

# KGC-RAG: Knowledge Graph Construction from Large Language Model Using Retrieval-Augmented Generation

Thin Prabhong<sup>1</sup>, Natthawut Kertkeidkachorn<sup>2</sup> and Areerat Trongratsameethong<sup>1</sup>

<sup>1</sup>*Chiang Mai University, Chiang Mai, Thailand*

<sup>2</sup>*Japan Advanced Institute of Science and Technology, Ishikawa, Japan*

## Abstract

The construction of Knowledge Graphs (KGs) has become increasingly important due to their ability to integrate and represent complex relationships across various domains, making them essential for applications like information retrieval and semantic search. Recently, Large Language Models (LLMs) have been utilized to enhance KGs creation by leveraging their advanced capabilities in understanding and generating human-like text. The Large Language Models for Knowledge Engineering (LLMKE) pipeline was introduced to combine knowledge probing with Wikidata entity mapping for knowledge engineering. Nevertheless, this approach has a limitation: it primarily relies on retrieval-augmented context drawn from the first paragraph and Wikipedia Infobox of the subject entity's page. This narrow focus can lead to incomplete knowledge representations, as relevant information is often spread throughout the text and linked pages. To address this issue, we propose the Knowledge Graph Construction from Large Language Model using Retrieval-Augmented Generation method (KGC-RAG). This method leverages web scraping to retrieve documents from the subject entity's Wikipedia page and to extend the search to include linked pages, thereby increasing the likelihood of capturing comprehensive and contextually rich information. We further enhance this approach by using LLMs in conjunction with cosine similarity to filter out irrelevant content, ensuring that only the most pertinent data are included in the relevant contexts. We conducted an experiment on datasets from ISWC 2024 LM-KBC Challenge and applied the meta-llama/Meta-Llama-3-8B-Instruct model as our pre-trained large language model along with the all-MiniLM-L6-v2 as our vector embedding model. We set a relevant score threshold of 0.5 to filter Wikipedia URLs. Our approach achieved macro average F1-scores of 0.695 and 0.698 on the validation and test sets, respectively. The implementation is available at <https://github.com/jaejeajay/LM-KBC2024>.

## 1. Introduction


Knowledge graphs [1] store information in a Subject-Predicate-Object format, providing more efficient semantic data storage compared to relational database. They are more easily understood by computers, offering flexibility and the ability to integrate diverse information, making them crucial for applications such as question answering and recommendation systems. A notable example of a significant knowledge graph repository is Wikidata [2], an open-source platform by the Wikimedia Foundation for storing factual data about the world.


---

*KBC-LM'24: Knowledge Base Construction from Pre-trained Language Models workshop at ISWC 2024*

✉ [thin.pra2013@gmail.com](mailto:thin.pra2013@gmail.com) (T. Prabhong); [natt@jaist.ac.jp](mailto:natt@jaist.ac.jp) (N. Kertkeidkachorn); [areerat.t@cmu.ac.th](mailto:areerat.t@cmu.ac.th) (A. Trongratsameethong)

🌐 <https://github.com/jaejeajay/LM-KBC2024> (T. Prabhong)

 © 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

Constructing knowledge graphs (KGs) is highly challenging [3] due to the need for access to vast and diverse information sources. Large Language Models (LLMs) are gaining popularity for their ability to perform a wide range of tasks, such as answering questions, providing information, translating languages, and engaging in conversations. Since LLMs are trained on extensive datasets, they serve as comprehensive repositories of knowledge. Extracting knowledge from LLMs can yield valuable information for various applications, including the construction of knowledge graphs, where LLMs can provide foundational data and relationships needed for building these complex structures.

However, challenges in extracting knowledge from LLMs [4] include the generation of fabricated answers and outdated information. Over time, facts can change or be disproven, but the data used to train LLMs remain factual only as of the time of training. Thus, optimizing LLMs for accurate and reliable responses is crucial.

Retrieval-Augmented Generation (RAG) [5] offers an effective solution to these challenges. By combining retrieval and generation processes, RAG enhances the accuracy of LLM outputs by retrieving up-to-date information from external databases or sources, thereby providing the necessary context for generating more factual and reliable responses.

In this study, we explore the construction of knowledge graphs using knowledge extracted from LLMs. We optimize LLM performance with RAG by improving the quality of relevant context through web scraping and web crawling. The relevance score is determined by the path names in the Wikipedia URLs. We utilized datasets from the ISWC 2024 LM-KBC Challenge [6], where the task limited LLM parameters to 10B. We selected the Llama-3-8B-Instruct model for our study, and the results demonstrate that our approach outperforms the baseline.

