

MAR: Matching-Augmented Reasoning for Enhancing Visual-based Entity Question Answering

Anonymous ACL submission

Abstract

A multimodal large language model (MLLM) may struggle with answering *visual-based (personal) entity questions (VEQA)*, such as “who is *A*?” or “who is *A* that *B* is talking to?” for various reasons, *e.g.*, the absence of the name of *A* in the caption or the inability of MLLMs to recognize *A*, particularly for less common entities. Furthermore, even if the MLLM can identify *A*, it may refrain from answering due to privacy concerns. In this paper, we introduce a novel methodology called Matching-Augmented Reasoning (MAR) to enhance VEQA. Given a collection of visual objects with captions, MAR preprocesses each object individually, identifying faces, names, and their alignments within the object. It encodes this information and stores their vector representations in vector databases. When handling VEQA, MAR retrieves matching faces and names and organizes these entities into a matching graph, where nodes represent entities and edges indicate their similarities. MAR then derives the answer to the query by reasoning over this matching graph. Extensive experiments show that MAR significantly improves VEQA compared with the state-of-the-art methods using MLLMs.

1 Introduction

Multimodal language models (MLLMs) (Cui et al., 2024) like GPT-4V (Zhang et al., 2023) and



Figure 1: **Data** (V : image, T : text) pair; **Query** (R : entity selection, Q : question) pair. (a) The advantages of MLLMs; (b) The limitations of MLLMs, and (c) Our proposal MAR.

LLaVA (Liu et al., 2023) have significantly improved visual question answering (VQA) by integrating text and images. However, they still face challenges in visual-based entity question answering (VEQA), a crucial subset of VQA that focuses on extracting information about specific entities, especially for personal entities.

MLLMs for VEQA: Advantages and Limitations.

040	In VEQA tasks, MLLMs excel at integrating visual	080
041	cues and textual information for effective reason-	081
042	ing and answer generation (Li et al., 2023b;	082
043	Liu et al., 2024). For instance, as depicted in	083
044	Figure 1(a), GPT-4V, when tasked with answer-	084
045	ing question Q_1 regarding the face in region R_1 ,	085
046	leverages the associated caption T_1 of image V_1	
047	to precisely identify the person within the red	
048	box as “Wang Yi”.	
049	However, MLLMs often struggle to recognize	
050	all details in images, particularly for less com-	
051	mon entities (Li et al., 2023b). For instance, in	
052	Figure 1(b), GPT-4V fails to answer question	
053	Q_2 about the person in the red rectangle R_2 due	
054	to the lack of information in the image caption	
055	T_2 and its limited knowledge base. Furthermore,	
056	even when an MLLM identifies an entity, it may	
057	withhold an answer due to privacy regulations.	
058	Despite rapid advancements of MLLMs , accu-	
059	rately identifying all (personal) entities in im-	
060	ages and adhering to privacy regulations make	
061	answering VEQA questions solely using MLLMs a	
062	significant challenge (Chen et al., 2024; Li et al.,	
063	2023a, 2024; Yu et al., 2023).	
064	Matching-Augmented Reasoning (MAR) . Given	
065	a collection of visual objects with captions,	
066	sourced from public or enterprise datasets with-	
067	out privacy concerns, MAR identifies the faces of	
068	entities within visual objects and the names of	
069	entities within captions by tools like CLIP (Rad-	
070	ford et al., 2021) and Deepface (Taigman et al.,	
071	2014). These entities are encoded with respec-	
072	tive visual and text encoders, and the resulting	
073	embeddings are stored in vector databases <i>e.g.</i> ,	
074	Meta Faiss (Douze et al., 2024). When a VEQA	
075	query is posed, MAR retrieves “similar” faces and	
076	names from the database and performs reason-	
077	ing over these matched pieces of information	
078	to generate an accurate response. Note, in this	
079	study, our focus is on personal entities. We plan	
	to extend our analysis to include additional types	080
	of entities in future research.	081
	As illustrated in Figure 1(c), if we can suc-	082
	cessfully match the face in image V_2 with the	083
	face in image V_1 , and if we know that the face	084
	in V_1 is “Yi Wang”, we can easily answer Q_2 .	085
	Contributions. Our notable contributions are	086
	summarized as follows.	087
	• We study VEQA , an important and com-	088
	monly used subset of VQA, but is under-	089
	explored. (Section 3)	090
	• We propose <i>matching graphs</i> that can cap-	091
	ture the relationships of the same enti-	092
	ties over multiple captioned visual objects.	093
	Based on a matching graph, we proposed	094
	matching-augmenting reasoning (MAR), to	095
	effective answer a VEQA . (Section 4)	096
	• Given that VEQA is a relatively new prob-	097
	lem, existing benchmarks are not suit-	098
	able. Therefore, we have constructed a new	099
	benchmark NewsPersonQA including 235k	100
	images and 6k QA pairs. (Section 5)	101
	• We conduct extensive experiments to show	102
	that MAR > MLLMs + RAG > MLLMs , where	103
	RAG is to feed the retrieved matching graph	104
	to MLLMs . (Section 6)	105
	2 Related Work	106
	VQA. VQA aims at reasoning over visual and	107
	textual content and cues to generate answers (Lu	108
	et al., 2021; Stengel-Eskin et al., 2022; Agrawal	109
	et al., 2023). It primarily utilizes approaches	110
	such as Fusion-based (Zhang et al., 2019), Multi-	111
	modal Learning (Ilievski and Feng, 2017), Mem-	112
	ory Networks (Su et al., 2018), Visual Atten-	113
	tion (Mahesh et al., 2023), etc., to discover and	114
	integrate information from text and images.	115
	MLLMs for VQA. MLLMs , such as GPT-	116
	4V (Zhang et al., 2023) and LLaVa (Liu et al.,	117

