

KGQUIZ: Evaluating the Generalization of Encoded Knowledge in Large Language Models

Anonymous EMNLP submission

Abstract

Large language models’ successes in knowledge-intensive tasks suggest that real-world knowledge is encoded in the model parameters. However, with the exception of a few probing tasks in limited domains, it is not well understood what knowledge is stored in the models and how to evaluate LLMs’ knowledge systematically, across a spectrum of domains and progressively complex task formats. To this end, we propose KGQUIZ, a benchmark to comprehensively probe the knowledge generalization abilities of LLMs. KGQUIZ is a scalable framework constructed from triplet-based knowledge, containing five tasks with increasing complexity: true/false, multiple choice, blank filling, factual editing, and open-ended text generation. We evaluated 10 open-source and black-box LLMs on the KGQUIZ benchmark across three domains: commonsense, encyclopedic, and biomedical knowledge. Unsurprisingly, LLMs demonstrate advanced abilities in simpler tasks, while tasks requiring more complex reasoning or using domain-specific facts still present a significant challenge. KGQUIZ provides a testbed to analyze such nuanced variations in performance across domains and task formats, and ultimately to understand, evaluate, and improve LLMs’ knowledge abilities across different knowledge domains and tasks.

1 Introduction

Large language models (LLMs) demonstrate great potential to encode and represent real-world knowledge in model parameters, advancing multiple tasks and settings, such as question answering (Lin et al., 2019; Feng et al., 2020; Yasunaga et al., 2021; Zhang et al., 2021; Yasunaga et al., 2022; Feng et al., 2023), dialogue generation (Dinan et al., 2019; Liu et al., 2021; Adolphs et al., 2022), summarization (Goyal et al., 2023; Zhang et al., 2023; Liu et al., 2023) and more. However, this knowledge is brittle, with LLMs producing hallucinations

(Pagnoni et al., 2021; Ji et al., 2022; Mallen et al., 2023; Bang et al., 2023), struggling to encode long-tail facts (Mallen et al., 2023), and falling short of abstaining when relevant information is not present in model parameters (Chen et al., 2022).

As a result, studies and benchmarks have been proposed to probe the knowledge abilities of LLMs (Petroni et al., 2019; Dhingra et al., 2022; Hendrycks et al., 2021a; Sung et al., 2021; Meng et al., 2022), as well as works assessing LLMs’ temporal knowledge capabilities, making them more aware of the changing nature of facts over time (Dhingra et al., 2022). A widely adopted KILT benchmark (Petroni et al., 2021) focuses on knowledge-intensive language tasks using a single snapshot of Wikipedia as the knowledge source. In this work, we identify two important but under-explored factors in probing the generalization of LLM parametric knowledge.

Knowledge Utilization: How do LLMs employ their parametric knowledge to perform tasks with varying formats or complexities? Previous works have primarily focused on limited task formats such as fill-in-the-blank questions to test the model’s knowledge abilities. However, the complexity or format of a task can itself influence a model’s knowledge abilities. This crucial aspect often goes unaddressed in the current literature, with benchmarks generally probing model performance in specific tasks, without considering the impact of the task’s complexity or format. Our proposed benchmark presents increasingly complex tasks, allowing for a comprehensive evaluation of LLM knowledge abilities which could provide constructive insights for model performance improvements and practical applications in specific scenarios.

Knowledge Breadth: Current methods predominantly consider Wikipedia or a specific domain like biomedical knowledge as the knowledge source. However, it has been observed that LLMs performance can vary significantly across differ-

ent knowledge domains - an aspect that has not been adequately addressed in the previous works. Hence, our benchmark evaluates model performance across three different domains: commonsense, encyclopedic, and biomedical. This approach helps identify the strengths and weaknesses of LLMs across different domains. Furthermore, our benchmark is designed to be easily extended to any specific domain knowledge with triplet data from that domain. This makes our benchmark distinct from the previous ones by allowing for a more diverse and adjustable evaluation of LLMs’ broad knowledge abilities.

To this end, we propose KGQUIZ, a comprehensive benchmark designed to evaluate the knowledge abilities of LLMs across knowledge utilization patterns in diverse knowledge domains. KGQUIZ is constructed with structured information from knowledge graphs (KGs) from three varying domains, namely ConceptNet (commonsense), YAGO (encyclopedic), and UMLS (biomedical). For each knowledge graph, KGQUIZ presents a collection of 41,000 knowledge-intensive questions, covering five tasks of increasing complexity, namely true/false, multiple choice, blank-filling, multi-hop factual editing, and open-ended text generation. Significantly, the same tasks are constructed for all domains, providing us with a comprehensive and comparative setting to assess LLMs’ abilities. These tasks respectively test LLMs’ abilities to judge factual correctness, select facts based on model confidence, retrieve entity names, perform factual editing, and generate long-form knowledge documents, presenting a holistic probe of LLM knowledge abilities in different usage contexts.

We evaluate 10 open-source and black-box LLMs on the KGQUIZ benchmark to better analyze which LLM cover what knowledge domain better, and under which utilization contexts. Our experiments demonstrate that: 1) **LLM performance greatly varies across knowledge domains.** For instance, on *Task 5: Open-Ended Text Generation*, ChatGPT (Ouyang et al., 2022), ChatGLM (Du et al., 2022), and TEXT-DAVINCI-003 (Ouyang et al., 2022) respectively perform best when it comes to YAGO, ConceptNet, and UMLS, three knowledge graphs representing varying knowledge domains. 2) **Knowledge utilization greatly impacts LLM’s ability to retrieve and employ factual knowledge.** For instance, ChatGPT’s perfor-

mance on biomedical knowledge drops by 30% from the fill-in-the-blank task to the factual editing task, suggesting that the additional multi-hop context in factual editing poses new challenges to LLM knowledge abilities. Together, our extensive experiments demonstrate that probing the knowledge abilities of LLMs is nuanced and multi-faceted, with the largest LLMs excelling in simple knowledge utilization tasks on general knowledge domains, while advanced knowledge contexts and domain-specific information remain open challenges. We envision KGQUIZ as a valuable testbed to understand, evaluate, and improve LLM knowledge abilities across varying knowledge domains and utilization contexts.

2 The KGQUIZ Benchmark

KGQUIZ employs knowledge graphs from diverse domains to construct five knowledge-intensive tasks with increasing complexity. We denote a knowledge graph as a set of triples \mathcal{T} , where the k -th triple is $\mathcal{T}_k = (h_k, r_k, t_k)$, and h_k , r_k and t_k represent the head entity, relation, and tail entity, respectively. We use \mathcal{E} and \mathcal{R} to denote the sets of all entities and relations in the knowledge graph.

2.1 Task 1: True-or-False

Given a structured knowledge fact, a true-or-false question asks whether the statement is factually correct or not, providing a simple and straightforward way to assess knowledge abilities (Clark et al., 2019). We adopt True-or-False questions as the first task of the KGQUIZ benchmark, aiming to evaluate to what extent could LLMs accurately verify the factuality of KG-based information with their inherent parametric knowledge.

Task Formulation We construct two sets of KG triples to represent positive and negative samples (\mathcal{T}_{pos} and \mathcal{T}_{neg}). For positive triple $(h, r, t) \in \mathcal{T}_{pos}$, we replace the tail entity t with another entity t' to generate a negative sample and add it to \mathcal{T}_{neg} . We then use the prompt for the positive or negative triple (h, r, t) : “Is the statement h r t True or False?”. We expect LLMs to answer with *True* or *False*, indicating their judgment of the knowledge statement based on their parametric knowledge.

Negative Sampling We propose four approaches to sample negative entities t' in the knowledge graph to obtain hard negative samples.

- **Random** We randomly sample an entity from a set of entities not connected to the head entity h

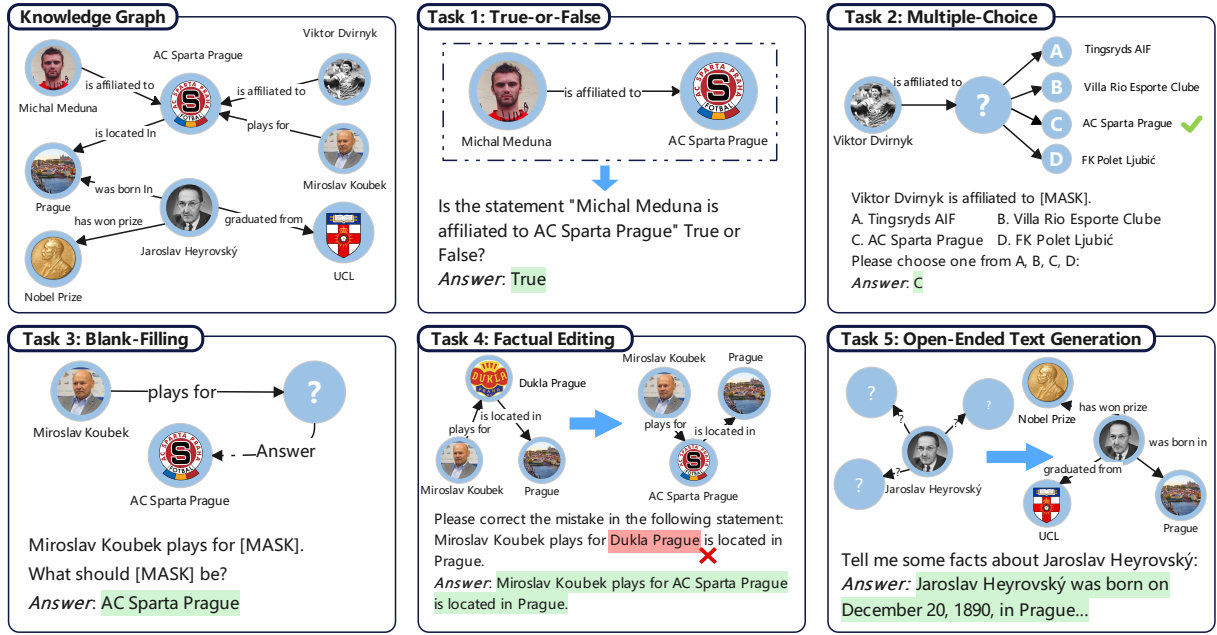


Figure 1: Overview of the KGQUIZ Benchmark, featuring five knowledge-intensive tasks with increasing complexity.

