiel, J. Krishnamurthy, N. Lao, K. Mazaitis, T. Mohamed, N. Nakashole, E. Platanios, A. Ritter, M. Samadi, B. Settles, R. Wang, D. Wijaya, A. Gupta, X. Chen, A. Saparov, M. Greaves and J. Welling, Never-Ending Learning, *Commun. ACM* 61(5) (2018), 103–115–. doi:10.1145/3191513.
- [16] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C.L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P.F. Christiano, J. Leike and R. Lowe, Training language models to follow instructions with human feedback, in: *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho and A. Oh, eds, 2022. http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html.*
- 46 [17] L. Page, S. Brin, R. Motwani, T. Winograd et al., The pagerank citation ranking: Bringing order to the web (1999).
 - [18] S. Pan, L. Luo, Y. Wang, C. Chen, J. Wang and X. Wu, Unifying Large Language Models and Knowledge Graphs: A Roadmap, *IEEE Trans. Knowl. Data Eng.* 36(7) (2024), 3580–3599. doi:10.1109/TKDE.2024.3352100.
- Indus, Knowi, Data Eng. 30(1) (2024), 5380–5399. doi:10.1109/TKDE.2024.5552100.
 [19] L.O. Prokhorenkova, G. Gusev, A. Vorobev, A.V. Dorogush and A. Gulin, CatBoost: unbiased boosting with categorical features, in: Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS
 2018, December 3-8, 2018, Montréal, Canada, S. Bengio, H.M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi and R. Garnett, eds,
 2018, pp. 6639–6649. https://proceedings.neurips.cc/paper/2018/hash/14491b756b3a51daac41c24863285549-Abstract.html.

- [20] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li and P.J. Liu, Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer, J. Mach. Learn. Res. 21 (2020), 140:1–140:67. http://jmlr.org/papers/v21/20-074.html.
 - [21] N. Reimers and I. Gurevych, Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks, in: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, 2019. https://arxiv.org/abs/1908. 10084.
- [22] L.F.R. Ribeiro, M. Schmitt, H. Schütze and I. Gurevych, Investigating Pretrained Language Models for Graph-to-Text Generation, in: Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI, 2021, A. Papangelis, P. Budzianowski, B. Liu, E. Nouri, A. Rastogi and Y.-N. Chen, eds, Association for Computational Linguistic, 2021, pp. 211—227. https://aclanthology.org/2021. nlp4convai-1.20.
- [23] A. Roberts, C. Raffel and N. Shazeer, How Much Knowledge Can You Pack Into the Parameters of a Language Model?, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, B. Webber, T. Cohn, Y. He and Y. Liu, eds, Association for Computational Linguistics, 2020, pp. 5418–5426. doi:10.18653/V1/2020.EMNLP-MAIN.437. https://doi.org/10.18653/v1/2020.emnlp-main.437.*
- [24] M. Salnikov, M. Lysyuk, P. Braslavski, A. Razzhigaev, V. Malykh and A. Panchenko, Answer Candidate Type Selection: Text-To-Text Language Model for Closed Book Question Answering Meets Knowledge Graphs, in: *Proceedings of the 19th Conference on Natural Language Processing (KONVENS 2023), September 19-21, 2023, Ingolstadt, Germany*, M. Georges, A. Herygers, A. Friedrich and B. Roth, eds, Association for Computational Lingustics, 2023, pp. 155–164. https://aclanthology.org/2023.konvens-main.16.
- [25] M. Salnikov, H. Le, P. Rajput, I. Nikishina, P. Braslavski, V. Malykh and A. Panchenko, Large Language Models Meet Knowledge Graphs to Answer Factoid Questions, in: *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation, PACLIC* 2023, *The Hong Kong Polytechnic University, Hong Kong, SAR, China, 2-4 December 2023*, C. Huang, Y. Harada, J. Kim, S. Chen, Y. Hsu, E. Chersoni, P. A, W.H. Zeng, B. Peng, Y. Li and J. Li, eds, Association for Computational Linguistics, 2023, pp. 635–644. https://aclanthology.org/2023.paclic-1.63.