## 2. Related Works

### 2.1. Using Large Language Models for Knowledge Engineering (LLMKE): A Case Study on Wikidata

In recent research, Zhang et al. [7] explored the use of LLMs for knowledge engineering tasks within the context of the ISWC 2023 LM-KBC Challenge. They utilized pre-trained LLMs to generate relevant objects in string format from given subject-relation pairs sourced from Wikidata, subsequently linking these objects to their respective Wikidata QIDs. The developed pipeline, known as Large Language Models for Knowledge Engineering (LLMKE), combines knowledge probing with Wikidata entity mapping and incorporates retrieval-augmented context to enhance predictions. This context is derived from external sources, such as Wikipedia, to refine the model's responses by comparing and integrating this information with the model's initial predictions.

Nonetheless, LLMKE operates under the assumption that the most relevant documents are primarily found in the first paragraph (Introduction) and the Wikipedia Infobox, which may not always be the case, as related information can be dispersed throughout various sections.

To address this limitation, we propose a method to enhance content extraction from Wikipedia by employing web scraping techniques to gather information from paragraphs, Infoboxes, and Wikitable on the Wikipedia page of the subject entity. Also, this approach extends the search to other linked Wikipedia pages, thereby capturing more comprehensive and relevant information.

One challenge of extending the search beyond the subject entity's Wikipedia page is the increased volume of documents, which can slow down the process. To mitigate this, we implemented a solution using LLM and cosine similarity scores to filter and prioritize relevant Wikipedia pages and documents, ensuring efficient and effective information retrieval.

## 2.2. Retrieval-Augmented Generation for Large Language Models

Yunfan Gao et al. [8] reviewed the RAG paradigm, categorizing it into three types: Naïve, Advanced, and Modular RAG. Naïve RAG involves a simple two-phase process of retrieving documents and generating responses based on them. Advanced RAG improves this by optimizing query modification and refining the retrieved context for better LLM performance. Modular RAG further breaks down the process into independent modules, allowing easier updates but requiring more resources for development.

This study focuses on utilizing the Naïve RAG process to explore its application in constructing knowledge graphs through knowledge extraction from LLMs, using the ISWC 2024 LM-KBC Challenge dataset [6]. The research integrates LLMs into the retrieval stage to aid in filtering relevant documents from web scraping, rather than limiting their role to the generation phase. Given that LLMs are trained on vast knowledge, the hypothesis is that incorporating them in the Retrieval process may outperform using relevance scores alone.

## 3. Methods

The Knowledge Graph Construction from Large Language Model using Retrieval-Augmented Generation method (KBC-RAG) provides an overview as shown in Figure 1. This method is composed of two processes: RAG and Entity Mapping. In the RAG process, relevant contexts are identified and used to enhance the LLM's knowledge extraction. In the entity mapping process, the answers from the LLM are used to construct knowledge graphs by mapping through API functions from Wikidata, resulting in completed knowledge tuples. The details of each process are as follows:

### 3.1. Retrieval-Augmented Generation (RAG)

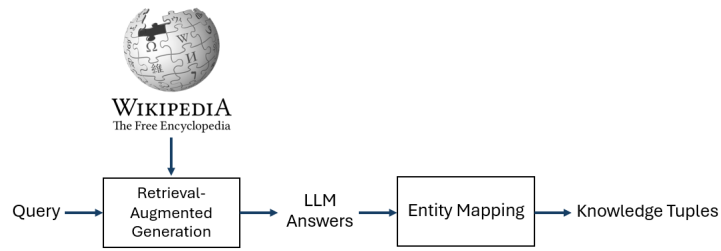
The overview of RAG process is presented in Figure 2. The goal of this process is to extract knowledge from LLM using RAG as an LLM optimization method. The outcome of this process will be LLM-generated answers, which will then be used for subsequent entity mapping process.