2023), have played a pivotal role in advancing VQA. By seamlessly integrating textual and visual information, these models have demonstrated a remarkable ability to understand and respond to complex queries about images.

RAG for VQA. However, in many cases, the cues within images and text are insufficient for reasoning and answering. Retrieval-augmented generation (RAG) (Lewis et al., 2021) has been studied for VQA, especially with Knowledge-Based VQA approaches that incorporate external knowledge to provide additional cues for answers (Khademi et al., 2023; Lin et al., 2022).

VEQA. In this paper, we investigate **VEQA**, a critical subset of VQA that concentrates on querying information about entities, especially persons. As will be shown in Section 6, **MLLMs** often struggle with such questions due to limited knowledge and privacy considerations. While RAG can enhance **MLLMs** for **VEQA** tasks, **MLLMs** still face challenges (or confused) in reasoning with multiple interconnected visual objects.

Data Matching. Data matching refers to the process of identifying, comparing, and merging records from multiple datasets to determine whether they correspond to the same entities (Christen and Christen, 2012). With the increasing multimodality of data, the concept of matching has been continually expanded from its original string matching (Text-Text) and entity matching (Tuple-Tuple) context. For instance, Image-Text Matching (Lee et al., 2018; Li et al., 2019), Image-Image matching (Zhu et al., 2018), etc. In fact, matching can aggregate more clues, enhance the reasoning ability of models, and possess strong interpretability (Zheng et al., 2022).

3 Problem

Captioned Visual Objects. We consider a *captioned visual object* O as a pair $O : (V, T)$

where V is an image, and T is an optional text description relative to the image V .

Figure 1(a) and Figure 1(b) provide two sample captioned visual objects, (V_1, T_1) and (V_2, T_2) , respectively.

Let $\mathbf{O} = \{O_1, O_2, \dots, O_n\}$ be a group of captioned visual objects, sourced from public or enterprise datasets without privacy concerns. Note that, such a group is common in practice, e.g., a collection of news articles.

Users can pose a Visual-based (Personal) Entity Question Answering (**VEQA**) on either a single captioned visual object (**Single-VEQA**) or a group of such objects (**Group-VEQA**), as defined below.

Single-VEQA. Given a captioned visual object $O : (V, T)$, this type of queries allows the user to provide a rectangle selection of the image and ask the question like “who is he/she” or “is he/she John”.

More formally, a **Single-VEQA** Q_s is a pair (R, Q) , where R is a rectangle selection over image V and Q is a natural language question.