as t' , formally $t' \in \mathcal{E} - \mathcal{E}(h)$, where $\mathcal{E}(h)$ denotes the set of entities connected to h .

- **Semantic Similarity** We hypothesize that semantically similar entities could provide a more challenging setting with harder negative examples. We first use the **Random** method to sample m negative entities. These sampled entities form the set \mathcal{E}_m . Then, we employ an encoder-based language model, denoted as $\text{enc}(\cdot)$, to encode the names of these entities. Finally, we use cosine similarity $\text{sim}(\cdot, \cdot)$ to select an entity t' that is most similar to t in the embedding space. Formally, $t' = \arg\max_{e \in \mathcal{E}_m} \text{sim}(\text{enc}(e), \text{enc}(t))$.
- **Relation Sharing** We hypothesize that entities that share the same relation, r , as the selected triple would provide a challenging adversarial setting. We first obtain the set of entities that also has relation r as $\mathcal{E}^{(r)}$, then randomly sample one entity from $\mathcal{E}^{(r)}$ as the negative sample t' .
- **Network Proximity** We hypothesize that entities that are close to h in the KG could also present a hard negative example. We obtain the set of entities that connected to h and randomly sample one entity from it as the negative sample t' .

Evaluation We use accuracy as the evaluation metric for the binary output of *True* or *False*.

2.2 Task 2: Multiple-Choice

Multiple-choice question is a classic test format widely adopted in many question-answering bench-

marks (Talmor et al., 2019; Hendrycks et al., 2021b; Robinson et al., 2022). Compared to true-or-false problems, multiple-choice questions additionally present distractors to test whether LLMs could evaluate the plausibility of different answer options based on their parametric knowledge.

Task Formulation We randomly sample a subset of the knowledge graph, denoted as \mathcal{T}_s . For $(h, r, t) \in \mathcal{T}_s$, we replace the tail entity t with $[MASK]$ and provide m answer options, including the correct entity t and $m - 1$ distractors. We follow the same negative sampling strategies in *Task 1: True-or-False* to obtain the distractors.

Evaluation We calculate accuracy as the evaluation metric.

2.3 Task 3: Blank-Filling

Blank-filling provides LLMs with part of the factual association and asks LLMs to fill in the missing information (Petroni et al., 2019). Compared to *Task 1: True-or-False* and *Task 2: Multiple-Choice*, the blank-filling task is more challenging as it requires LLMs to retrieve the correct answer without any options and it becomes very hard for LLMs to guess the correct answer.

Task Formulation We randomly sample one subset of the knowledge graph, denoted as \mathcal{T}_s . For $(h, r, t) \in \mathcal{T}_s$, we replace the tail entity t with $[MASK]$. The model is asked to generate the correct answer to replace $[MASK]$.

Evaluation We denote the model output as t_o and we use the following metrics for evaluation:

- **LCS:** We denote the Longest Common Subsequence of t_o and t as s , and LCS is defined as:
$$\text{LCS} = \frac{\text{Len}(s)}{\max\{\text{Len}(t_o), \text{Len}(t)\}}$$
- **F1-score:** We denote the set of common tokens in both t_o and t as C . We denote the F1-score of t_o and t as $F1 = \frac{2PR}{P+R}$, where $P = \frac{|C|}{|t_o|}$, $R = \frac{|C|}{|t|}$.
- **Semantic Match:** We measure semantic similarity between the model’s output and the correct answer using cosine similarity on embeddings obtained via InstructGPT Ada LLM $\text{enc}(\cdot)$. This gives us the $\text{AdaScore}(t_o, t) = \text{sim}(\text{enc}(t_o), \text{enc}(t))$. A threshold θ of Adascore is based on a held-out validation set (detailed in Appendix B) to determine whether the model-generated answer and the ground truth are a semantically exact match. Concretely, we define the semantic match metric as $\text{SM}(t_o, t) = 1$ if $\text{AdaScore}(t_o, t) \geq \theta$, else 0.

2.4 Task 4: Factual Editing

In addition to examining the knowledge memorization of LLMs, we also explore whether LLMs could process sentence contexts to identify factual inconsistencies and revise accordingly. To this end, we design the Factual Editing task, where LLMs are asked to detect inconsistency and correct factual errors in multi-hop knowledge statements. While previous works have also explored LLMs’ potential in factual editing (Balachandran et al., 2022; Chen et al., 2023), we uniquely focus on a multi-hop format where one of the hop features inconsistent factual information.

Task Formulation Given a knowledge graph, we first sample a k -hop path, and we use a structured format to present the multi-hop knowledge path as $d = (h_1, r_1, e_1, r_2, \dots, t_k)$.¹ Then we randomly replace one of the entities in the path (denoted as e_s) with e' sampled with the negative sampling strategies in Task 1: *True-or-False* to obtain d' . This task prompts LLMs to correct the factual inconsistency in d' .

Evaluation We adopt the same set of evaluation metrics as Task 3: *Blank-Filling*, namely LCS, F1-SCORE, and SEMANTIC MATCH, to compare the ground truth entity e_s and the revised entity given by LLMs.

¹To avoid confusion, we denote e_m as the tail entity t_m of the m -th triple in the knowledge path. At the same time, it also serves as the head entity h_{m+1} of the $(m+1)$ -th triple in the knowledge path.

2.5 Task 5: Open-Ended Text Generation

The Open-Ended Text Generation Task focuses on assessing the ability of LLMs to generate multiple factual associations about a given entity and evaluate whether it matches the information in existing knowledge graphs. This comparison aims to measure the ability of LLMs to generate accurate and comprehensive factual knowledge of a particular entity. In addition, while tasks in previous works mostly focus on a single factual association (Talmor et al., 2019; Hendrycks et al., 2021b), we propose the Open-Ended Text Generation task to encourage the knowledge abilities of LLMs in multi-fact and knowledge synthesis settings.

Task Formulation We randomly sample one subset of KG, denoted as \mathcal{T}_s . For $(h, r, t) \in \mathcal{T}_s$, we ask the model to “Tell me some facts about h ”. We denote all triplets containing h in the knowledge graph as $\mathcal{G} = \{(h, r_g, t_g) \in \mathcal{T}\}$.

Evaluation We evaluate Open-Ended Text Generation generation by comparing the model outputs with the information about entity h denoted as \mathcal{G} . Concretely, we first prompt a GPT-3.5 LLM to turn the given model output in natural language into a list of fact triplets $\mathcal{O} = \{(h, r_o, t_o)\}$ inspired by previous works (Josifoski et al., 2023; Min et al., 2023), where we further evaluate this approach in Appendix B. Using the semantic match metric SM in Task 3: *Blank-Filling*, we define the Precision and Recall between model predictions \mathcal{O} and ground truth \mathcal{G} as: Precision = $\frac{|\mathcal{O} \cap \mathcal{G}|}{|\mathcal{O}|}$, Recall = $\frac{|\mathcal{O} \cap \mathcal{G}|}{|\mathcal{G}|}$, where $\mathcal{O} \cap \mathcal{G}$ denotes the set of triples that are both in model predictions and the knowledge graph with SM = 1.

3 Experiment Settings

Knowledge Domains We consider KGs from three distinct domains in our experiments: ConceptNet (Speer et al., 2017) for commonsense, YAGO (Mahdisoltani et al., 2015) for encyclopedic, and UMLS (Bodenreider, 2004) for the biomedical domain.²

Models and Settings We evaluate both black-box and open-source LLMs on the KGQUIZ benchmark. For black-box LLMs, we adopt InstructGPT (Ouyang et al., 2022) (TEXT-ADA-001, TEXT-BABAGGE-001, TEXT-CURIE-001, and TEXT-DAVINCI-003) and ChatGPT (GPT-3.5-TURBO) through the OpenAI API. For open-source LLMs,

²Details about size of each KG is in Appendix B.