#### 3.1.1. Entity ID Query

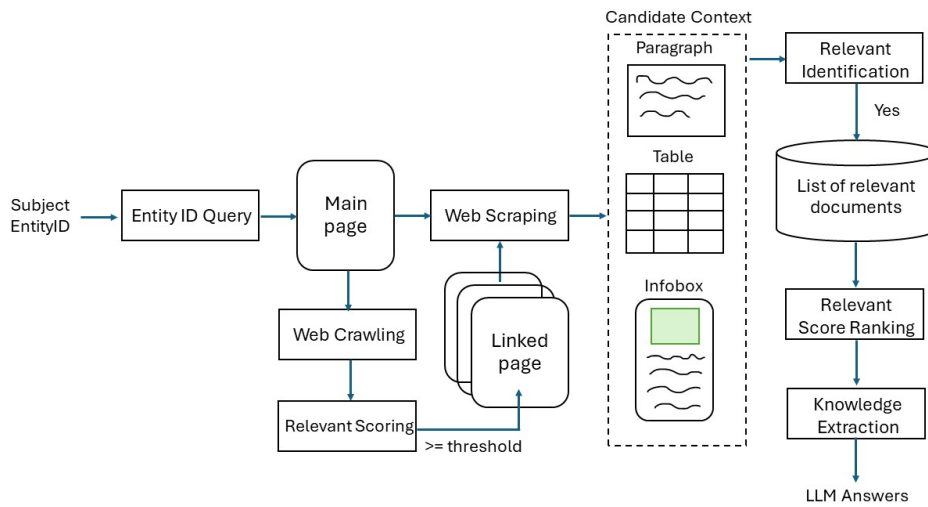
The entity ID query step involves locating the Wikipedia URLs of the subject entities using the provided SubjectEntityID, which is essential for the web scraping process. The SubjectEntityID is used to search for the URLs in Wikidata through an API called SPARQLWrapper <sup>1</sup>, with a preference for English URLs. The outcome of this step is the subject entity's Wikipedia URL, which is used to access the subject entity's Wikipedia page, referred to as the main page in

---

<sup>1</sup><https://sparqlwrapper.readthedocs.io/en/latest/main.html>



**Figure 1:** The overview of the KGC-RAG method



**Figure 2:** The overview of Retrieval-Augmented Generation (RAG)

this study. For example, if SubjectEntity is “Uruguay”, SubjectEntityID is “Q77”, and relation is “countryLandBordersCountry”, after using SubjectEntityID to search via SPARQLWrapper, the resulting Wikipedia URL for SubjectEntity would be <https://en.wikipedia.org/wiki/Uruguay>.

### 3.1.2. Web Scraping

In web scraping step, the focus is on scraping data from the Wikipedia pages. The data targeted for scraping includes three tags: *p* (Paragraph), Infobox, and Wikitable. The outputs of this step are the scraped data from Wikipedia, which will later be assessed for relevance in the relevant identification step.

**Table 1**  
Templates of Prompt for web scraping

Relation	Prompt for web scraping
awardWonBy	“Who has won ” + SubjectEntity + “ ?”
companyTradesAtStockExchange	“Which stock exchange does ” + SubjectEntity + “ trade on ?”
countryLandBordersCountry	“Which country share land border with ” + SubjectEntity + “ ?”
personHasCityOfDeath	“What is the city of death of ” + SubjectEntity + “ ?”
seriesHasNumberOfEpisodes	“How many episodes does series ” + SubjectEntity + “ has ?”

### 3.1.3. Web Crawling

The web crawling step is designed to broaden the scope of web scraping by extending the search from the main page to other linked Wikipedia pages, thereby increasing the chances of identifying relevant context. This is achieved by identifying *<a>* tags within the main page. The result of this step is a collection of linked Wikipedia URLs linked to the main page. For example, from the Wikipedia page of the SubjectEntity “Uruguay” at <https://en.wikipedia.org/wiki/Uruguay>, examples of linked Wikipedia pages include [https://en.wikipedia.org/wiki/Economy\\_of\\_Uruguay](https://en.wikipedia.org/wiki/Economy_of_Uruguay), [https://en.wikipedia.org/wiki/Religion\\_in\\_Uruguay](https://en.wikipedia.org/wiki/Religion_in_Uruguay), and so on.

### 3.1.4. Relevant Scoring

The primary aim of relevant scoring step is to filter out unnecessary linked URLs from the web crawling process by using a Relevant score, calculated via cosine similarity. This approach helps streamline the RAG process. The filtering method involves calculating the cosine similarity between the path name of each URL linked to the main page and the Prompt for web scraping, utilizing the vector model all-MiniLM-L6-v2. The templates of the prompt for web scraping are shown in Table 1. For instance, from the URL [https://en.wikipedia.org/wiki/Economy\\_of\\_Uruguay](https://en.wikipedia.org/wiki/Economy_of_Uruguay), only the part “Economy of Uruguay” is used to calculate the Relevant Score. URLs with a cosine similarity score above the threshold are selected for web scraping. The outcome of this step is a refined list of linked URLs for web scraping.

