

Fine-Grained Retrieval-Augmented Generation for Visual Question Answering

Zhengxuan Zhang¹, Yin Wu¹, Yuyu Luo^{1,2}, Nan Tang^{1,2*}

¹The Hong Kong University of Science and Technology (Guangzhou)

²The Hong Kong University of Science and Technology

{zzhang393, ywu450}@connect.hkust-gz.edu.cn {yuyuluo, nantang}@hkust-gz.edu.cn

Abstract

Visual Question Answering (VQA) focuses on providing answers to natural language questions by utilizing information from images. Although cutting-edge multimodal large language models (MLLMs) such as GPT-4o achieve strong performance on VQA tasks, they frequently fall short in accessing domain-specific or the latest knowledge. To mitigate this issue, retrieval-augmented generation (RAG) leveraging external knowledge bases (KBs), referred to as KB-VQA, emerges as a promising approach. Nevertheless, conventional unimodal retrieval techniques, which translate images into textual descriptions, often result in the loss of critical visual details. This study presents fine-grained knowledge units, which merge textual snippets with entity images stored in vector databases. Furthermore, we introduce a knowledge unit retrieval-augmented generation framework (KU-RAG) that integrates fine-grained retrieval with MLLMs. The proposed KU-RAG framework ensures precise retrieval of relevant knowledge and enhances reasoning capabilities through a knowledge correction chain. Experimental findings demonstrate that our approach significantly boosts the performance of leading KB-VQA methods, achieving improvements of up to 10%.

1 Introduction

Knowledge-based Visual Question Answering (KB-VQA) extends traditional Visual Question Answering (VQA) by incorporating external knowledge to answer questions where image information alone is insufficient (Marino et al., 2019; Lin et al., 2022; Wen et al., 2024). However, traditional methods often face limitations in their ability to perform complex reasoning over both visual content and external knowledge sources, as they typically rely on predefined retrieval mechanisms or specific training data (Wu and Mooney, 2022; Yang et al., 2023).

Recently, the emergence of multimodal large language models (MLLMs), such as GPT-4 (Achiam et al., 2023) and LLaVA (Liu et al., 2023a), has introduced new possibilities for VQA. Unlike previous methods, MLLMs serve not only as powerful reasoning engines but also as vast knowledge repositories, with information learned from world knowledge during pretraining (Wang et al., 2024; Liu et al., 2023b). This dual capability enables more nuanced answers. However, the knowledge acquired during training is general and (maybe outdated) world knowledge, limiting the model’s ability to respond to domain-specific and update-to-date queries. As shown in Figure 1(a), when using GPT-4 to ask a question about the bridge in the image, it fails to provide an answer due to a lack of relevant knowledge and LLaVA even hallucinated and provided a “false” answer.

At this point, it becomes necessary to employ KB-VQA approaches, by retrieving information from a database – a process also known as Retrieval-Augmented Generation (RAG) in the context of LLMs (Fan et al., 2024). This typically involves converting images into captions and then performing passage-level retrieval combined with the query. However, this method struggles to handle fine-grained information for question answering, and during the image-to-text modality conversion process, some visual details are inevitably lost. As shown in Figure 1(b), a unimodal, coarse-grained approach fails to retrieve the relevant knowledge.

Intuitively, in order to accurately find the knowledge corresponding to this bridge, it is necessary to identify the corresponding images through its visual features and then look through the information behind it, as illustrated in Figure 1(c).

Following this approach, we propose a “Knowledge Unit” component to bridge the query and specific knowledge. Specifically, we propose a Knowledge Unit Retrieval-Augmented Generation

* Nan Tang is the corresponding author

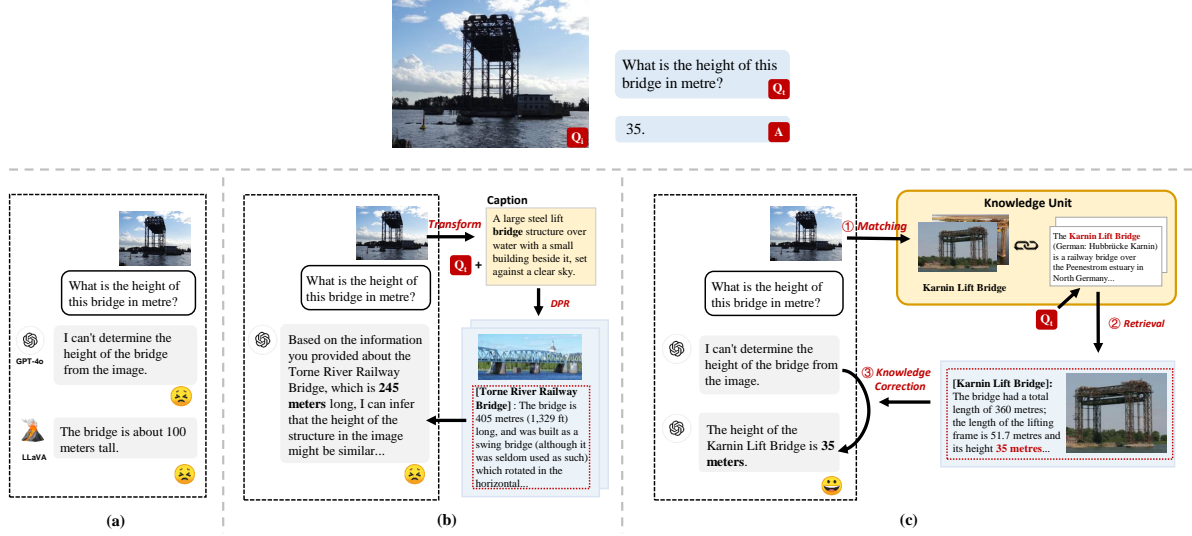


Figure 1: Sample VQA Solution with MLLMs: (a) Direct answer without additional knowledge. (b) Single-modality coarse-grained RAG (KB-RAG). (c) Our proposal **KU-RAG**.