Model	YAGO			ConceptNet			UMLS		
	F1-score	LCS	Sem. Match	F1-score	LCS	Sem. Match	F1-score	LCS	Sem. Match
ADA	2.26	18.24	61.67	1.24	11.76	45.43	5.72	19.43	55.52
BABBAGE	2.60	17.63	60.48	2.07	12.06	64.67	10.37	21.68	71.43
CURIE	<u>5.38</u>	19.63	71.54	3.32	15.11	78.68	<u>10.90</u>	<u>26.04</u>	84.70
DAVINCI	14.02	28.65	73.00	6.27	27.40	91.19	8.28	23.81	<u>87.88</u>
TURBO	4.47	11.83	52.33	<u>5.56</u>	14.42	80.48	19.44	28.18	89.27
GPT-J	0.56	10.75	24.55	1.20	4.53	39.07	9.38	11.74	73.17
OPT	0.66	10.75	27.33	0.75	4.40	45.55	6.88	11.21	73.52
CHATGLM	3.53	<u>21.50</u>	<u>72.27</u>	2.35	<u>20.15</u>	<u>88.07</u>	4.04	19.45	58.71
LLAMA	1.24	11.43	35.97	1.03	3.42	25.96	7.44	9.31	76.64
ALPACA	3.16	10.37	41.52	1.92	6.25	56.55	10.63	13.61	81.88

Table 1: LLM performance on *Task 3: Blank-Filling*. DAVINCI leads on YAGO and ConceptNet, while TURBO performs best on UMLS, indicating that LLM knowledge abilities vary greatly depending on the knowledge domain.

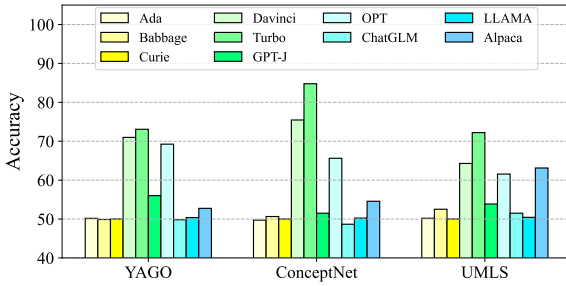


Figure 2: Model performance on *Task 1: True-or-False*. Larger LMs are generally better at judging factual correctness, while the same LM performs differently across varying knowledge domains.

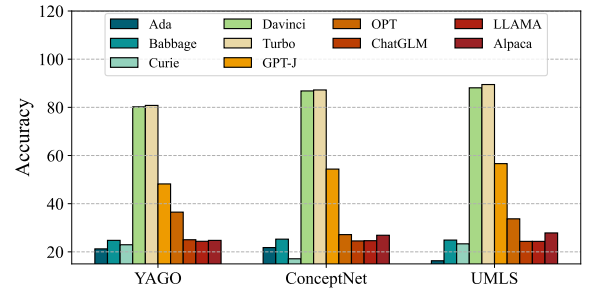


Figure 3: LLM performance on *Task 2: Multiple-Choice*. DAVINCI and TURBO consistently outperform other models, indicating their knowledge abilities under the multiple-choice knowledge utilization format.

we adopt GPT-J (Wang and Komatsuzaki, 2021), OPT (6.7B) (Zhang et al., 2022), ChatGLM (Du et al., 2022), LLAMA (7B) (Touvron et al., 2023), and Alpaca (Taori et al., 2023) in the experiments. We use a temperature of $\tau = 0$ to ensure deterministic model behavior, while employing a one-shot setting for Task 1, 2, and 4 and zero-shot setting for Task 3 and 5 due to the nature of the task. We use *Semantic Similarity* to sample negative examples in our subsequent experiments.³ Additional details on experiment setup in Appendix B.

4 Results

4.1 Task 1: True-or-False

As depicted in Figure 2, among the assessed LLMs, four of them (TEXT-DAVINCI-003, GPT-3.5-TURBO, ChatGLM) performed substantially better than random chance (50%) on all KGs. Notably, GPT-3.5-TURBO achieved the best overall performance, showcasing its ability to discern correct from incorrect knowledge statements. Ob-

³The specific effect of these four strategies and our choice for *Semantic Similarity* is detailed in section 5.1.

servation of improved performance with larger model sizes suggests that models with more parameters can encode more knowledge and leverage the stored knowledge to accurately identify the veracity of knowledge statements. Additionally, Even in the simple binary task, many LLMs show accuracy close to 50%, indicating difficulty in distinguishing true from false statements. This suggests a need for further improvement in LLMs’ knowledge abilities, particularly for those with smaller parameter sizes.

4.2 Task 2: Multiple-Choice

Figure 3 showcases that TEXT-DAVINCI-003 and GPT-3.5-TURBO consistently outperformed other LLMs in understanding and applying knowledge across all KGs and domains. An observation from tasks comparison revealed TEXT-DAVINCI-003 and GPT-3.5-TURBO’s improved performance in *Task 2: Multiple-Choice* compared to *Task 1: True-or-False*. However, Alpaca’s relative performance dwindled in Task 2, suggesting that the specific knowledge utilization format significantly influences an LLM’s ability to retrieve potentially cor-

Model	YAGO			ConceptNet			UMLS		
	F1-score	LCS	Sem. Match	F1-score	LCS	Sem. Match	F1-score	LCS	Sem. Match
ADA	2.50	<u>14.51</u>	86.76	0.12	14.65	83.84	2.50	<u>18.11</u>	59.85
BABBAGE	2.90	9.47	90.68	0.02	10.42	86.53	2.90	17.78	60.03
CURIE	6.21	8.93	<u>91.20</u>	0.10	<u>15.92</u>	83.14	6.21	19.76	60.24
DAVINCI	16.99	20.58	91.77	5.15	17.31	<u>93.25</u>	<u>5.44</u>	7.28	<u>64.19</u>
TURBO	<u>12.29</u>	13.24	91.06	0.51	1.28	93.32	0.88	8.93	59.05
GPT-J	0.03	0.17	90.34	0.00	0.22	93.21	0.20	0.71	59.98
OPT	0.01	0.06	90.37	0.00	0.06	93.24	0.30	0.88	59.96
CHATGLM	4.94	1.32	89.66	0.14	4.57	90.62	0.42	2.58	76.26
LLAMA	0.03	0.04	90.33	0.00	0.00	93.20	0.43	1.81	59.98
ALPACA	6.80	12.27	90.20	<u>0.87</u>	14.84	93.20	1.46	8.66	59.93

Table 2: LLM performance on *Task 4: Factual Editing*. Model performance is generally higher than blank-filling, indicating the helpfulness of additional context. Improved performance compared to blank-filling task is observed, further emphasizing the influence of the form of knowledge utilization on LLMs’ abilities to utilize knowledge.

rect answers.

4.3 Task 3: Blank-Filling

Compared to true-or-false and multiple-choice questions, blank filling requires LLMs to retrieve the correct answer from their parametric knowledge without relying on any options. In Table 1, the overall low LCS scores reflect that LLMs’ generated answers rarely match the exact target answer. Moreover, the models’ abilities differ significantly, with TEXT-DAVINCI-003 excelling in two domains (YAGO and ConceptNet) but GPT-3.5-TURBO performing better in the biomedical domain (UMLS). Additionally, we observe a noticeable decrease in performance in the biomedical domain, suggesting that the models may not be as proficient in handling domain-specific knowledge.

4.4 Task 4: Factual Editing

Compared to blank-filling, *Task 4: Factual Editing* involves identifying and rectifying factual inconsistencies within given knowledge statements. According to the results in Table 2, the additional context indeed aids some models in generating accurate responses on certain KGs (YAGO and ConceptNet), with TEXT-DAVINCI-003 and GPT-3.5-TURBO scoring well for YAGO and ConceptNet respectively, and ChatGLM excelling in UMLS utilization. It highlights that tasks like dialogue generation and summarization, which usually come with relevant context, may work better with LLMs. However, when provided only with a short question, QA models may get confused easily. Also, The task-wise change in top-performing models reemphasizes that the form of knowledge utilization impacts an LLM’s knowledge abilities significantly. Tasks that provide more contextual infor-

Model	YAGO		ConceptNet		UMLS	
	Precision	Recall	Precision	Recall	Precision	Recall
ADA	75.84	34.89	90.93	24.90	59.45	19.47
BABBAGE	84.66	35.34	<u>95.01</u>	18.84	<u>81.52</u>	22.93
CURIE	85.69	38.64	96.59	22.46	83.43	26.80
DAVINCI	76.39	53.96	88.12	<u>41.55</u>	77.48	46.06
TURBO	<u>77.28</u>	57.63	89.39	40.53	75.94	<u>43.89</u>
GPT-J	11.97	8.78	24.11	12.07	10.72	5.96
OPT	14.06	7.72	16.89	5.26	10.35	5.43
CHATGLM	71.00	<u>54.54</u>	88.05	46.49	63.59	39.72
LLAMA	39.17	29.29	36.78	11.78	26.14	11.85
ALPACA	22.96	17.77	28.63	13.94	12.69	7.53

Table 3: Model performance on *Task 5: Open-Ended Text Generation*. Different from previous tasks, generating long and open-ended statements about given entities poses new challenges to LLMs.

mation, like dialogue generation and summarization, may work better with LLMs; however, shorter questions may cause confusion and impact LLM performance.