### 3.1.5. Relevant Identification

The relevant identification step focuses on filtering scraped data from Wikipedia by leveraging LLM responses to specific questions. The result is a list of relevant documents, which will be used in the subsequent step for determining the relevant score.

We begin by filtering each paragraph, Infobox, and Wikitable data using the Meta-Llama-3-8B-Instruct Model through question-based evaluation. The format of question is: *Is this information “[Paragraph/Infobox/Wikitable]” able to answer the question: “[Prompt for web scraping]”?* Information deemed relevant by the LLM is added to the list of relevant documents. If the LLM determines that it can answer the question, it will return a response of 1; otherwise, it will return 0. The answer from the LLM is 1, meaning that this paragraph will be included in the list of relevant documents.

**Table 2**  
Templates of Prompt for Knowledge Extraction

Relation	Prompt for Knowledge Extraction
awardWonBy	“Provide a name only list of all award winners in ” + SubjectEntity + “ with no explanation and name with comma ?”
companyTradesAtStockExchange	“Which stock exchange does ” + SubjectEntity + “ trade on ? answering with no explanation and name with comma? If None, answer None.”
countryLandBordersCountry	“Which countries share land borders with ” + SubjectEntity + “ with country name only with comma? If None, answer None.”
personHasCityOfDeath	“What is the city of death of ” + SubjectEntity + “? answering with one city name only with no explanation. If there is no place of death mentioned, answer None”
seriesHasNumberOfEpisodes	“How many total episodes of series ” + SubjectEntity + “ ? answering with only one number ?”

### 3.1.6. Relevant Score Ranking

In relevant score ranking step, the goal is to retrieve the relevant context from the list of relevant documents by applying the same method used in the relevant scoring step, combined with relevant score ranking. Each document is compared to the web scraping prompt by calculating cosine similarity using the vector model *all-MiniLM-L6-v2*, as in the relevant scoring step, to determine its relevant score. The documents are then ranked from highest to lowest, and the Top *K* documents are selected to form the relevant context, which will be crucial for knowledge extraction. The key difference between this step and the relevant scoring step lies in their objectives. In relevant scoring, the focus is on identifying relevant linked URLs using a threshold-based filtering technique, while this step selects relevant documents to create the relevant context through a Top *K* filtering approach. The Top *K* method is not used in the relevant scoring step to avoid restricting the number of links, thereby allowing broader web scraping and increasing the chances of finding relevant documents.

### 3.1.7. Knowledge Extraction

The purpose of knowledge extraction step is to extract knowledge from LLM by using relevant context. The result of this step will be LLM-generated answers that are ready for entity mapping. For example, for the question “Which countries share land borders with Uruguay with country name only with comma? If None, answer None.” the answer generated by the LLM is “Argentina, Brazil”.

To ask questions using the prompt for knowledge extraction, as shown in Table 2, and the obtained relevant context, the input message format consists of three parts:

1. **Relevant Context:** The first message provides the relevant context to the LLM in the following format:

{“role”: “system”, “content”: “Using this context to answer the question: ” + [Relevant Context]}.

**Table 3**  
Domain and Range for Each Relation

Relation	Domain	Range
awardWonBy	Award	Human
companyTradesAtStockExchange	Company	Stock Exchange
countryLandBordersCountry	Country	Country
personHasCityOfDeath	Human	City
seriesHasNumberOfEpisodes	Series	Number

2. **Behavior Setting:** The second message sets the behavior of the Chatbot of LLM to ensure responses are as required:

*{“role”: “system”, “content”: “You are a chatbot who always responds an answer in english with comma and no explanation. If u don’t know the answer, answer None”}.*

3. **Question Prompt:** The third message is used to ask the question and is formatted as:

*{“role”: “user”, “content”: [Prompt for Knowledge Extraction]}.*

This structured approach ensures that the LLM has the necessary context and guidance to provide accurate and concise answers.