Group-VEQA. Given a group of captioned visual objects \mathbf{O} , we support two types of queries Q_g : (1) a simple natural language query Q , such as “how many news contain Donald Trump”; and (2) a natural language query with a selected face, i.e., a pair (R, Q) , such as “in which news the selected person appears”.

We will simply use Q to represent either a **Single-VEQA** or a **Group-VEQA** query, when it is clear from the context.

4 Algorithms for VEQA

In this section, we will first discuss solely using **MLLMs** for **VEQA** in Section 4.1. We will then discuss coarse-grained retrieval-augmented generation (RAG) in Section 4.2. We then propose a new concept, called matching graphs, which

can provide fine-grained information among retrieved objects in Section 4.3. Based on matching graphs, we describe fine-grained RAG in Section 4.4 and matching-augmented reasoning (MAR) in Section 4.5.

4.1 MLLMs for VEQA

Given a VEQA query Q , a crude solution is to directly prompt Q to a MLLM as:

$$Q \rightarrow \text{MLLM} \rightarrow \text{answer}$$

Figure 2(a) depicts this solution.

4.2 Coarse-Grained RAG for VEQA

Alternatively, we can retrieve top- k captioned visual objects and feed them to MLLMs as:

$$(Q, \text{top-}k \text{ objects}) \rightarrow \text{MLLM} \rightarrow \text{answer}$$

Figure 2(b) illustrates this approach, which we refer to as *coarse-grained RAG*. This method is characterized by its transmission of entire retrieved objects to the MLLMs. Unfortunately, current MLLMs perform poorly in reasoning with multiple interconnected retrieved visual objects.

4.3 Matching Graphs

To improve the performance of RAG models, it’s beneficial to focus on fine-grained information rather than entire objects. By identifying specific entities (*e.g.*, faces, names) and their connections within each object, we can provide a more meaningful context for reasoning.

Matching Graphs. A matching graph $G(N, E)$ contains a set N of nodes and a set E of undirected edges. Each node $n \in N$ has two labels $\text{face}(n)$ and $\text{name}(n)$, where $\text{face}(n)$ is a face image, and $\text{name}(n)$ is a set of possible names.

If we are certain about a person’s name, we will use a square bracket *e.g.*, $\text{name}(n) = [\text{Yi Wang}]$ for the selected face in Figure 1(a); if we are not sure about a person’s name, we will

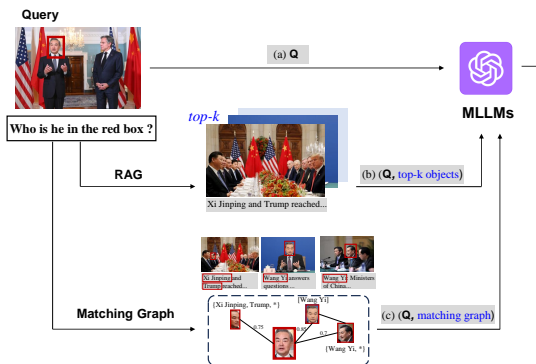


Figure 2: Different algorithms for VEQA. (a) MLLMs. (b) Coarse-grained RAG. (c) Fine-grained RAG.

use a curly bracket to indicate possible names *e.g.*, $\text{name}(n) = \{\text{Xi Jinping, Trump, *}\}$ for the selected face in Figure 1(b), where $*$ is a wildcard meaning that n ’s name could be something other than Xi Jinping and Trump.

Each undirected edge $e(n_i, n_j) \in E$ indicates that the two faces corresponding to n_i (*i.e.*, $\text{face}(n_i)$) and n_j (*i.e.*, $\text{face}(n_j)$) are likely to be the same person. Each edge has a weight $\text{weight}(e) \in [0, 1]$, indicating the *similarity* of the two faces.

Matching Graph Construction. It consists of two steps: offline index construction (for all data objects) and online matching graph construction (for each query).

Offline Index Construction. We first preprocess each captioned visual object $O(V, T)$ as follows.

- **Face identification.** We use Meta DeepFace (Taigman et al., 2014) to extract face entities as (f_1, f_2, \dots, f_k) from image V .
- **Name identification.** We use spaCy (Honnibal et al., 2020) to extract name entities as (x_1, x_2, \dots, x_m) from text T .