4.5 Task 5: Open-Ended Text Generation

Open-ended generation tasks present a more complex challenge to LLMs as it requires not just specific factual associations, but the generation of a consistent paragraph about a certain entity encapsulating assorted facts and knowledge. As observed in Table 3, TEXT-DAVINCI-003 tops the chart with the highest AdaScore_s score across all three KGs, denoting its proficient ability to produce well-structured and factually accurate knowledge paragraphs. TEXT-CURIE-001 stands out with the highest Precision score, indicating its preference to generate knowledge closely in line with the respective knowledge graph. From a Recall perspective, the best performances are achieved by GPT-3.5-TURBO, ChatGLM, and TEXT-DAVINCI-003 on the three respective KGs. These findings emphasize that the knowledge domain significantly affects

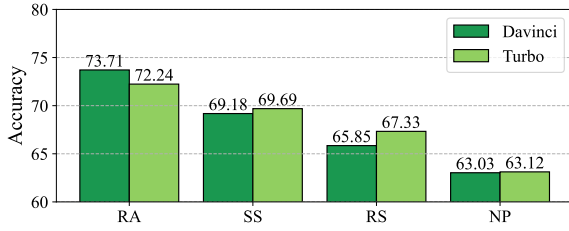


Figure 4: Performance on *Task 1: True-or-False* with varying negative sampling. We used *Semantic Similarity* (SS) as the main strategy due to its intermediate difficulty, while KGQUIZ also supports *Random* (RA), *Relation Sharing* (RS), and *Network Proximity* (NP).

the performance of LLMs in knowledge-intensive tasks, underscoring the need for comprehensive evaluations of LLMs’ knowledge abilities that consider varying knowledge domains.

5 Analysis

5.1 Negative Sampling Strategy

In section 2.1, we propose and formalize four negative sampling methods to generate questions in the KGQUIZ benchmark. In order to investigate their impact on the difficulty of the task, we use the four negative sampling strategies, *Random* (RA), *Semantic Similarity* (SS), *Relation Sharing* (RS), and *Network Proximity* (NP) to generate questions for *Task 1: True-or-False* based on the YAGO knowledge graph. We evaluate TEXT-DAVINCI-003 and GPT-3.5-TURBO as shown in Figure 4. These results show that different negative sampling methods *do* impact on the difficulty of the problem, ranging from easy to difficult in the following order: *Random*, *Semantic Similarity*, *Relation Sharing*, and *Network Proximity*. It is also demonstrated that whether LLMs can select the correct answer is impacted by the plausibility of negative examples.

In particular, we employed *Semantic Similarity* as an intermediate strategy presenting reasonable complexity. This strategy, while challenging, does not make the task excessively difficult. Furthermore, while we propose this specific strategy, KGQUIZ benchmark supports the flexibility of adopting other negative sampling settings.

5.2 Consistency Study

In this study, we investigate the robustness towards minor changes in prompts and knowledge statements. We select 100 questions from the YAGO knowledge graph in *Task 1: True-or-False* and evaluate with five different prompts and instructions

Model	ADA	BABBAGE	CURIE	DAVINCI	TURBO
Fleiss Kappa	-0.187	-0.057	-0.168	0.285	0.645

Table 4: Fleiss Kappa’ over five paraphrased prompts for different LLMs on *Task 1: True-or-False* with the YAGO knowledge graph. LLMs exhibit varying robustness to framing-level changes in knowledge statements.

Question	Prediction	Gold
Bob Hawke graduated from ____	Oxford University	University of Oxford
Rosemary Sutcliff has won ____ prize	The Carnegie Medal	Carnegie Medal (literary award)
Taito Corporation is located in ____	Tokyo, Japan	Shibuya, Tokyo

Table 5: Qualitative analysis of *Task 3: Blank-Filling*, suggesting that our proposed *Semantic Match* presents a more nuanced metric for knowledge probing.

(more details in Appendix C.8). We measure response consistency of the five black-box LLMs using the Fleiss Kappa measure (Fleiss, 1971). Table 4 shows that LLMs have varying robustness towards prompt formats: TURBO has the highest score at 0.645, suggesting a moderate level of agreement. DAVINCI exhibits a lower but still positive value of 0.285. However, ADA, BABBAGE, and CURIE show negative Fleiss Kappa values, indicating poor agreement and suggesting that model responses are less consistent towards minor changes in knowledge probing instructions. This study highlights that the robustness to minor changes in knowledge-intensive prompts is in itself part of LLM’s knowledge abilities.

5.3 Exact Match vs. Semantic Match

We conduct qualitative analysis for *Task 3: Blank-Filling* and present a few examples in Table 5. It is demonstrated that answers generated by LLMs do not exactly match the gold label, where the exact match (EM) metric would treat the answer as incorrect. However, the generated responses are semantically equivalent. For instance, in the first example, the word order is different but both answers convey the same meaning. Similarly, in the third example, “Tokyo, Japan” is more general than the gold answer “Shibuya, Tokyo” but it still provides the correct location information. While the exact match metric would treat them as incorrect, under

our proposed *Semantic Match*, all four answers are deemed as correct, indicating that *Semantic Match* presents a better evaluation metric in LLM knowledge probing given the nuanced nature of entity names (Li et al., 2020).

6 Related Work

LLM Knowledge Probing Research into what knowledge is stored in LLMs has drawn significant interest. Pioneering work like LAMA (Petroni et al., 2019), TempLAMA (Dhingra et al., 2022), MMLU (Hendrycks et al., 2021a) quantitatively measured the factual knowledge in these models. Other approaches have expanded these probing techniques, exploring topics like few-shot learning and 2-hop relational knowledge (He et al., 2021). Furthermore, open-domain question-answering benchmarks like Natural Questions (Kwiatkowski et al., 2019), and TriviaQA (Joshi et al., 2017) have been used to measure the practical knowledge abilities of these models, aligning the probing tasks with real-world applications.

Improving LLM Knowledge Abilities Efforts to enhance LLM’s knowledge abilities include augmenting language models with KGs for structured, factual knowledge (Mihaylov and Frank, 2018) and using retrieval-augmented methods like RAG (Lewis et al., 2020), REALM (Guu et al., 2020), and REPLUG (Shi et al., 2023) to incorporate external documents as a dynamic knowledge source. Further, REMEDI (Hernandez et al., 2023) aims to create a finer control over knowledge in LLMs by understanding fact encodings in the model’s internal representation system. In parallel, the framework Cook (Feng et al., 2023) suggests using specialized language models to provide modular and up-to-date knowledge in a collaborative process.

Extracting Knowledge from LLMs The extraction of knowledge from LLMs has become an emerging topic in the research community. Some works focus on constructing KGs from the LLMs (Cohen et al., 2023; Trajanoska et al., 2023). For example, Crawling Robots (Cohen et al., 2023) uses a robot role-play setting to extract named entities and relations by encoding them into actions. Other works utilize the prompt-based paradigm, where they generate knowledge probes in the form of structured prompts (Liu et al., 2022; Yu et al., 2023). These tools aim to extract and organize the knowledge within an LLM in a human-readable

and interpretable way. Furthermore, other techniques involve augmenting training data with recitation tasks to express internally represented knowledge explicitly (Sun et al., 2023).

Investigating the Limitation of LLM Knowledge Abilities As LLMs have shown promise in knowledge-based tasks, researchers have also started examining the limitations of these models’ knowledge abilities. This includes their ability to handle conflicted information (Xie et al., 2023), recall abilities (Mallen et al., 2023), and self-evaluating skills (Kadavath et al., 2022). By investigating these limitations, researchers aim to not only devise ways to address them but also shed light on how LLMs can operate more effectively in more sophisticated tasks, particularly in professional domains (Sung et al., 2021; Meng et al., 2022).

In summary, while considerable work has been done in probing the knowledge abilities of LLMs, improving these abilities, extracting knowledge, and investigating their limitations, two major aspects have seen less consideration: knowledge utilization and knowledge breadth. These areas are vital for understanding and evaluating the performance of LLMs in more real-world, complex scenarios. Therefore, this calls for a more comprehensive approach, which our proposed KGQUIZ benchmark aims to address, making strides towards a future where LLMs exhibit robust knowledge abilities applicable to a wider range of domains and utilization contexts.

7 Conclusion

We propose KGQUIZ, a benchmark for probing the knowledge generalization abilities of Large Language Models (LLMs). Unlike previous work, our benchmark focuses on two often-overlooked aspects: the complexity of knowledge utilization and the breadth of knowledge domains. Our benchmark uses structured information from knowledge graphs (KGs) across three diverse domains, and it consists of several tasks representing increasingly complex forms of knowledge utilization. Our experimental results illustrate varying performances of several LLMs across different domains and tasks, underscoring the multi-faceted nature of knowledge abilities in LLMs. We envision KGQUIZ as a comprehensive testbed to evaluate, understand, and improve the knowledge abilities of LLMs across varying domains and tasks.