(KU-RAG) method, which is a multimodal, fine-grained, zero-shot retrieval approach covering both data storage and retrieval. As shown in the figure, our method matches the image from the question with images in the database, identifying the relevant knowledge (e.g., “Karnin Lift Bridge”) and subsequently retrieving fine-grained knowledge chunks related to it to answer the question more effectively. Finally, we designed a Knowledge Correction Chain (KCC) to assist in answer generation. The KCC integrates retrieved information into the MLLM’s reasoning process and verifies the accuracy of the knowledge generated by MLLMs.

Our contributions are summarized as follows.

- We introduce *knowledge units* that consist of fine-grained multimodal data fragments (e.g., text fragments, entity images, and so on) from external knowledge bases. We store these knowledge units in vector databases, in order to effectively retrieve relevant knowledge for a given visual-based query.
- We propose a knowledge unit retrieval-augmented generation (KU-RAG) method, which retrieves fine-grained knowledge units, employs a knowledge correction chain (KCC) during query inference, and achieves zero-shot for combining retrieved knowledge units with MLLMs.
- The experimental results on various KB-VQA benchmarks demonstrate that our method ef-

fectively enhances the performance of MLLM on this task.

2 Related Work

2.1 Knowledge-based Visual Question Answering

Knowledge-based Visual Question Answering (KB-VQA) uses external knowledge like Wikipedia to assist in answering image-related questions (Wu et al., 2022). Early methods used a retriever-reader framework, retrieving relevant text to predict answers, but struggled with fine-grained questions about visual entities (Hu et al., 2023a), events (Yang et al., 2023), and visual information (Chen et al., 2023).

With the rise of LLMs, some methods combined implicit LLM knowledge with database retrieval, converting images into tags or captions and using GPT for knowledge retrieval (Gui et al., 2021; Chen et al., 2024; Wu et al., 2025). However, a gap exists between queries and LLM knowledge. Hu et al. (2023b) proposed a prompt-guided image captioning method to adjust captions based on queries. Although some methods reduce information loss through visual features (Salaberria et al., 2023) or enriched prompts (Shao et al., 2023), they don’t fully resolve the issue. Our approach integrates multimodal information for both retrieval and reasoning, with the LLM acting as both a knowledge base and a reasoner.

2.2 Multimodal Retrieval-augmented Generation

Multimodal Retrieval-augmented Generation (RAG) improves LLMs by retrieving relevant external documents to address domain-specific knowledge and avoid hallucinations (Gao et al., 2023). The process retrieves documents and generates answers by integrating the information, similar to KB-VQA’s retrieval. RAG has evolved with models like GRAG, which improves retrieval relevance (Hu et al., 2024), FiD-RAG, which fuses multiple documents during generation (Izacard and Grave, 2020), and DPR-RAG, which uses dense retrieval to find relevant fragments quickly (Karpukhin et al., 2020).

Unlike the aforementioned methods, our goal is to build a multimodal, fine-grained model that combines the internal knowledge of LLMs to generate high-quality answers in zero-shot scenarios.

3 Our Methodology: VQA with Knowledge Unit RAG and MLLMs

In this section, we will provide a detailed explanation of the **Knowledge Unit Retrieval-Augmented Generation (KU-RAG)** method. We construct a novel framework based on KU-RAG for KB-VQA, as illustrated in Figure 2. We will introduce the KB-VQA task in Section 3.2, then present our proposed “knowledge unit” in Section 3.2, followed by a complete overview of KU-RAG in Section 3.3.

3.1 Task Definition

Visual Question Answering (VQA): Given a question Q , which consists of an image Q_i and a textual question Q_t related to the content of the image, the task of VQA is to generate an answer A based on the information available in the image and the text. In this setup, the system aims to understand both the visual and textual aspects of the input and provide a relevant response.

Knowledge-Based Visual Question Answering (KB-VQA): In KB-VQA, the goal extends beyond the basic VQA task by incorporating external knowledge K stored in knowledge bases (KBs) to answer the question. This external knowledge, which can be categorized as either image knowledge K_i or text knowledge K_t , is retrieved based on the question Q and is used to generate a more informed and accurate answer A .

3.2 Knowledge Unit

3.2.1 Definition

We introduce a new structure called **Knowledge Unit (KU)**. Each KU serves as a knowledge carrier or object generated in combination with the query, such as entities, events, rules, topics, etc., designed to bridge the gap between the query and the database during the actual question-answering process.