After pre-processing, we have constructed all possible nodes for all possible matching graphs. We then use pre-trained CLIP (Radford et al., 2021) to convert each identified face and each

259	identified person names into its vector representation, and store them in two separate vector database: faceDB and nameDB .		
260			
261			
262	<u>Iterative Online Matching Graph Construction.</u>		
263	Given a VEQA query, we construct a matching graph as follows.		
264			
265	[Step 1: Initialization.] The user starts with a <i>seed node</i> (for Single- VEQA) or a group of <i>seed nodes</i> for (Group- VEQA). Each seed node contains a face and its candidate names that could be empty.		
266			
267			
268			
269			
270			
271	[Step 2: Graph Expansion.] For each node in the graph, we search either similar faces from faceDB with vector similarity above a given threshold σ_f , or similar names from nameDB with vector similarity above a given threshold σ_n . For each added node, the edge weight is set as face similarity.		
272			
273			
274			
275			
276			
277			
278			
279	[Step 3: Iterative Search and Termination.] When there are new nodes added in Step 2, we will loop Step 2. The process terminates when either there is no new nodes can be added or we have done k iterations. From our empirical findings, we set $k = 2$, which is enough to retrieve useful nodes (<i>e.g.</i> , 10 nodes) and edges for reasoning.		
280			
281			
282			
283			
284			
285			
286			
287			
288	4.4 Fine-Grained RAG for VEQA		
289	Given the fine-graph matching graph relative to a query Q , we prompt it to MLLMs as:		
290			
291	(Q , matching graph) \rightarrow MLLM \rightarrow answer		
292	Figure 2(c) shows this approach, which we refer to as <i>fine-grained RAG</i> . It works as follows.		
293	[Step 1: Image Stitching.] Most MLLMs (<i>e.g.</i> , LLaVA) only support only single-image input, thus we simply combine multiple retrieved visual objects into one visual object V .		
294			
295	[Step 2: Image Annotation.] We annotate each node n_i in the matching graphs on V in a red		
296			
297			
298			
299			
300			
301			
	box, resulting in an annotated image V' .		302
	[Step 3: Matching Graph Serialization.] Each node n_i and edge $e(n_i, n_j)$ will be serialized as:		303
	$\mathbf{ser}(n_i) = \mathbf{face}(n_i), \mathbf{name}(n_i)$		305
	$\mathbf{ser}(e) = n_i, n_j, \mathbf{weight}(e)$		306
	Serializing a matching graph $g(N, E)$ is to serialize all nodes and edges as:		307
	$\mathbf{ser}(g) = \mathbf{ser}(N), \mathbf{ser}(E)$		308
	We then prompt Q , V' , and $\mathbf{ser}(g)$ to MLLMs . In order to enable it to consider information from its own model simultaneously, we also designed an Original knowledge-aware Prompt (OP): “Please tell me Q . If you are unsure, read the following.”		309
			310
			311
			312
			313
			314
			315
			316
	4.5 MAR for VEQA		317
	MAR for Single-VEQA. This type of queries asks the name of a single entity. Given a matching graph $g(N, E)$ where $n^* \in N$ is the seed node, our method works as follows.		318
	[Step 1: Remove Uncertain Nodes.] For each node $n_i \in N \setminus \{n^*\}$, if its name is uncertain, we remove n_i and its associated edges, which will result in a modified graph $g(N', E')$.		319
	[Step 2: Name Aggregation for n^* .] We count all distinct names in the modified matching graph g' , each associated with a weight as $\sum_{e(n_i, n^*) \in E'} \mathbf{weight}(e)$.		320
	[Step 3: Name Identification for n^* .] We pick the name with the highest weight, as the answer to the Single- VEQA query.		321
			322
			323
			324
			325
			326
			327
			328
			329
			330
			331
			332
			333
			334
			335
	MAR for Group-VEQA. This type of queries ask for aggregated information of nodes whose names are queried in the query, <i>e.g.</i> , “which image/how many images have person A”. Given a matching graph $g(N, E)$, it works as follows.		336
	[Step 1: Name Identification for Each Node.] It first identifies the name of each node, as discussed above.		337
	[Step 2: Answer Aggregation.] It aggregates		338
			339
			340
			341
			342
			343
			344
			345

the information of each node to answer the given Group-VEQA.