Limitations

LM and KG selection Due to computational and budget constraints, we restricted our study to ten representative LLMs and three knowledge graphs each from a different domain. As we plan to make KGQUIZ publicly accessible, further investigation into the performance of a broader range of LLMs on assorted knowledge graphs is left for future endeavors.

Evaluation Metrics Being the case that LLMs might not fully adhere to the context in our prompts, we were required to deploy human-crafted string-processing functions to preprocess the content the models generated, to evaluate the results. This step is susceptible to errors that may lead to inaccurate results. Additionally, the Semantic Match method we utilized is also not without error. Two semantically similar entities could have wildly different referents, which could lead to assessment errors. Addressing the issue of fuzzy match (semantic match) is a direction for future research.

Knowledge Coverage Due to the vast scale of real-world knowledge, we are unable to evaluate whether all the content generated by the model is completely factual in our benchmark. We can only assess whether the content generated by the model aligns with the knowledge stored in the knowledge graphs. However, the coverage of real-world knowledge by the knowledge graph is limited, leading to potential errors in our evaluation. However, as our benchmark is scalable, we can mitigate this limitation to some extent by generating corresponding tasks (questions) using broader (or more applicable) and more up-to-date knowledge graphs.

Prompt Effectiveness The prompts we utilized for each question may not necessarily be the most effective. Given the constraints of our budget, we were unable to execute extensive testing on all plausible prompts. Therefore, for *Task 1: True-or-False*, *Task 2: Multiple-Choice* *Task 4: Factual Editing*, we chose the method of incorporating one in-context example to aid model understanding of the task instructions.

Ethics Statement

Privacy As KGs encompass a wealth of knowledge on a multifarious range of topics, it can include sensitive or private information. The potential for an LLM, that effectively covers and utilizes

this knowledge domain, could generate responses disclosing personal details of individuals or organizations. This introduces privacy concerns and reinforces the need for developing privacy-conscious approaches when leveraging and assessing LLMs and KGs.

Accessibility In making KGQUIZ publicly accessible, we aspire to propel further research on LLMs' knowledge abilities. However, the use of this benchmark may necessitate significant resources due to the inherent complexities of large language models. Similarly, evaluating black-box LLMs could incur significant costs, potentially creating barriers to access to the benchmark for researchers with limited computational resources or budget, contributing to elevated entry barriers in this field.

References

- Leonard Adolphs, Kurt Shuster, Jack Urbanek, Arthur Szlam, and Jason Weston. 2022. [Reason first, then respond: Modular generation for knowledge-infused dialogue](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7112–7132, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Vidhisha Balachandran, Hannaneh Hajishirzi, William Cohen, and Yulia Tsvetkov. 2022. [Correcting diverse factual errors in abstractive summarization via post-editing and language model infilling](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9818–9830, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenhao Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *ArXiv*, abs/2302.04023.
- O. Bodenreider. 2004. [The unified medical language system \(UMLS\): integrating biomedical terminology](#). *Nucleic Acids Research*, 32(90001):267D–270.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Durán, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'13, page 2787–2795, Red Hook, NY, USA. Curran Associates Inc.
- Anthony Chen, Panupong Pasupat, Sameer Singh, Hongrae Lee, and Kelvin Guu. 2023. *Purr*: Efficiently

699	editing language model hallucinations by denois-	Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2023.	755
700	ing language model corruptions. <i>arXiv preprint</i>	News summarization and evaluation in the era of	756
701	<i>arXiv:2305.14908</i> .	gpt-3 .	757
702	Hung-Ting Chen, Michael Zhang, and Eunsol Choi.	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu-	758
703	2022. Rich knowledge sources bring complex knowl-	pat, and Ming-Wei Chang. 2020. Realm: Retrieval-	759
704	edge conflicts: Recalibrating models to reflect con-	augmented language model pre-training. In <i>Proceeed-</i>	760
705	flicting evidence . In <i>Proceedings of the 2022 Con-</i>	<i>ings of the 37th International Conference on Machine</i>	761
706	<i>ference on Empirical Methods in Natural Language</i>	<i>Learning</i> , ICML’20. JMLR.org.	762
707	<i>Processing</i> , pages 2292–2307, Abu Dhabi, United		
708	Arab Emirates. Association for Computational Lin-	Tianxing He, Kyunghyun Cho, and James Glass. 2021.	763
709	guistics.	An empirical study on few-shot knowledge probing	764
		for pretrained language models .	765
710	Christopher Clark, Kenton Lee, Ming-Wei Chang,	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	766
711	Tom Kwiatkowski, Michael Collins, and Kristina	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	767
712	Toutanova. 2019. Boolq: Exploring the surprising	2021a. Measuring massive multitask language under-	768
713	difficulty of natural yes/no questions. In <i>Proceedings</i>	standing . In <i>International Conference on Learning</i>	769
714	<i>of the 2019 Conference of the North American Chap-</i>	<i>Representations</i> .	770
715	<i>ter of the Association for Computational Linguistics:</i>		
716	<i>Human Language Technologies, Volume 1 (Long and</i>	Dan Hendrycks, Collin Burns, Steven Basart, Andy	771
717	<i>Short Papers)</i> , pages 2924–2936.	Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-	772
718	Roi Cohen, Mor Geva, Jonathan Berant, and Amir	hardt. 2021b. Measuring massive multitask language	773
719	Globerson. 2023. Crawling the internal knowledge-	understanding. <i>Proceedings of the International Con-</i>	774
720	base of language models . In <i>Findings of the Asso-</i>	<i>ference on Learning Representations (ICLR)</i> .	775
721	<i>ciation for Computational Linguistics: EACL 2023</i> ,		
722	pages 1856–1869, Dubrovnik, Croatia. Association	Evan Hernandez, Belinda Z. Li, and Jacob Andreas.	776
723	for Computational Linguistics.	2023. Inspecting and editing knowledge representa-	777
		tions in language models .	778
724	Bhuwan Dhingra, Jeremy R. Cole, Julian Martin	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan	779
725	Eisenschlos, Daniel Gillick, Jacob Eisenstein, and	Su, Yan Xu, Etsuko Ishii, Yejin Bang, Wenliang Dai,	780
726	William W. Cohen. 2022. Time-aware language mod-	Andrea Madotto, and Pascale Fung. 2022. Survey of	781
727	els as temporal knowledge bases . <i>Transactions of the</i>	hallucination in natural language generation. <i>ACM</i>	782
728	<i>Association for Computational Linguistics</i> , 10:257–	<i>Computing Surveys</i> , 55:1 – 38.	783
729	273.		
730	Emily Dinan, Stephen Roller, Kurt Shuster, Angela	Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke	784
731	Fan, Michael Auli, and Jason Weston. 2019. Wizard	Zettlemoyer. 2017. TriviaQA: A large scale distantly	785
732	of wikipedia: Knowledge-powered conversational	supervised challenge dataset for reading comprehen-	786
733	agents . In <i>International Conference on Learning</i>	sion . In <i>Proceedings of the 55th Annual Meeting of</i>	787
734	<i>Representations</i> .	<i>the Association for Computational Linguistics (Vol-</i>	788
735	Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding,	<i>ume 1: Long Papers)</i> , pages 1601–1611, Vancouver,	789
736	Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Gln:	Canada. Association for Computational Linguistics.	790
737	General language model pretraining with autoregres-		
738	sive blank infilling. In <i>Proceedings of the 60th An-</i>	Martin Josifoski, Marija Sakota, Maxime Peyrard,	791
739	<i>nuual Meeting of the Association for Computational</i>	and Robert West. 2023. Exploiting asymmetry	792
740	<i>Linguistics (Volume 1: Long Papers)</i> , pages 320–335.	for synthetic training data generation: Synthie and	793
741	Shangbin Feng, Weijia Shi, Yuyang Bai, Vidhisha Bal-	the case of information extraction. <i>arXiv preprint</i>	794
742	achandran, Tianxing He, and Yulia Tsvetkov. 2023.	<i>arXiv:2303.04132</i> .	795
743	Cook: Empowering general-purpose language mod-		
744	els with modular and collaborative knowledge .	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	796
745	Yanlin Feng, Xinyue Chen, Bill Yuchen Lin, Peifeng	Henighan, Dawn Drain, Ethan Perez, Nicholas	797
746	Wang, Jun Yan, and Xiang Ren. 2020. Scalable multi-	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	798
747	hop relational reasoning for knowledge-aware ques-	Tran-Johnson, Scott Johnston, Sheer El-Showk,	799
748	tion answering . In <i>Proceedings of the 2020 Con-</i>	Andy Jones, Nelson Elhage, Tristan Hume, Anna	800
749	<i>ference on Empirical Methods in Natural Language</i>	Chen, Yuntao Bai, Sam Bowman, Stanislav Fort,	801
750	<i>Processing (EMNLP)</i> , pages 1295–1309, Online. As-	Deep Ganguli, Danny Hernandez, Josh Jacobson,	802
751	sociation for Computational Linguistics.	Jackson Kernion, Shauna Kravec, Liane Lovitt, Ka-	803
752	Joseph L. Fleiss. 1971. Measuring nominal scale agree-	mal Ndousse, Catherine Olsson, Sam Ringer, Dario	804
753	ment among many raters. <i>Psychological Bulletin</i> ,	Amodei, Tom Brown, Jack Clark, Nicholas Joseph,	805
754	76:378–382.	Ben Mann, Sam McCandlish, Chris Olah, and Jared	806
		Kaplan. 2022. Language models (mostly) know what	807
		they know .	808