For a piece of knowledge, the three most important factors are its image, its name, and detailed textual knowledge. Therefore, we designed each KU as a triplet, consisting of Knowledge Image (K_i), Knowledge Name (K_n), and Knowledge Text (K_t):

$$KU = \{K_i, K_n, K_t\} \quad (1)$$

In the KB-VQA task, image-image or name-name matching is typically used to determine which piece of knowledge a given image belongs to. Hence, we encapsulate K_i and K_n into the ‘**Matching End**’ to link the query and the KU. The purpose of KB-VQA often involves querying the knowledge behind an image, so we refer to K_t as the ‘**Detail End**’.

3.2.2 Construction

Knowledge Predefinition. Firstly, we should determine how to extract the knowledge unit with the application scenario and consider the query and database. For example, in an object recognition QA system, different entities can serve as knowledge units; in an event query system, different events can serve as knowledge units; in a corporate rules and regulations query system, a rule-based knowledge unit should be constructed. It is important to note that knowledge units do not necessarily need to be atomic or as fine-grained as possible. Its division should be determined based on the granularity of the data in the query and the database. For instance, a general animal knowledge QA system may only require the general species of an animal (e.g., “cat”), whereas a specific species QA system may require the specific species name (e.g., “Persian cat”).

Knowledge Segment. Since subsequent steps involve the storage of the knowledge unit, storing textual knowledge K_t within the detail end at the document level, which is a coarse-grained storage method, is highly detrimental to knowledge

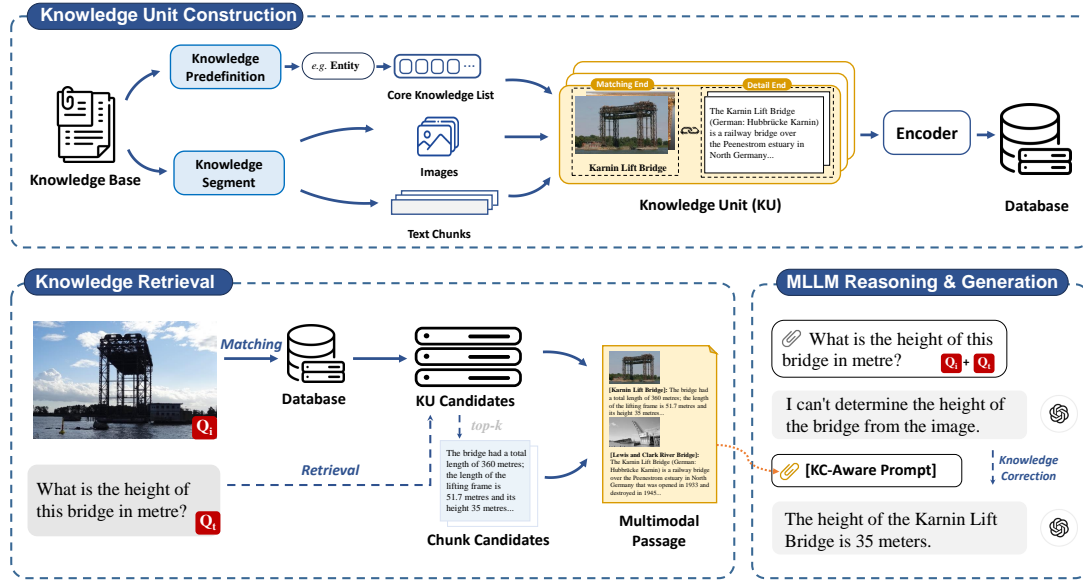


Figure 2: The KU-RAG Framework

retrieval (Chen et al., 2021). Additionally, during reasoning, this can be limited by the LLM’s maximum token capacity, leading to incomplete information. Therefore, we first need to segment the raw data, breaking it down into finer-grained units.

- **Textual data:** Similar to general RAG methods, we next segment all text passage $P = (p_1, p_2, \dots, p_n)$ in a knowledge base to obtain the smallest retrieval units. Each passage P contains n sentences, i.e., $P_k = (s_1, s_2, \dots, s_n)$. Considering the importance of knowledge coherence in the KB-VQA task, we adopt a combination of sentence splitting and maximum token limit. In each chunk, as many sentences as possible are retained without exceeding the maximum token limit, and the remaining sentences are assigned to the next chunk. That is, $C = (c_1, c_2, \dots, c_i)$, where $c_j = (s_1, s_2, \dots, s_j)$. i represents the i -th chunk, and j represents the j -th sentence in the chunk.

- **Image data:** For images pertaining to the same piece of information (such as all images within a news article), we directly extract them, treating each image individually.

Knowledge Assembly. After segmenting the text and generating chunks, the next step is to assemble these units with the unprocessed image information to form a knowledge unit. In simple terms, we assemble this multimodal knowledge into knowledge

units by leveraging the original structural properties of the knowledge and performing an inverted index on the text. To facilitate understanding, we illustrate this entity-type and event-type knowledge unit shown in Figure 3.

3.2.3 Storage

Next, we need to store the knowledge contained within the knowledge units. We encode each chunk into a vector using a text encoder and store them in a Faiss database (Douce et al., 2024), which we denote as D_t . Considering the need to handle multimodal data in the framework and the possibility of longer text within the chunks, we use Long-CLIP (Zhang et al., 2024) as the vector encoder.

$$V_{c_i} = \text{Encoder}(c_i), D_c = (V_{c1}, V_{c2}, \dots, V_{cn}) \quad (2)$$

For the images in the knowledge base, we also encode each image using Long-CLIP to obtain visual features and store them in the Meta Faiss vector database, denoted as D_i .

$$V_{i_j} = \text{Encoder}(i_j), D_i = (V_{i1}, V_{i2}, \dots, V_{in}) \quad (3)$$

3.3 Knowledge Unit Retrieval Framework

In this section, we will introduce our knowledge unit retrieval framework and detail how to achieve knowledge retrieval through knowledge unit and apply it to the KB-VQA task.