5 A New NewsPersonQA Benchmark

The problem of VEQA needs to address complex interactions between multiple visual and textual data. Despite its growing importance, existing benchmarks fall short in adequately representing the diverse challenges posed by VEQA tasks. Particularly in the domain of News QA, where the accurate identification and understanding of both common and uncommon persons are crucial, current datasets (*e.g.*, GoodNews (Biten et al., 2019) and NewsQA (Trischler et al., 2016)) do not provide the necessary depth and breadth. To bridge this gap, based on GoodNews (Biten et al., 2019), we are constructing a new benchmark, namely NewsPersonQA, that encompasses a wide range of scenarios, including both well-known and obscure individuals.

The construction of the dataset entails the generation of QA pairs from the raw data in GoodNews, which consists of images and captions. This process involves two main steps: data preprocessing and QA pair construction.

Data Preprocessing: Raw data undergoes preprocessing, which includes structuring news data, extracting faces from images, annotating original images, and recognizing named entities in captions. The processed data is then randomly distributed into groups. Each group contains thousands of images and is categorized into Single-VEQA (100 groups) and Group-VEQA (10 groups) queries.

Single-VEQA Question Generation: We begin by counting the frequency of each person’s name within each group. To ensure the availability of clues for answering, we select names that appear at least three times in captions. We then mask these names in the captions to gener-

Category	Count
Total Images	235,912
Totally Extracted Faces	336,075
Totally Extracted Names	379,313
Single-VEQA Queries	4,937
Group-VEQA Queries	1,004
Total Queries	5,941

Table 1: Statistics of NewsPersonQA

ate QA pairs. For example: **Question:** “Who is the person labeled ‘face n ’ in the red box?” **Answer:** “name”. In total, approximately 5,000 queries of this type are generated, about 50 per group.

Group-VEQA Question Generation: Similarly, we count the occurrences of names within each group and store the image names as a set, denoted as S . To prevent exceeding the maximum token limit of MLLMs in the answers and to facilitate clearer visualization of experimental results, we limit each person’s name to a maximum of 5 appearances within the same group. We then randomly mask part of the captions corresponding to the images in the set to increase the difficulty and encourage MLLMs to generate correct answers through retrieved content. The format of QA pairs is **Question:** “Which photos are of the person named ‘name’?” **Answer:** S . The number of queries of this type is approximately 1,000.

Table 1 shows the statistics of NewsPersonQA.

6 Experiment

Methods. For answering VEQA queries, we selected two well-known and highly capable MLLMs to serve as baselines.

- **LLaVA:** This model utilizes CLIP-ViT-L-336px with an MLP projection. We refer to the 1.5 version with 7 billion parameters as LLaVA-7b and the version with 13 billion parameters as LLaVA-13b.

417 - **GPT-4V**: Recognized as OpenAI’s most
 418 powerful general-purpose MLLM to date, GPT-
 419 4V boasts 1.37 trillion parameters.

420 - **Human**: This represents the human-
 421 annotated results, showcasing the level of cog-
 422 nitive ability and performance that humans can
 423 achieve on this task.

424 + **FRAG**: MLLMs struggle with reasoning
 425 over coarse-grained RAG that consists of mul-
 426 tiple captioned visual objects. Therefore, we
 427 provide only fine-grained RAG (FRAG), *i.e.*,
 428 matching graph, to the above-mentioned models
 429 and human evaluators.

430 **Implementation.** The experiments were con-
 431 ducted in a zero-shot setting using RTX 4090
 432 GPUs. For GPT-4V, we used the interface of the
 433 GPT-4-vision-preview model. It’s worth noting
 434 that GPT-4V often refrains from answering per-
 435 son identify questions without additional clues
 436 due to policy reasons. However, with the incor-
 437 poration of matching graph techniques, it can
 438 leverage weak signals and combine them with its
 439 own knowledge base. In the case of Group-**VEQA**
 440 queries, a maximum of 10 cases are recalled and
 441 then filtered for subsequent processing.