809	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	Zaiqiao Meng, Fangyu Liu, Ehsan Shareghi, Yixuan Su,	864
810	field, Michael Collins, Ankur Parikh, Chris Alberti,	Charlotte Collins, and Nigel Collier. 2022. Rewire-	865
811	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	then-probe: A contrastive recipe for probing biomed-	866
812	ton Lee, Kristina Toutanova, Llion Jones, Matthew	ical knowledge of pre-trained language models . In	867
813	Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob	<i>Proceedings of the 60th Annual Meeting of the As-</i>	868
814	Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natu-	<i>sociation for Computational Linguistics (Volume 1:</i>	869
815	ral questions: A benchmark for question answering	<i>Long Papers)</i> , pages 4798–4810, Dublin, Ireland. As-	870
816	research . <i>Transactions of the Association for Compu-</i>	sociation for Computational Linguistics.	871
817	<i>tational Linguistics</i> , 7:452–466.		
818	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	Todor Mihaylov and Anette Frank. 2018. Knowledge-	872
819	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	able reader: Enhancing cloze-style reading compre-	873
820	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	hension with external commonsense knowledge . In	874
821	täschel, Sebastian Riedel, and Douwe Kiela. 2020.	<i>Proceedings of the 56th Annual Meeting of the As-</i>	875
822	Retrieval-augmented generation for knowledge-	<i>sociation for Computational Linguistics (Volume 1:</i>	876
823	intensive nlp tasks . In <i>Advances in Neural Informa-</i>	<i>Long Papers)</i> , pages 821–832, Melbourne, Australia.	877
824	<i>tion Processing Systems</i> , volume 33, pages 9459–	Association for Computational Linguistics.	878
825	9474. Curran Associates, Inc.		
826	Jing Li, Aixin Sun, Jianglei Han, and Chenliang Li.	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike	879
827	2020. A survey on deep learning for named entity	Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer,	880
828	recognition. <i>IEEE Transactions on Knowledge and</i>	Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023.	881
829	<i>Data Engineering</i> , 34(1):50–70.	Factscore: Fine-grained atomic evaluation of factual	882
		precision in long form text generation. <i>arXiv preprint</i>	883
		<i>arXiv:2305.14251</i> .	884
830	Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	885
831	Ren. 2019. KagNet: Knowledge-aware graph net-	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	886
832	works for commonsense reasoning . In <i>Proceedings</i>	Sandhini Agarwal, Katarina Slama, Alex Gray, John	887
833	<i>of the 2019 Conference on Empirical Methods in Natu-</i>	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	888
834	<i>ral Language Processing and the 9th International</i>	Maddie Simens, Amanda Askell, Peter Welinder,	889
835	<i>Joint Conference on Natural Language Processing</i>	Paul Christiano, Jan Leike, and Ryan Lowe. 2022.	890
836	<i>(EMNLP-IJCNLP)</i> , pages 2829–2839, Hong Kong,	Training language models to follow instructions with	891
837	China. Association for Computational Linguistics.	human feedback . In <i>Advances in Neural Information</i>	892
		<i>Processing Systems</i> .	893
838	Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Pe-	Artidoro Pagnoni, Vidhisha Balachandran, and Yulia	894
839	ter West, Ronan Le Bras, Yejin Choi, and Hannaneh	Tsvetkov. 2021. Understanding factuality in abstrac-	895
840	Hajishirzi. 2022. Generated knowledge prompting	tive summarization with FRANK: A benchmark for	896
841	for commonsense reasoning . In <i>Proceedings of the</i>	factuality metrics . In <i>Proceedings of the 2021 Con-</i>	897
842	<i>60th Annual Meeting of the Association for Computa-</i>	<i>ference of the North American Chapter of the Asso-</i>	898
843	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	<i>ciation for Computational Linguistics: Human Lan-</i>	899
844	3154–3169, Dublin, Ireland. Association for Compu-	<i>guage Technologies</i> , pages 4812–4829, Online. As-	900
845	tational Linguistics.	sociation for Computational Linguistics.	901
846	Shilei Liu, Xiaofeng Zhao, Bochao Li, Feiliang	Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick	902
847	Ren, Longhui Zhang, and Shujuan Yin. 2021. A	Lewis, Majid Yazdani, Nicola De Cao, James Thorne,	903
848	three-stage learning framework for low-resource	Yacine Jernite, Vladimir Karpukhin, Jean Maillard,	904
849	knowledge-grounded dialogue generation. In <i>Confer-</i>	Vassilis Plachouras, Tim Rocktäschel, and Sebastian	905
850	<i>ence on Empirical Methods in Natural Language</i>	Riedel. 2021. KILT: a benchmark for knowledge	906
851	<i>Processing</i> .	intensive language tasks . In <i>Proceedings of the 2021</i>	907
852	Yixin Liu, Alexander R. Fabbri, Pengfei Liu, Dragomir	<i>Conference of the North American Chapter of the</i>	908
853	Radev, and Arman Cohan. 2023. On learning to	<i>Association for Computational Linguistics: Human</i>	909
854	summarize with large language models as references .	<i>Language Technologies</i> , pages 2523–2544, Online.	910
		Association for Computational Linguistics.	911
855	Farzaneh Mahdisoltani, Joanna Asia Biega, and	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	912
856	Fabian M. Suchanek. 2015. Yago3: A knowledge	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	913
857	base from multilingual wikipe-dias. In <i>Conference on</i>	Alexander Miller. 2019. Language models as knowl-	914
858	<i>Innovative Data Systems Research</i> .	edge bases? In <i>Proceedings of the 2019 Confer-</i>	915
859	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,	<i>ence on Empirical Methods in Natural Language Pro-</i>	916
860	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.	<i>cessing and the 9th International Joint Conference</i>	917
861	When not to trust language models: Investigating	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	918
862	effectiveness of parametric and non-parametric mem-	pages 2463–2473, Hong Kong, China. Association	919
863	ories .	for Computational Linguistics.	920

921	Joshua Robinson, Christopher Michael Rytting, and David Wingate. 2022. Leveraging large language models for multiple choice question answering. <i>arXiv preprint arXiv:2210.12353</i> .	976
922		977
923		978
924		979
		980
925	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen tau Yih. 2023. Replug: Retrieval-augmented black-box language models .	981
926		982
927		983
928		984
929	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In <i>Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI'17</i> , page 4444–4451. AAAI Press.	985
930		986
931		987
932		988
933		
934	Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. 2023. Recitation-augmented language models . In <i>The Eleventh International Conference on Learning Representations</i> .	989
935		990
936		991
937		992
		993
938	Mujeen Sung, Jinhyuk Lee, Sean Yi, Minji Jeon, Sungdong Kim, and Jaewoo Kang. 2021. Can language models be biomedical knowledge bases? In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4723–4734, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	994
939		995
940		996
941		997
942		998
943		999
944		1000
		1001
945	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge . In <i>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)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	1002
946		1003
947		1004
948		1005
949		
950		
951		
952		
953		
954	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca .	1006
955		1007
956		1008
957		1009
958		1010
959	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models .	
960		
961		
962		
963		
964		
965	Milena Trajanoska, Riste Stojanov, and Dimitar Trajanov. 2023. Enhancing knowledge graph construction using large language models .	
966		
967		
968	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/kingoflolz/mesh-transformer-jax .	
969		
970		
971		
972	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2023. Adaptive chameleon or stubborn sloth: Unraveling the behavior of large language models in knowledge clashes .	
973		
974		
975		

A Discussion

Performance of LLMs Across Different Knowledge Domains Our comprehensive exploration of ten large-scale language models utilizing KGQUIZ revealed that these models exhibited far from uniform performance across diverse knowledge domains and contexts. For instance, the most advanced model, TEXT-DAVINCI-003 displayed varying performance across different knowledge graphs and tasks. Broadly speaking, the performance of this model was the highest on the YAGO knowledge graph, consistently surpassing other models in tasks like true-or-false and multiple-choice. However, when faced with the UMLS knowledge graph representing the biomedical domain, the model showed a significant decline in performance, with ChatGLM and GPT-3.5-TURBO taking the lead instead. These findings emphasize the model’s struggles with domain-specific knowledge. Similar trends were also observed with other models like Alpaca, which performed poorly on the multiple-choice task, but displayed a notable improvement on the blank-filling task. Such performance variations across knowledge domains serve as an interesting direction for future research, aiming to investigate the reasons behind such contrasts in LLM performance across diverse knowledge realms.