As shown in Figure 2, our framework is divided into three modules: **Knowledge Unit Construc-**

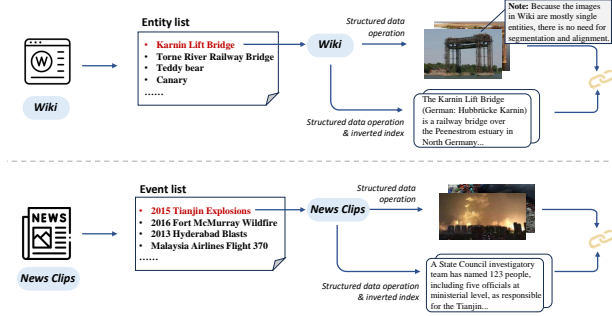


Figure 3: The construction of knowledge units (KU) with entity and event. The entity list comes from the work of Hu et al. (2023a), while the events are sourced from E-VQA (Yang et al., 2023).

tion, Knowledge Retrieval, and MLLM Reasoning & Generation. The knowledge unit construction module mainly transforms raw knowledge into knowledge units and stores them in the database, as illustrate in Section 3.2. The knowledge retrieval module processes the original query, matches it with the corresponding knowledge units, finds the relevant knowledge, and integrates it into a structured, MLLM-readable passage. Finally, combining the original question and the retrieved knowledge, the MLLM Reasoning & Generation module analyzes and generates the answer.

3.3.1 Query Processing

For the user’s input query Q , it is first necessary to preprocess and rewrite it to reduce interference during the retrieval process. To find the region in the image related to the question, we propose a query-aware instance segmentation method. Specifically, we first use YOLO (Redmon et al., 2016) to perform instance segmentation on the image, obtaining n segmented instance objects $O = (o_1, o_2, \dots, o_n)$. We then encode these instances using Long-CLIP, resulting in corresponding vectors $V_o = (v_{o_1}, v_{o_2}, \dots, v_{o_n})$:

$$V_{o_i} = \text{Encoder}(o_i) \quad (4)$$

Simultaneously, we encode the textual query Q_t into a vector and compute the similarity between it and each vector in V_o to find the object related to the query:

$$V_{q_t} = \text{Encoder}(Q_t), S_n = \text{Sim}(V_{q_t}, V_{o_i}) \quad (5)$$

Here, $S_n = (s_1, s_2, \dots, s_n)$ represents the similarity values corresponding to each O . We select

the object o with the highest similarity that exceeds a threshold γ for subsequent retrieval, with its vector denoted as V_{q_i} . Of course, sometimes the query may not only be related to one object but also to areas outside the objects or the entire scene in the image. Therefore, if no object meets the criteria or multiple objects meet the criteria, we will encode the entire image and use this encoding for subsequent retrieval:

$$V_{q_i} = \text{Encoder}(Q_i) \quad (6)$$

3.3.2 Knowledge Unit Matching

Next, we use the obtained visual features to match the corresponding knowledge unit. We select the top k knowledge unit items with the highest similarity, denoted as the set KU' , where each ku' contains j indices in its detail end.

$$KU = \text{Matching}(V_{q_i}), KU' = \text{top-k}(KU) \quad (7)$$

With $KU' = (ku'_1, ku'_2, \dots, ku'_k)$. The indices of each ku'_i are represented as $C_i = (c_{i,1}, c_{i,2}, \dots, c_{i,j})$. Finally, we obtain the combined index set of the knowledge unit:

$$C_{ku} = (C_1, C_2, \dots, C_n) \quad (8)$$

To integrate KU information into the query while highlighting the importance of certain content words, we rewrite Q_t as: " $Q'_t = Q_t [\text{SEP}] [KU \text{ name}] [\text{SEP}] \text{keywords}$ ", and encode it as:

$$V'_{q_t} = \text{Encoder}(Q'_t) \quad (9)$$

Here, “KU name” refers to the name of the matching segment of the retrieved KU, and “keywords” is a list of content words extracted from Q_t , separated by commas. “[SEP]” is a special token used to separate different parts. Next, we combine the features of Q_t and calculate the similarity to obtain the top k chunks related to the query, denoted as C' :

$$C' = \text{top-k}(\text{Sim}(V_{q_t}, C_{ku})) \quad (10)$$

with $C' = (c'_1, c'_2, \dots, c'_k)$.

3.3.3 MLLM Reasoning and Generation

After retrieving the relevant blocks, the next step is to provide the retrieved information to the MLLM to assist with reasoning and generation. The specific steps are as follows:

Modal Aligning and Fusing: First, based on the retrieval results C' , we find the corresponding knowledge unit KU' for each chunk and combine its matching end information to form an image with the structure $'[Image][[Name][Chunk Text]]'$, where the image corresponding to the i -th chunk is denoted as I'_i . Notably, if multiple chunks correspond to the same image, to enhance the connection between knowledge and improve processing efficiency, we merge the texts of these chunks into a single image in the format $'[Image][[Name][Chunk Text_1] \dots [Chunk Text_n]]'$.