### 3.2. Entity Mapping

The goal of entity mapping process is to construct a knowledge graph using the answers from the LLM obtained in the previous session. The result of this process will be completed knowledge tuples based on the given subject entity and relation. For example, after receiving the answer from the question *“Which countries share land borders with Uruguay with country name only with comma? If None, answer None.”* from the LLM as *“Argentina, Brazil”*, we map the answer using *wbsearchentities*. The final result is a completed tuple: *{“SubjectEntity”: “Uruguay”, “Relation”: “countryLandBordersCountry”, “ObjectEntitiesID”: [“Q414”, “Q155”]}*.

We begin by mapping the answers to create a knowledge graph using *wbsearchentities*, an API function provided by the Wikidata Query Service <sup>2</sup>. This function searches for entities in Wikidata using labels or keywords and returns a list of matching entities ranked by relevance. The method used to select entities for linking as objects is “Choose First” for all relations. After selecting an entity, its validity is verified by checking the range of each relation. As shown in Table 3, if the entity’s “instance of” is within the same range of relations as the subject, the entity is selected for mapping.

In cases where responses from LLM exhibit ambiguity, it may be due to the LLM relying too heavily on context. For instance, in the relation *companyTradesAtStockExchange*, stock exchange information for a company on Wikipedia is often presented in the format “Traded as”, such as *“West Japan Railway Company traded as TYO: 9021”*, where TYO is the stock exchange name and 9021 is the stock exchange code. The LLM might provide an answer like “TYO: 9021”.

<sup>2</sup>[https://www.wikidata.org/wiki/Wikidata:SPARQL\\_query\\_service](https://www.wikidata.org/wiki/Wikidata:SPARQL_query_service)

**Table 4**  
Number of unique subject-entities in the data splits

Relation	Train	Val	Test	Special features
awardWonBy	10	10	10	Many objects per subject
companyTradesAtStockExchange	100	100	100	Null values possible
countryLandBordersCountry	68	68	68	Null values possible
personHasCityOfDeath	100	100	100	Null values possible
seriesHasNumberOfEpisodes	100	100	100	Object is numeric

To address this ambiguity, we employ hard coding to filter and refine the responses from the LLM, thereby improving the accuracy and clarity of the output.

## 4. Experiments and Results

### 4.1. Datasets

The datasets used in this study were obtained from ISWC 2024 LM-KBC Challenge [6]. Details of the datasets are shown in Table 4. Each dataset consists of 378 unique subject-entities and 5 relations, with object entities referenced from Wikidata. Additionally, each relation has special features that describe the characteristics of the object entities, which vary by relation. For example, the special feature of the relation *countryLandBordersCountry* is the possibility of Null values. For instance, New Zealand, being an island nation with no land borders, would have a Null value for the object entity corresponding to the given subject entity “New Zealand” and the given relation *countryLandBordersCountry*.

### 4.2. Experiment Settings

The LLM used in this study is Llama-3-8B-Instruct. Its parameter count also falls within the required limit of the task. Wikipedia was the primary data source, and the vector model chosen for cosine similarity calculations was all-MiniLM-L6-v2. This model was selected for its compact size, speed, and reasonable performance, particularly given the resource-intensive nature of web scraping and crawling, which may involve processing numerous linked URLs and a large volume of documents. It helps efficiently manage this challenge. The relevant score threshold for filtering Wikipedia URLs is set at 0.5 to select URLs similar to the web scraping prompt, which functions as a user query. A cosine similarity of 0.5 reflects a moderate level of similarity, balancing relevance and coverage to avoid overly filtering out pertinent URLs. Each document was divided into segments of 4,500 tokens, and a Top  $K$  approach, with  $K=20$ , was applied to combine documents into a single context within the relevant score ranking pipeline. In the relevant identification step, we used  $\text{max\_new\_tokens}=1$ ,  $\text{temperature}=0.1$ , and  $\text{top\_p}=0.9$ .  $\text{max\_new\_tokens}=1$  was chosen because this step of the pipeline only returns a 0 or 1 for filtering. In the knowledge extraction step, we set  $\text{max\_new\_tokens}=3000$ ,  $\text{temperature}=0.1$ , and  $\text{top\_p}=0.9$ , with  $\text{max\_new\_tokens}=3000$  being an estimated value for the maximum length of the LLM’s responses. The choice of using  $K=20$ ,  $\text{temperature}=0.1$ , and  $\text{top\_p}=0.9$  was



based on preliminary experiments, which demonstrated better results compared to not setting these values. Moreover, the values of `temperature=0.1` and `top_p=0.9` are suitable for tasks requiring logical consistency.