442 **Metrics.** For Single-**VEQA** queries, we use accu-
 443 racy (**Acc**) as an evaluation metric. Furthermore,
 444 we assess the accuracy only for instances where
 445 relevant clues are successfully retrieved (*e.g.*,
 446 the case of Figure 1(c)), which is denoted as
 447 Acc^{hit} . For Group-**VEQA** queries, we employ
 448 recall (**Recall**) as the metric.

449 6.1 Single-VEQA Queries

450 The main results from the Single-**VEQA** queries
 451 are summarized in Table 2, which leads to the
 452 following insights:

453 1. **Model Parameter Size**: LLaVA-13b
 454 demonstrates higher accuracy (27.93%) com-
 455 pared to LLaVA-7b (22.26%), suggesting that

Models	Acc (%)	Acc ^{hit} (%)
Human	3.36	5.19
Human + FRAG	47.01	98.31
LLaVA-7b	22.26	27.53
LLaVA-7b + FRAG	31.19	62.81
LLaVA-13b	27.93	32.86
LLaVA-13b + FRAG	31.13	62.34
GPT-4V	-	-
GPT-4V + FRAG	34.84 (4.2)	68.31 (2.6)
MAR	39.09	79.65

Table 2: Result for Single-**VEQA** Queries. (Note: GPT-4V could not answer these queries directly due to policy constraints. Values within parentheses are those GPT-4V still refuses to answer.)

456 a model’s recognition ability is positively cor-
 457 related with its parameter size, which to some
 458 extent reflects its knowledge base.

459 2. **Impact of Matching Graph**: Incorporat-
 460 ing a matching graph leads to an 8.9% improve-
 461 ment in accuracy for LLaVA-7b and a 3.2%
 462 improvement for LLaVA-13b. GPT-4V, with
 463 matching, achieves a character recognition accu-
 464 racy of 34.83%.

465 3. **Comparative Improvement**: The en-
 466 hancement from matching is more pronounced
 467 for LLaVA-7b than for LLaVA-13b, indicating
 468 that while matching can compensate for differ-
 469 ences in parameters, a model’s inherent capabil-
 470 ities still set an upper limit on its performance.

471 To further understand the impact of matching
 472 on the **models’ reasoning abilities**, we analyzed
 473 examples of successfully recalled clues:

474 i. **Human Performance**: Human identifica-
 475 tion accuracy reaches 98.31% when incorporat-
 476 ing matching clues, setting a high benchmark
 477 for model performance.

478 ii. **Algorithmic Strength**: Our algorithm sur-
 479 passes others in analytical capabilities, achiev-
 480 ing an accuracy 11% higher than GPT-4V with
 481 matching in non-human results. However, there
 482 remains a gap compared to human performance.

483 iii. **Model Comparison**: Among LLaVA-7b,
 484 LLaVA-13b, and GPT-4V with matching, GPT-

Models	Recall
LLaVA-7b + FRAG	22.06%
LLaVA-13b + FRAG	40.05%
GPT-4V + FRAG	65.04%
MAR	70.85%

Table 3: Result for Group-VEQA Queries.

485 4V exhibits the best performance with an accu-
 486 racy of 68%, attributed to its superior analytical
 487 and reasoning abilities.

488 6.2 Group-VEQA Queries

489 Group-VEQA queries focus on identifying all per-
 490 tinent clues for more reliable reasoning. The
 491 result is shown in Table 3.

492 Our method achieves the highest recall rate
 493 at 70.85%, outperforming GPT-4V, LLaVA-7b,
 494 and LLaVA-13b combined with matching by
 495 5.81%, 30.81%, and 48.79%, respectively. This
 496 indicates that our approach excels in retrieval
 497 tasks compared to MLLMs, likely due to the ef-
 498 fectiveness of rule-based methods in managing
 499 excessive information. Additionally, the per-
 500 formance of baseline MLLMs diminishes with
 501 reduced parameter sizes, suggesting a positive
 502 correlation between their analytical reasoning
 503 abilities and parameter sizes.