Images Stitching: Next, we stitch all the images obtained in the previous step to generate a multimodal passage with both image and text information. This multimodal passage MP is as:

$$MP = (I'_1, I'_2, \dots, I'_n) \quad (11)$$

Knowledge Correction Chain: At this stage, a key challenge is effectively managing the relationships among the ‘information in the query’, ‘the knowledge retrieved’, and ‘the inherent knowledge of the MLLM’, as well as ensuring a fine-grained correspondence between text and images.

In our experiments, we found that when combining the retrieved knowledge with the question for the MLLM to answer, the model tended to prioritize the retrieved information while neglecting its own knowledge. We also attempted to use guiding prompts, such as “*Based on your own knowledge first...*” and “*Focus on the first image and ignore other image...*” to encourage the MLLM to consider its own knowledge before referring to the retrieved information, but the results were unsatisfactory (as demonstrated in Section 4.5).

To address this issue, we design a **Knowledge Correction Chain (KCC)** that guides MLLMs in reasoning through multi-turn dialogue and reading comprehension. In details, we first input the question Q to MLLM to obtain the original answer A_0 :

$$A_0 = \text{MLLM}(Q) \quad (12)$$

The purpose of this step is to obtain the pure knowledge of the MLLM regarding the query, without being influenced by the retrieved information mentioned above. Finally, we input the passage MP into the MLLM with a knowledge correction aware (KC-aware) prompt and get the final answer A :

[KC-aware prompt]: The initial answer has already been provided. The new image information may either be related or unrelated to the previous input. If this new information conflicts with the initial answer, please update the response accordingly. If no changes are needed, simply output the initial answer again.

$$A = \text{MLLM}(MP, \text{Prompt}, (Q, A_0)) \quad (13)$$

In short, the idea of KCC is to shift the MLLM’s focus from analyzing the relationship between “information in the query”, “the inherent knowledge of the MLLM”, and “the knowledge retrieved” to allow “the knowledge retrieved” to correct the MLLM’s responses, fostering a reflective process. We have also attempted to use a single prompt to have the MLLM generate and then reflect on its answer, but it still gets influenced by the retrieved information.

In this way, we can fully utilize multimodal information and handle the fine-grained correspondences between them, enhancing the MLLMs’ ability to reason and answer questions in complex scenarios. In addition, we have designed a method for managing KU. More details can be referred to in Appendix A.

4 Experiment

4.1 Dataset

To validate the effectiveness of our method, we selected four representative KB-VQA datasets, each with its own focus areas:

- **OVEN** (Hu et al., 2023a): An Open-domain Visual Entity Recognition dataset, primarily examining the ability to recognize the names of visual entities.
- **INFO SEEK** (Chen et al., 2023): An extension of the OVEN dataset, focusing on the coarse-grained knowledge behind entities, environments, etc., in images. It requires identifying the image and then discovering the knowledge behind it.
- **OK-VQA** (Marino et al., 2019): A classic KB-VQA dataset focusing on open-domain knowledge, featuring images paired with open-ended questions.
- **E-VQA** (Yang et al., 2023): An event-centric dataset, primarily evaluating the ability to recognize events and the knowledge behind them.

Table 1: Main results of the experiment. And [†] indicates that the result is from experiments conducted on the full version of the test set, sourced respectively from Hu et al. (2023a) and Chen et al. (2023). Except for the SOTAs, which are trained, other methods are performed in a **zero-shot** scenario.

Model	Dataset			
	OVEN _s	INFO SEEK _s	OK-VQA	E-VQA
SOTA (Trained)	21.70 [†]	22.10 [†]	66.10	19.42
LlaVa NEXT-7b	9.51	6.37	73.33	10.51
LlaVa NEXT-7b + KU-RAG	10.80	9.09	73.07	11.00
Llama 3.2-11b	18.50	19.07	68.53	15.04
Llama 3.2-11b + KU-RAG	18.83	19.04	68.89	15.22
GPT-4o	22.30	36.05	75.52	15.17
GPT-4o + KU-RAG	26.50	38.35	77.23	26.16

Table 2 shows some characteristics of each dataset. The more stars in question granularity, the finer the question. The higher the popularity of knowledge, the more general it is, meaning the MLLM is more likely to have learned it during pre-training. Note that since our method is conducted in a zero-shot setting, we only selected the test sets of these datasets. Due to the large size of the original test sets for OVEN and INFO SEEK, we sampled some examples using an arithmetic sequence for testing, and the term ‘s’ is used in Table 2 and subsequent experimental results to indicate this.

Table 2: Characteristics of different dataset.

Dataset	Tests Number	Knowledge Source	Knowledge Granularity	Knowledge Popularity
OVEN _s	23,650	Wiki	**	**
INFO SEEK _s	11,600	Wiki	***	**
OK-VQA	5,064	Wiki	**	***
E-VQA	1,819	News	***	*

4.2 Baseline

For the selection of baselines, we chose representative multimodal large language models (MLLMs) with different parameter sizes, including the closed-source model **GPT-4o**, as well as the open-source models **Llava NEXT-7b** (Liu et al., 2024) and **Llama 3.2-11b**. Due to variations in the formats and objectives of each dataset, there is no single unified state-of-the-art (SOTA) model across all of them. To ensure a fair comparison, we select the best-performing model for each dataset as its respective SOTA baseline (Chen et al., 2022, 2023; Hu et al., 2023a). For detailed information about the baselines, please refer to Appendix B.1.