### 4.3. Baseline

Our study used the baseline of meta-llama/Meta-Llama-3-8B-Instruct from ISWC 2024 LM-KBC Challenge<sup>1</sup>. The baseline was derived from the use of a Masked Language Model, an Autoregressive Language Model, and Llama-3 chat models.

### 4.4. Evaluation Metrics

In our study, we used three evaluation metrics: Macro-Precision (Macro-P), Macro-Recall (Macro-R), and Macro-F1. Macro-P is the average precision score across all classes, calculated by determining the precision for each class and averaging the values. Macro-R represents the average recall score across all classes, computed similarly by averaging the recall for each class. Finally, Macro-F1 is the average F1 score across all classes, derived by calculating the F1 score for each class and averaging the results.

### 4.5. Results

We conducted experiments on our system using the validation set, with the results presented in Table 5. The experimental results demonstrated that our approach outperformed the baseline, achieving an average F1-score of 0.695. For the test set, we submitted the results to the ISWC 2024 LM-KBC Challenge. The test results, shown in Table 6, indicate an average F1-score of 0.698, which is consistent with the performance on the validation set.

## 5. Discussion

### 5.1. Data Source Quality

The source of context information plays a critical role in determining the coverage score of the context [9]. If a data source provides limited information, the coverage score will be lower, which subsequently affects the responses generated by the LLM [10]. ISWC 2024 KM-KBC Challenge this year is represented by five distinctive relations, particularly the *awardWonBy* and *companyTradesAtStockExchange* relations. For the *awardWonBy* relation, some subjects on Wikipedia do not display the award winners on the main award page, but instead on a separate “List of laureates” page. Similarly, for the *companyTradesAtStockExchange* relation, information about stock exchanges may not be directly provided or may be absent from the Wikipedia page. Therefore, utilizing effective data sources and retrieval methods is crucial.

### 5.2. Knowledge Discrepancy

While performing the task, we observed that although the responses generated by the LLM were consistent with the context derived from Wikipedia, they did not align with the information

**Table 5**

The Results of Our System and Baseline (meta-llama/Meta-Llama-3-8B-Instruct) on Validation Set

Method	Relation	macro-p	macro-r	macro-f1
Baseline	awardWonBy	0.238	0.028	0.045
	companyTrades			
	AtStockExchange	0.540	0.703	0.474
	countryLand			
	BordersCountry	0.961	0.912	0.919
	personHas			
	CityOfDeath	0.700	0.600	0.460
	seriesHas			
	NumberOfEpisodes	0.493	0.160	0.155
	<b>All Relations</b>	<b>0.638</b>	<b>0.552</b>	<b>0.455</b>
KGC-RAG	awardWonBy	0.771	0.103	0.125
	companyTrades			
	AtStockExchange	0.842	0.795	0.678
	countryLand			
	BordersCountry	0.921	0.900	0.850
	personHas			
	CityOfDeath	0.750	0.910	0.660
	seriesHas			
	NumberOfEpisodes	0.725	0.700	0.697
	<b>All Relations</b>	<b>0.799</b>	<b>0.801</b>	<b>0.695</b>

**Table 6**

The results of our system on test set

Relation	macro-p	macro-r	macro-f1
awardWonBy	0.825	0.021	0.032
companyTradesAtStockExchange	0.797	0.755	0.638
countryLandBordersCountry	0.910	0.903	0.854
personHasCityOfDeath	0.695	0.920	0.647
seriesHasNumberOfEpisodes	0.800	0.770	0.770
<b>All Relations</b>	<b>0.792</b>	<b>0.810</b>	<b>0.698</b>

referenced in Wikidata. This discrepancy affected the system’s performance. For example, in the training data for the subject: *Love & Hip Hop: New York*, relation: *seriesHasNumberOfEpisodes*, the LLM predicted that the series had 143 episodes, based on Wikipedia. However, Wikidata indicated that the series had 82 episodes, which was outdated. Although both Wikipedia and Wikidata are open-source platforms where users can update information, the discrepancy in information between them still exists. Helping to update information or reporting issues to the community could help reduce this discrepancy.