LLM Performance Across Knowledge Utilization Contexts KGQUIZ has laid emphasis on knowledge utilization patterns along with knowledge domains, providing a comprehensive overview of the knowledge abilities of LLMs. This has enabled a detailed analysis of the models’ performance across different knowledge-intensive tasks. A fascinating observation is the influence of task complexity and format on model performance. Alpaca exhibited a significant improvement from *Task 1: True-or-False* to *Task 2: Multiple-Choice*, while the performance of models like TEXT-CURIE-001 dipped. This pattern suggests various models adapt differently to varying complexity and the nature of knowledge utilization at hand. Such insights could be valuable to refine LLM’s understanding and handling of tasks, thus warranting further exploration.

Provide Comprehensive Insight for LLM Evaluation and Comparison KGQUIZ is specifically designed to offer a rich set of metrics and contexts for in-depth evaluation and comparison of

LLMs’ performance across various knowledge domains and utilization contexts. By presenting a fine-grained and multi-perspective analysis, KGQUIZ contributes to a thorough understanding of the strengths and weaknesses of individual LLMs. This not only enables researchers and users to make informed decisions when selecting the best-suited model for a specific task, but also paves the way for the evidence-based development of more capable and versatile LLMs in the future.

Guidance for Future Development of LLMs

The performance heterogeneity of LLMs that we observed across varied tasks indicates the challenges certain tasks pose for these models. For instance, LLMs, despite their robust performance on simpler tasks such as True-or-False, struggle to meet the challenge of the increasing complexity of tasks like Factual Editing, emphasizing their limitations in context-rich, advanced knowledge reasoning. Moving forward, these observations can provide valuable insights for future advancements in the field. Identifying specific areas that require attention and improvement can guide developers to iteratively refine model architectures, enrich training data, and adopt more effective pre-training and fine-tuning methods.

B KGQUIZ Details

Knowledge Graph Details In our experiments, we postulate that the performance of LLMs in knowledge-intensive tasks is greatly influenced by diverse knowledge domains. Thus, we consider knowledge graphs from three distinct domains in our experiments: commonsense, encyclopedic, and domain-specific. For commonsense KGs, we leverage the ConceptNet knowledge graph with 1,103,036 entities, 47 relations, and 3,098,674 triples. For encyclopedic KGs, we adopt the YAGO knowledge graph with 123,182 entities, 37 relations, and 1,089,040 triples. For domain-specific KGs, we mainly consider the biomedical domain and adopt the UMLS knowledge graph with 297,554 entities, 98 relations, and 1,212,586 triples. By conducting our evaluations across knowledge graphs that span a range of domains, we aim to provide a comprehensive assessment of how the knowledge abilities of LLMs fare across diverse knowledge domains.

In-Context Examples Through experiments, we discovered that for the majority of LLMs, their per-

formance in a zero-shot setting is unusually low on some tasks. We think this is because they are unable to precisely comprehend the question’s meaning (instructions), and they cannot produce output in the format we expect. Therefore, to preserve fairness without compromise, we have incorporated an in-context example into the prompts of each question for *Task 1: True-or-False*, *Task 2: Multiple-Choice*, and *Task 4: Factual Editing*, which will enable a better assessment of the model’s knowledge abilities.

Threshold for Semantic Match For three knowledge graphs, we randomly selected 1,000 entities each. For each entity, we prompted GPT-4 to generate five entities with the same reference and five entities with different references. As a result, we obtained a total of $3 \times 1,000 \times 5$ positive samples and $3 \times 1,000 \times 5$ negative samples. For each sample pair, we calculated their AdaScore. We chose a threshold so that if a positive sample’s AdaScore is above the threshold or a negative sample’s AdaScore is below the threshold, the sample pair is correctly classified; otherwise, it is misclassified. We selected the threshold that minimized the number of misclassified samples as the Semantic Match threshold.

LLM-based Triplets Extraction We find that it is difficult to measure the similarity between a piece of text and a set of triples. However, evaluating the similarity between two sets of triplets is much easier. So in KGQUIZ Benchmark, we prompt a GPT-3.5 LLM to turn the given model output in natural language into a set of fact triplets. To verify the reliability of this method, we manually evaluate 20 (essay, triplets) pairs. (essay: the TEXT-DAVINCI-003’s output text; triplets: the extracted triplets from the model output with our method.) In our human evaluation, the triplets extracted by this method have a precision of 0.87 and a recall of 0.86, demonstrating that our approach has high reliability. The problem with this method is that it extracts triples that do not have the target entity as the head, and the extracted triples do not conform to the format. We expect that providing more in-context examples can help alleviate these issues.

Task Settings For *Task 1: True-or-False*, we construct 10k examples for each knowledge graph and adopt semantic similarity as the default negative sampling method. For *Task 2: Multiple-Choice*, we use four answer options as the default setting and

construct 10k examples for each knowledge graph. For *Task 3: Blank-Filling*, we randomly sample 10k triplets for each knowledge graph to generate the blank-filling questions. Moving on to *Task 4: Factual Editing*, we construct 10k knowledge walks for each knowledge graph with the default walk length $k = 3$. Lastly, for *Task 5: Open-Ended Text Generation*, we select 1k entities in each knowledge graph and ask LLMs to perform open-ended generation⁴.

C Analysis (cont.)

C.1 Ranking

Model	Task					Domain			Avg.
	T1	T2	T3	T4	T5	YAGO	CPNet	UMLS	
ADA	8.3	9.7	6.1	5.1	4.8	†6.5	6.8	7.1	6.5
BABBAGE	7.0	6.0	5.0	5.0	3.8	5.7	5.5	†4.8	5.7
CURIE	8.7	9.3	<u>2.8</u>	4.0	2.7	†5.2	6.1	5.2	5.2
DAVINCI	<u>2.0</u>	<u>2.0</u>	1.7	1.6	3.0	†1.9	2.0	2.3	1.9
TURBO	1.0	1.0	3.0	<u>3.9</u>	<u>2.8</u>	<u>†2.3</u>	<u>2.4</u>	<u>2.3</u>	<u>2.3</u>
GPT-J	7.0	7.3	8.7	7.7	9.0	8.0	†7.6	8.1	8.0
OPT	9.0	7.0	8.0	7.8	9.8	†8.2	8.5	8.3	8.2
CHATGLM	4.7	3.0	4.0	7.1	3.8	4.3	†4.0	5.3	4.3
LLAMA	4.0	5.7	8.9	8.1	7.3	7.2	7.1	†6.1	7.2
ALPACA	3.3	4.0	6.9	4.8	7.8	5.6	†4.9	5.6	5.6

Table 6: Overall average rankings of ten LLMs on KGQUIZ across five tasks and three knowledge domains. **Bold**, underline represents the highest and the second highest ranking on each task (or KG). † denotes the KG on which each model has its best performance.

In addition to performance metrics, we calculate the average ranking of ten LLMs across the five tasks on three knowledge graphs and report results in Table 6. We observe that models with larger parameter sizes and more extensive training data, such as TEXT-DAVINCI-003 and GPT-3.5-TURBO, typically perform better in most tasks and domains, securing the first and second average rankings, respectively. Notably, individual models tend to excel in specific knowledge domains. For instance, TEXT-DAVINCI-003 achieves the highest average ranking on YAGO, which covers a wide array of general world knowledge, while GPT-3.5-TURBO shines on UMLS, which pertains primarily to the biomedical domain. This further underlines the value of our KGQUIZ benchmark in facilitating a more comprehensive evaluation of LLMs’ knowledge abilities, unraveling variations in performance across knowledge utilization and knowledge domains.

⁴For some tasks, we use in-context examples. More details in Appendix B

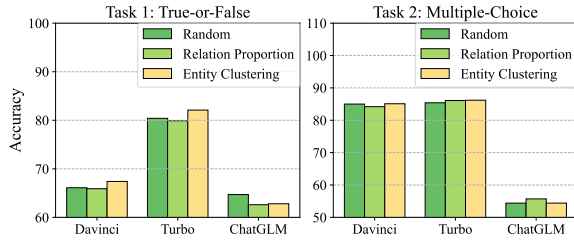


Figure 5: Comparison of model performance across different question sampling methods. Models are evaluated on 1,000 *Task 1: True-or-False* questions and 1,000 *Task 2: Multiple-Choice* questions sampled via three different methods: Random, Relation Proportion, and Entity Clustering. The results show the model’s performance is not significantly affected by the sampling method, indicating that randomly sampled triples are sufficient for capturing the features of the entire knowledge graph.

C.2 Question Sampling

In KGQUIZ, for each task, we generate questions by randomly sampling triplets (or head entities) from the KG, while whether the randomly sampled subsets is represented of the whole KG remain underexplored. To this end, we design two additional ways to sample a problem subset:

- **Relation Proportion:** We first calculate the proportion of relations in the KG, then sample triplets based on the relation distribution. This ensures that the proportion of relations in the sampled triples is consistent with the proportion of relations in the entire knowledge graph.
- **Entity Clustering:** First, we use knowledge graph embedding model TransE (Bordes et al., 2013) to obtain the embedding for each entity, then we use K-means to obtain 10 clusters of entities. We sample triplets based on the proportions of the number of entities in each cluster.