4.3 Experimental Setting

Our experiments were conducted on RTX 4090 GPUs. GPT-4o uses the base version of the API

interface, while LlaVa and Llama conduct experiments using the Hugging Face transformers library¹. In our methods’ settings, OVEN, INFO SEEK, and OK-VQA all use entities as the knowledge unit, while E-VQA uses events as the knowledge unit. For the recall of knowledge units and chunks, the top-k is set to 3. For the experiment evaluation, we used accuracy as the metric.

4.4 Main Result

As shown in Table 1, we have the following findings.

Zero-shot Capability of MLLMs: MLLM demonstrates remarkable zero-shot capability, especially in image understanding and reasoning. Compared to the previous SOTA model for KB-VQA, GPT-4o shows improvements of 0.6%, 13.95%, and 9.42% on the OVEN, INFO SEEK, and OK-VQA datasets, respectively. This is mainly due to the extensive world knowledge accumulated during MLLM’s pre-training phase. However, since the E-VQA dataset involves less popular news knowledge, MLLM’s performance in this area is not as strong as that of the specially trained SOTA model.

Superior Performance of MLLM+KU-RAG: Our method, MLLM+KU-RAG, performs excellently across all datasets. In a zero-shot scenario, even without reviewing the training set knowledge, GPT-4o+KU-RAG outperforms the existing SOTA models by 4.8%, 16.25%, 11.13%, and 6.74% on the four datasets, respectively, validating the powerful performance of our approach.

Enhancement of MLLM by KU-RAG: Combining KU-RAG with GPT-4o results in performance improvements of 4.2%, 2.3%, 1.7%, and 10.99% on the respective tasks. The largest improvement is seen in the E-VQA dataset, as it in-

¹<https://huggingface.co/docs/transformers/main>

volves less popular knowledge, and the new knowledge provided by KU-RAG significantly enhances model performance. In contrast, the improvement on the OK-VQA dataset is smaller because it involves open-domain general knowledge, which MLLM may have already encountered during pre-training, allowing it to answer effectively.

However, for smaller MLLMs like Llava NEXT-7b and Llava 3.2-11b, the improvement from RAG is not very significant, with increases of less than 3% across datasets. This may be due to the limited parameter size, which makes them less capable than GPT-4o in handling the relationship between original and retrieved information.

4.5 Ablation Result

To validate the effectiveness of each component in our proposed method, we designed ablation experiments comparing the following models:

- **w/o KCC:** This model omits the knowledge correction chain (KCC), relying instead on the model’s analysis of the question and the retrieved information in a single-turn Q&A setup.
- **w/o KU:** This model removes the fine-grained retrieval approach (*i.e.*, knowledge unit), converting the information from images into captions and using a text-only retrieval modality.

Additionally, we included the full implementation of the GPT-4o+KU-RAG method, as well as a standalone GPT-4o. The experimental results are shown in Figure 4. From the figure, we can draw the following conclusions:

Effectiveness of GPT-4o+KU-RAG: The GPT-4o+KU-RAG method consistently achieves the highest performance, which demonstrates that every part of this method is indispensable.

Impact of Removing KCC: Removing KCC and using single-turn dialogue markedly reduces model performance across four datasets, with decreases of 6%, 18.79%, 8.3%, and 7.97%, respectively. Except for the E-VQA dataset, the model’s performance is inferior to using only GPT-4o. This likely occurs because the model struggles to effectively focus on the original question’s image and manage the logical relationships between the query information, its own knowledge, and the retrieved knowledge. Consequently, some questions, which the model could originally answer correctly, are

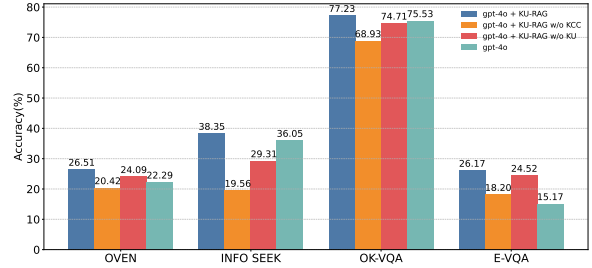


Figure 4: The Results of Ablation Study

answered incorrectly due to interference from the injected information.

Impact of Removing KU: Removing KU and adopting a coarse-grained, single-modality retrieval approach results in a slight performance drop across datasets, with the most significant decrease observed in the INFO SEEK dataset (9.04%). This is partly because INFO SEEK requires matching detailed image content and background knowledge, and converting the original image to captions loses a substantial amount of visual information. As illustrated by examples in Figure 1, it’s challenging to accurately match “Karnin Lift Bridge” using just the text “bridge,” let alone find corresponding background knowledge. Furthermore, introducing incorrect knowledge adds noise, impeding the MLLM’s reasoning process and leading to erroneous results. The smallest performance drop is observed in the E-VQA dataset (1.65%), likely because, in this dataset, the images primarily serve to supplement information, allowing text-only retrieval to still achieve reasonably good matches.