**Table 7**

The average coverage scores on validation set

Relation	Without Web Crawling	With Web Crawling
awardWonBy	0.361	0.137
companyTradesAtStockExchange	0.213	0.675
countryLandBordersCountry	0.731	0.967
personHasCityOfDeath	0.480	0.715
seriesHasNumberOfEpisodes	0.860	0.790
<b>All Relations</b>	0.529	0.657

### 5.3. Relevant Context

High-quality context can improve the performance of question-answering in LLMs [11]. We employed LLMs and relevant scores to filter the quality of documents obtained from web scraping and web crawling. After consolidating documents into a single context, we evaluate the quality of the relevant context using the coverage score. The coverage score indicates how well the relevant context contains substrings that are object entities from the given subject entity and relation. It is calculated by dividing the number of object entity substrings found in the relevant context by the total number of object entities for that given subject entity and relation. The coverage scores for each relation on the validation set are shown in Table 7.

It was observed that for certain relations, particularly awardWonBy, the coverage score was not high. This may be due to the fact that Wikipedia often presents award winners in tables that include additional information. Although the award winner’s name is present, the presence of other noise in the data can impact the relevant score [12]. When comparing the coverage score with the macro F1 score for each relation in Table 5, it was found that the macro F1 score is generally lower or similar. This implies that increasing the coverage score of the context could potentially enhance the quality of responses generated by the LLM.

Additionally, comparing the difference in coverage score between using web crawling and not using it reveals that web crawling generally improves the coverage score for most relations. However, in some cases, it does not. This may be due to the presence of noise in the documents obtained from web scraping. Expanding the scope of web scraping increases the likelihood of encountering noisy data, which can affect the cosine similarity score. As a result, some documents with higher cosine similarity scores may be selected over those that actually contain the correct answers to the query.

## 6. Conclusion

This study aims to construct KGs from LLM by using RAG to optimize knowledge extraction from LLM and to involve the large language model in finding relevant documents using ISWC 2024 LM-KBC Challenge datasets [6]. We achieved an F1 score of 0.695 for the validation set and 0.698 for the test set. The results demonstrate that the quality of context is crucial for optimizing LLM performance in answering questions. Additionally, the LLM can effectively assist in screening relevant documents, which is a key factor in constructing an accurate and

high-quality knowledge graph. For future investigations, it is recommended to explore the implementation of automatic relevant document retrieval instead of relying solely on question-answering combined with relevant scores.

## Acknowledgement

This work was supported by JSPS Grant-in-Aid for Early-Career Scientists (Grant Number 24K20834).

## References

- [1] H. Paulheim, Knowledge graph refinement: A survey of approaches and evaluation methods, *Semantic Web* 8 (2017) 489–508.
- [2] D. Vrandečić, M. Krötzsch, Wikidata: A free collaborative knowledge base, *Communications of the ACM* 57 (2014) 78–85.
- [3] L. Jansen, R. van Hirtum, Constructing knowledge graphs from text: A survey of methods and tools, *Journal of Computer Science and Technology* 35 (2020) 1001–1020.
- [4] R. Binns, V. Veitch, N. Shadbolt, Evaluating the reliability of large language models for knowledge extraction, in: *Proceedings of the ACM Conference on Information and Knowledge Management (CIKM)*, 2021.
- [5] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, et al., Retrieval-augmented generation for knowledge-intensive nlp tasks, in: *Advances in Neural Information Processing Systems*, volume 33, 2020, pp. 9459–9474.
- [6] J.-C. Kalo, S. Razniewski, T.-P. Nguyen, B. Zhang, Knowledge base construction from pre-trained language models 2022, in: *Semantic Web Challenge on Knowledge Base Construction from Pre-trained Language Models*, CEUR-WS, 2024.
- [7] B. Zhang, I. Reklos, N. Jain, A. M. Peñuela, E. Simperl, Using large language models for knowledge engineering (llmke): A case study on wikidata, in: *Proceedings of the ISWC 2023*, 2023.
- [8] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, M. Wang, H. Wang, Retrieval-augmented generation for large language models: A survey, 2024.
- [9] X. Chen, J. Zhu, Enhancing context coverage for question answering with multiple knowledge sources, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020, pp. 2340–2350.
- [10] H. Zhao, M. Eskenazi, Understanding and mitigating the impact of data source limitations on llm performance, in: *Proceedings of NeurIPS*, 2021.
- [11] T. Kwiatkowski, J. Palomaki, M. Redfield, M. Edward, M. Collins, et al., Natural questions: a benchmark for question answering research, *Transactions of the Association for Computational Linguistics* 7 (2019) 453–466.
- [12] Y. Jernite, J. M., Analyzing the effects of noise on neural network performance in natural language processing, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017, pp. 2892–2900.