We generated 1,000 *Task 1: True-or-False* questions and 1,000 *Task 2: Multiple-Choice* questions on ConceptNet using these two methods respectively. According to Figure 5, we find that after changing to these two sampling methods that can theoretically better represent the features of the knowledge graph, the performance of each model did not change significantly (compared to random sampling). This indicates that randomly sampled triples can also reflect the features of the entire knowledge graph and the corresponding results are representative.

C.3 Knowledge Gap between LLMs and KGs

We conduct qualitative analysis on *Task 5: Open-Ended Text Generation* model outputs and present GPT-3.5-TURBO’s generated results and gold standard answers in Table 8. GPT-3.5-TURBO generated a total of 19 knowledge statements, of which 9 can be matched with triplets in YAGO. Among the remaining 10 knowledge statements that cannot be matched to YAGO, 8 of them are also found to be correct after manual annotation. This indicates that there is a knowledge gap between the parametric knowledge of LLMs and the structured knowledge of KGs. This also further emphasizes the necessity of considering knowledge utilization when discussing the role of KGs in augmenting LLMs. If general information about an entity is what we need, LLMs could provide mostly correct and factual answers; if LLMs need to perform tasks with the exact information in KGs, KG-augmented approaches could still be effective.

C.4 Negative Sampling Evaluation

Regarding the four negative sampling methods we proposed, a potential issue is that the sampled data may not be genuine negative samples. Therefore, in order to investigate the effectiveness of our negative sampling methods, we manually evaluated 20 samples for each method. In our manual evaluation, all the sampled examples were indeed true negative samples, which validated the effectiveness of our negative sampling methods.

C.5 Number of Options

Although extra answer options could serve as context information aid LLMs (as we analyzed in Section 4.2, we hypothesize that an increasing amount of distractors might sway LLMs away from the correct answer. To this end, we study the impact of the number of options on the difficulty of *Task 2: Multiple-Choice*. We follow the settings in Section 3 but change the number of options to 2, 3, 5, and 10 respectively. We present the performance of TEXT-DAVINCI-003 and GPT-3.5-TURBO on YAGO in Figure 6. We find that, although a small number of options providing extra context can give the model hints to answer questions, as the number of options increases, the model’s performance gradually declines due to the increasing number of distractors.

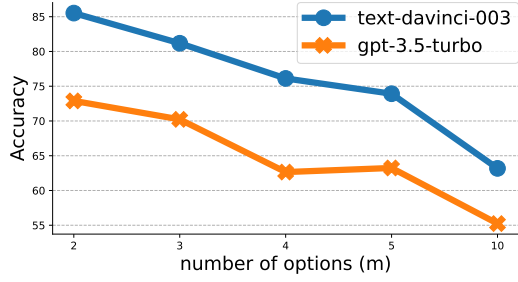


Figure 6: Impact of the number of answer options on LLM performance. The figure illustrates the performance of TEXT-DAVINCI-003 and GPT-3.5-TURBO on *Task 2: Multiple-Choice* (Multiple-Choice) using YAGO knowledge graph, with varying numbers of answer options (2, 3, 4, 5, and 10). The results show that as the number of options increases, the model’s performance declines, indicating that a higher number of distractors make the task more challenging. However, a small number of options can provide helpful contextual information.

C.6 Number of Hops

Task 4: Factual Editing investigates whether LLMs can correct factual mistakes in multi-hop knowledge reasoning chains. We additionally investigate whether the number of hops would affect the difficulty of the factual editing task. We generate 2-hop, 3-hop and 5-hop questions with triplets in YAGO and present the performance of textdavinci-003 and GPT-3.5-TURBO, shown in Figure 7. We observe that as the number of hops increases, the performance of textdavinci-003 improves, with the highest Semantic Match score (86.49) at 5 hops. This indicates that additional context from more hops can be beneficial in identifying and correcting factual inconsistencies in knowledge statements for this model. For GPT-3.5-TURBO, When the number of hops increases from 2 to 3, the performance of the model improves significantly. However, when the number of hops increases to 5, the performance of the model declines slightly but is still higher than that of 2 hops. This once again confirms that the impact of additional context from more hops on LLM performance in the factual editing task depends on the model.

C.7 Direct Triplets Generation

We use TEXT-DAVINCI-003 and GPT-3.5-TURBO to directly generate factual triplets about a certain entity (by giving an in-context example) and reported the precision and recall in Table 7. It can be observed that although the precision has im-

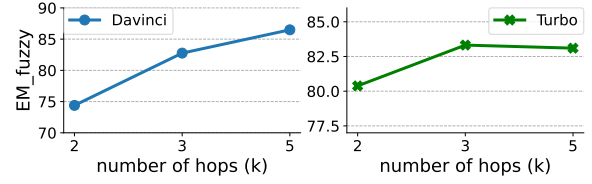


Figure 7: Effect of the number of hops on LLM performance in the Factual Editing task. The figure shows the Semantic Match scores for TEXT-DAVINCI-003 and GPT-3.5-TURBO on 2-hop, 3-hop, and 5-hop questions generated from YAGO KG. As the number of hops increases, the performance of TEXT-DAVINCI-003 improves, while the performance of GPT-3.5-TURBO exhibits a mixed pattern, indicating that the impact of the hop count on LLM performance varies depending on the model.

Model	Text		Triplets	
	Precision	Recall	Precision	Recall
DAVINCI	76.39	53.96	85.21	37.58
TURBO	77.28	57.63	91.42	37.21

Table 7: Comparison of precision and recall for open-ended text generation and direct triplets generation using TEXT-DAVINCI-003 and GPT-3.5-TURBO models. Direct Triplets Generation resulted in higher precision but lower recall than Open-Essay generation.

proved, the recall has dropped significantly. We analyzed that this is due to the model generating only a few high-confidence triplets when directly asked for triplets, which led to the aforementioned results. However, for other smaller-scale models, directly generating factual triplets is not feasible, as they cannot adequately understand the prompt’s instructions, resulting in poor performance.

C.8 Consistency Study

In Section 5.2, we investigate the robustness towards minor changes in prompts and knowledge statements. We present the five different prompts we used in Table 9.

Head	Gold	Matched	Factual	Unfactual
Mike Judge	{created, King of the Hill} {was born in, Guayaquil} {graduated from, University of California, San Diego} {directed, Office Space} {directed, Idiocracy} {directed, Extract (film)} {created, Office Space} {created, Idiocracy} {created, Extract (film)} {acted in, Office Space} {has gender, male} {lives in, Austin, Texas}	{creates, King of the Hill} {was born in, Guayaquil} {graduated from, University of California} {directs, Office Space} {directs, Idiocracy} {directs, Extract} {produces, Office Space} {produces, Idiocracy} {produces, Extract}	{creates, Beavis and Butt-Head} {creates, The Goode Family} {grew up in, New Mexico} {worked for, tech companies in Silicon Valley} {created, Frog Baseball} {won prize, Primetime Emmy Award} {won prize, Annie Award} {is known for, dry and satirical humor}	{started career as, programmer} {won prize, Peabody Award}
John Howard Northrop	{'was born in', 'Yonkers, New York'} {'graduated from', 'Columbia University'} {'works at', 'Rockefeller University'} {'has won prize', 'Nobel Prize in Chemistry'} {'died in', 'Wickenburg, Arizona'} {'works at', 'University of California, Berkeley'} {'has won prize', 'Daniel Giraud Elliot Medal'} {'has academic advisor', 'Thomas Hunt Morgan'} {'has won prize', 'National Medal of Science'} {'has gender', 'male'} {'is citizen of', 'United States'}	{'was born in', 'Yonkers'} {'earned a degree from', 'Columbia University'} {'worked at', 'Rockefeller Institute for Medical Research'} {'won the Nobel Prize in Chemistry in', '1946'} {'passed away in', 'Wickenburg'}	{'was a', 'biochemist'} {'shared the Nobel Prize with', 'James Sumner and Wendell Stanley'} {'worked on', 'isolation and crystallization of enzymes'} {'helped establish biochemistry as', 'a science'} {'conducted research on', 'enzymes'}	{'earned a PhD from', 'University of California'}

Table 8: Comparison between the generated answers by the GPT-3.5-TURBO model and the gold standard answers from the YAGO knowledge graph. The matched and factual columns indicate how well the model’s answers align with the ground truth and also highlight the factual answers not present in the knowledge graph, reflecting the knowledge gap between LLMs and KGs. The unfactual column shows model-generated answers that are not accurate.

ID	Prompt
1	Is the statement “[<i>Insert statement here</i>]” True or False?
2	Given the statement “[<i>Insert statement here</i>]”, is this factually correct? Please answer with True or False.
3	Assess the validity of this claim: “[<i>Insert statement here</i>]”. Respond with only True or False.
4	Is the following statement factually accurate? “[<i>Insert statement here</i>]” Provide your answer as either True or False.
5	Can you confirm if this statement is true or false? “[<i>Insert statement here</i>]”. Reply with just True or False.

Table 9: Five prompt templates we used to investigate the robustness towards minor changes in prompts and knowledge statements. We use the sampled knowledge statement to replace [*Insert statement here*] in each template and obtain 5 different prompts for the same knowledge statement.