5 Conclusion

In this paper, we introduce the Knowledge Unit Retrieval-Augmented Generation (KU-RAG) method, aimed at enhancing MLLMs by incorporating fine-grained retrieval of domain-specific knowledge. To improve the effectiveness of retrieval, we propose the concept of “knowledge units”, which allows for more targeted access to relevant information. Furthermore, we design a knowledge correction chain strategy to verify and refine the retrieved knowledge, which can mitigate errors and hallucinations, enhancing the overall reliability and coherence of the generated answers in VQA tasks. Our experimental results demonstrate significant performance gains across multiple KB-VQA benchmarks, highlighting the effectiveness of our approach. Future research directions could explore dynamic knowledge updates strategy to improve multimodal retrieval and semantic integration.

Limitations

Our method exhibits exceptional performance enhancement on large-scale MLLMs such as GPT-4o. However, as evident from the main experimental results, the degree of improvement is comparatively less pronounced when applied to smaller parameter MLLMs. This observation highlights a potential limitation of our approach in effectively scaling down its benefits for smaller models. To address this, future research could focus on refining and optimizing KU-RAG to better accommodate smaller language models, aiming to unlock significant performance gains across a broader range of model sizes.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. 2022. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*.
- Xilun Chen, Kushal Lakhota, Barlas Oğuz, Anchit Gupta, Patrick Lewis, Stan Peshterliev, Yashar Mehdad, Sonal Gupta, and Wen-tau Yih. 2021. Salient phrase aware dense retrieval: can a dense retriever imitate a sparse one? *arXiv preprint arXiv:2110.06918*.
- Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, Soravit Changpinyo, Alan Ritter, and Ming-Wei Chang. 2023. Can pre-trained vision and language models answer visual information-seeking questions? *arXiv preprint arXiv:2302.11713*.
- Zhenfang Chen, Qinhong Zhou, Yikang Shen, Yining Hong, Zhiqing Sun, Dan Gutfreund, and Chuang Gan. 2024. Visual chain-of-thought prompting for knowledge-based visual reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 1254–1262.
- Yang Ding, Jing Yu, Bang Liu, Yue Hu, Mingxin Cui, and Qi Wu. 2022. Mukea: Multimodal knowledge extraction and accumulation for knowledge-based visual question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5089–5098.
- Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library. *arXiv preprint arXiv:2401.08281*.
- Wenqi Fan, Yajuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6491–6501.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.
- Liangke Gui, Borui Wang, Qiuyuan Huang, Alex Hauptmann, Yonatan Bisk, and Jianfeng Gao. 2021. Kat: A knowledge augmented transformer for vision-and-language. *arXiv preprint arXiv:2112.08614*.
- Hexiang Hu, Yi Luan, Yang Chen, Urvashi Khandelwal, Mandar Joshi, Kenton Lee, Kristina Toutanova, and Ming-Wei Chang. 2023a. Open-domain visual entity recognition: Towards recognizing millions of wikipedia entities. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12065–12075.
- Yuntong Hu, Zhihan Lei, Zheng Zhang, Bo Pan, Chen Ling, and Liang Zhao. 2024. Grag: Graph retrieval-augmented generation. *arXiv preprint arXiv:2405.16506*.
- Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A Smith, and Jiebo Luo. 2023b. Promptcap: Prompt-guided image captioning for vqa with gpt-3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2963–2975.
- Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.
- Yuanze Lin, Yujia Xie, Dongdong Chen, Yichong Xu, Chenguang Zhu, and Lu Yuan. 2022. Revive: Regional visual representation matters in knowledge-based visual question answering. *Advances in Neural Information Processing Systems*, 35:10560–10571.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. Llava-next: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. *Visual instruction tuning*. *Preprint, arXiv:2304.08485*.

Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2023b. Gpt understands, too. *AI Open*.

Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204.

Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788.

Ander Salaberria, Gorka Azkune, Oier Lopez de Lacalle, Aitor Soroa, and Eneko Agirre. 2023. Image captioning for effective use of language models in knowledge-based visual question answering. *Expert Systems with Applications*, 212:118669.

Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. 2023. Prompting large language models with answer heuristics for knowledge-based visual question answering. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 14974–14983.

Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo Yuan, Quanzeng You, and Hongxia Yang. 2024. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning. *arXiv preprint arXiv:2401.06805*.

Haoyang Wen, Honglei Zhuang, Hamed Zamani, Alexander Hauptmann, and Michael Bendersky. 2024. Multimodal reranking for knowledge-intensive visual question answering. *arXiv preprint arXiv:2407.12277*.

Jialin Wu, Jiasen Lu, Ashish Sabharwal, and Roozbeh Mottaghi. 2022. Multi-modal answer validation for knowledge-based vqa. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 2712–2721.

Jialin Wu and Raymond J Mooney. 2022. Entity-focused dense passage retrieval for outside-knowledge visual question answering. *arXiv preprint arXiv:2210.10176*.

Qiaofeng Wu, Wenlong Fang, Weiyu Zhong, Fenghuan Li, Yun Xue, and Bo Chen. 2025. Dual-level adaptive incongruity-enhanced model for multimodal sarcasm detection. *Neurocomputing*, 612:128689.

Zhenguo Yang, Jiale Xiang, Jiuxiang You, Qing Li, and Wenyan Liu. 2023. Event-oriented visual question answering: The e-vqa dataset and benchmark. *IEEE Transactions on Knowledge and Data Engineering*, 35(10):10210–10223.

Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. 2024. Long-clip: Unlocking the long-text capability of clip. *arXiv preprint arXiv:2403.15378*.

A Database Management

Since our approach operates at the database level, compared to traditional methods, we must also consider issues related to data management. Additionally, there is no need to retrain the entire framework after adding, deleting, or updating data or knowledge.

Knowledge Addition: When new knowledge is introduced, we directly chunk the corresponding text. The encoded text vectors are then stored in the existing text vector database D_c . Similarly, for images within the knowledge base, we encode them and store the resulting vectors in the image vector database D_i .

Knowledge Unit Management: After performing operations of the raw data, it is necessary to consider whether corresponding actions need to be taken for the knowledge unit.

1) KU Addition and Update: When new raw data is added, it is essential to assess whether a new KU needs to be introduced. This process primarily involves two steps: first, matching the new knowledge with existing KU. If the similarity of the matching results exceeds the threshold α , the index containing the new keywords is added to the corresponding KU through keyword matching. If no matching result exceeds the threshold, a new KU is constructed according to the KU construction rules and the new keywords present in the chunk.

2) KU Deletion: After deleting a chunk from the raw data, it is necessary to check whether the related KU is still valid to reduce storage usage. Specifically, after deleting the chunk indexed as i , all KUs containing this index should be checked. If a certain KU has an empty detail end (*i.e.*, no remaining values), that KU can be deleted.

B Experimental Supplement

B.1 Baseline

For the selection of baselines, we chose representative MLLMs with different parameter sizes. Due to variations in the formats and objectives of each dataset, there is no single unified state-of-the-art (SOTA) model across all of them. To ensure a fair comparison, we select the best-performing model for each dataset as its respective SOTA baseline.

Model	Dataset			
	OVEN _s	INFO SEEK _s	OK-VQA	E-VQA
Llama	18.50	19.07	68.53	15.04
Llama + KU-RAG w/o KCC	18.56	18.19	68.25	23.15
Llama + KU-RAG	18.83	19.04	68.89	15.22
LlaVa	9.51	6.37	73.33	10.51
LlaVa + KU-RAG w/o KCC	7.01	5.75	66.35	12.79
LlaVa + KU-RAG	10.08	9.09	73.07	11.00
GPT-4o	22.30	36.05	75.52	15.17
GPT-4o + KU-RAG w/o KCC	20.42	19.56	68.93	18.20
GPT-4o + KU-RAG	26.50	38.35	77.23	26.16

Table 3: Further Analysis on KCC.

- **SOTA:** For **OVEN** dataset, we use PaLI-17B (Chen et al., 2022), as reported by Hu et al. (2023a). For **INFO SEEK** and **OK-VQA** datasets, the SOTA model is PaLI-X (Chen et al., 2022), as reported in the work of Chen et al. (2023). For **E-VQA** dataset, we adopt the best results of the SOTA model MuKEA (Ding et al., 2022), as reported Yang et al. (2023).

For the MLLMs:

- **GPT-4o:** A closed-source MLLM launched by OpenAI, with powerful multimodal content understanding and reasoning capabilities.
- **LlaVa NEXT-7b** (Liu et al., 2024): The 7b parameter version of the latest LlaVa model, an open-source MLLM.
- **Llama 3.2-11b:** The 11b parameter version of the latest Llama model proposed by Meta AI, also an open-source MLLM.

B.2 Setup and Environment

The experiments were conducted in a zero-shot setting using 4 RTX 4090 GPUs, with PyTorch version 2.3.0. For GPT-4o, we used the interface of the GPT-4o base model. For LLaVa NEXT-7b and Llama 3.2-11b models, we used the Hugging Face Transformers library, version 4.46.3. The LLaVa NEXT-7b model used the weight file ‘llava-v1.6-mistral-7b-hf’, while the Llama 3.2-11b model loaded the weight file ‘Llama-3.2-11B-Vision-Instruct’.

B.3 Further Analysis on KCC

We conducted a further analysis on the impact of removing the Knowledge Correction Chain (KCC), as shown in Table 3. Based on the results, for the

OVEN, INFO SEEK, and OK-VQA datasets, the performance of the models, whether it’s the smaller models like Llama and LlaVa, or the powerful large language model GPT-4o, all decreased to varying degrees after removing KCC. This suggests that KCC plays a significant role in enhancing model performance, especially for GPT-4o, where the performance drop after removing KCC is the most substantial, indicating that the gain from KCC is strongest for this model.

Additionally, after removing KCC, the models are influenced by the external knowledge retrieved, which leads them to disregard their own internal knowledge. For the OVEN, INFO SEEK, and OK-VQA datasets, the performance without KCC did not surpass the performance of the models used individually. This indicates that KCC helps balance the relationship between the retrieved knowledge and the model’s own knowledge, thereby improving overall performance.

It’s worth noting that for the E-VQA dataset, Llama and LlaVa experienced a performance drop when KCC was introduced. The main reason for this is that the knowledge in the E-VQA dataset is more specific and less general, so the answers to the questions mainly rely on the knowledge retrieved. However, due to the smaller model sizes of Llama and LlaVa, they struggle to effectively manage the conflict between their internal knowledge and the retrieved knowledge, leading to a decline in performance. In contrast, GPT-4o, with its larger model capacity, is able to handle these conflicts more effectively, resulting in better performance on the E-VQA dataset.