- [26] P. Sen, A.F. Aji and A. Saffari, Mintaka: A Complex, Natural, and Multilingual Dataset for End-to-End Question Answering, in: *Proceedings of the 29th International Conference on Computational Linguistics*, N. Calzolari, C.-R. Huang, H. Kim, J. Pustejovsky, L. Wanner, K.-S. Choi, P.-M. Ryu, H.-H. Chen, L. Donatelli, H. Ji, S. Kurohashi, P. Paggio, N. Xue, S. Kim, Y. Hahm, Z. He, T.K. Lee, E. Santus, F. Bond and S.-H. Na, eds, International Committee on Computational Linguistics, Gyeongju, Republic of Korea, 2022, pp. 1604–1619. https://aclanthology.org/2022.coling-1.138.
- [27] A. Shimorina and C. Gardent, Handling Rare Items in Data-to-Text Generation, in: *Proceedings of the 11th International Conference on Natural Language Generation, Tilburg University, The Netherlands, November 5-8, 2018, E. Krahmer, A. Gatt and M. Goudbeek, eds, Association for Computational Linguistics, 2018, pp. 360–370. doi:10.18653/V1/W18-6543. https://doi.org/10.18653/v1/w18-6543.*
- [28] A. Singhal, Introducing the Knowledge Graph: things, not strings, 2012, Official Google Blog. https://blog.google/products/search/ introducing-knowledge-graph-things-not/.
- [29] K. Song, X. Tan, T. Qin, J. Lu and T. Liu, MPNet: Masked and Permuted Pre-training for Language Understanding, in: Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan and H. Lin, eds, 2020. https://proceedings.neurips.cc/paper/2020/hash/ c3a690be93aa602ee2dc0ccab5b7b67e-Abstract.html.
- [30] A. Toroghi, W. Guo, M.M.A. Pour and S. Sanner, Right for Right Reasons: Large Language Models for Verifiable Commonsense Knowledge Graph Question Answering, *CoRR* abs/2403.01390 (2024). doi:10.48550/ARXIV.2403.01390. https://doi.org/10.48550/arXiv.2403. 01390.
- [31] A. Vijayakumar, M. Cogswell, R. Selvaraju, Q. Sun, S. Lee, D. Crandall and D. Batra, Diverse Beam Search for Improved Description of Complex Scenes, *Proceedings of the AAAI Conference on Artificial Intelligence* 32(1) (2018). doi:10.1609/aaai.v32i1.12340. https: //ojs.aaai.org/index.php/AAAI/article/view/12340.
- [32] D. Vrandečić and M. Krötzsch, Wikidata: A Free Collaborative Knowledgebase, *Commun. ACM* 57(10) (2014), 78–85–. doi:10.1145/2629489.
- [33] Y. Wang, M. Hu, Z. Huang, D. Li, D. Yang and X. Lu, KC-GenRe: A Knowledge-constrained Generative Re-ranking Method Based on Large Language Models for Knowledge Graph Completion, in: *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy, N. Calzolari, M. Kan, V. Hoste, A. Lenci, S. Sakti and N. Xue, eds, ELRA and ICCL, 2024, pp. 9668–9680. https://aclanthology.org/2024.lrec-main.845.*
- [34] Y. Wei, Q. Huang, J.T. Kwok and Y. Zhang, KICGPT: Large Language Model with Knowledge in Context for Knowledge Graph Completion, CoRR abs/2402.02389 (2024). doi:10.48550/ARXIV.2402.02389. https://doi.org/10.48550/arXiv.2402.02389.
- [35] Y. Wu, N. Hu, S. Bi, G. Qi, J. Ren, A. Xie and W. Song, Retrieve-Rewrite-Answer: A KG-to-Text Enhanced LLMs Framework for Knowledge Graph Question Answering, *CoRR* abs/2309.11206 (2023). doi:10.48550/ARXIV.2309.11206. https://doi.org/10.48550/arXiv. 2309.11206.
- [36] C. Zhang, Y. Lai, Y. Feng and D. Zhao, A review of deep learning in question answering over knowledge bases, AI Open 2 (2021), 205–215. doi:10.1016/J.AIOPEN.2021.12.001. https://doi.org/10.1016/j.aiopen.2021.12.001.