504 6.3 Further Study - The Influence of 505 Multi-Source Information

506 In principle, the effective recognition of per-
 507 sonal information by a model depends on three
 508 main sources: its inherent knowledge, clues
 509 from the query, and clues from retrieved data.
 510 Our FRAG framework leverages these sources
 511 to guide accurate answers. As demonstrated
 512 in Table 4, when recall is accurate, LLaVA-7b
 513 correctly answers 42.86% of cases post-FRAG,
 514 while LLaVA-13b achieves 39.18%.

515 However, in practice, the presence of noise
 516 in the recalled information and the potential in-
 517 ability of MLLMs to effectively integrate FRAG
 518 information with the model’s original knowl-
 519 edge may lead to incorrect answers. As shown

Models	Acc ^{hit} (%)
LLaVA-7b	
w/o FRAG ✗ → with FRAG ✓	42.86
w/o FRAG ✓ → with FRAG ✗	7.32
LLaVA-13b	
w/o FRAG ✗ → with FRAG ✓	39.18
w/o FRAG ✓ → with FRAG ✗	9.44

Table 4: Study on Successfully Recalled Data.

520 in Table 4, LLaVA-7b+FRAG and LLaVA-
 521 13b+FRAG respectively provide incorrect an-
 522 swers in 7.32% and 9.44% of cases that could
 523 have been answered correctly before FRAG.

524 To assess the impact of the prompt on the
 525 model’s original knowledge, we conducted ab-
 526 lation experiments by removing the Original-
 527 knowledge-aware Prompt (OP), as shown in Ta-
 528 ble 5. The accuracy of LLaVA-7b, LLaVA-13b,
 529 and GPT-4V combined with FRAG decreased by
 530 6.05%, 1.72%, and 4.51% respectively. These
 531 results highlight the importance of the model’s
 532 own knowledge as a crucial clue in the reason-
 533 ing process and underscore its significance in
 534 achieving accurate outcomes.

Models	Acc
LLaVA-7b with matching	31.19%
w/o OP	25.14%
LLaVA-13b with matching	31.13%
w/o OP	29.41%
GPT-4V with matching	39.09%
w/o OP	34.58%

Table 5: Original-knowledge-aware Prompt (OP) ablation study result

535 7 Conclusion

536 In this paper, we explore a novel VEQA problem
 537 that focuses on aggregating clues from multiple
 538 captioned visual objects. We introduce match-
 539 ing graphs designed to capture the relationships
 540 between identical entities across various visual
 541 objects. Extensive experiments demonstrate the
 542 high accuracy of our method. While our work
 543 has primarily focused on matching person enti-
 544 ties, future research can aim to extend match-
 545 ing-augmented reasoning to other tasks.

546	Limitations		
547	Currently, our framework primarily relies on	Pan Zhou, Yao Wan, and Lichao Sun. 2024. Mllm-	586
548	similarity for face matching and does not con-	as-a-judge: Assessing multimodal llm-as-a-judge	587
549	sider factors such as age-related changes and	with vision-language benchmark. <i>arXiv preprint</i>	588
550	facial blurring. This may result in inaccuracies	<i>arXiv:2402.04788</i> .	589
551	in matching certain nodes, representing a fu-	Peter Christen and Peter Christen. 2012. <i>The data</i>	590
552	ture research direction. Additionally, in real-	<i>matching process</i> . Springer.	591
553	world applications, news is dynamic. Efficient	Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang	592
554	retrieval and expansion strategies for a growing	Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu,	593
555	data lake pose challenges as the dataset evolves,	Zichong Yang, Kuei-Da Liao, et al. 2024. A	594
556	warranting further investigation.	survey on multimodal large language models	595
557	Ethics Statement	for autonomous driving. In <i>Proceedings of the</i>	596
558	The authors declare that they have no conflict of	<i>IEEE/CVF Winter Conference on Applications of</i>	597
559	interest. Our work aims to enhance the answer	<i>Computer Vision</i> , pages 958–979.	598
560	generation of visual question answering by re-	Matthijs Douze, Alexandr Guzhva, Chengqi Deng,	599
561	trieving entity-related clues. While improving	Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel	600
562	the accuracy of answer generation, our method	Mazaré, Maria Lomeli, Lucas Hosseini, and	601
563	significantly saves resources as it does not re-	Hervé Jégou. 2024. <i>The faiss library</i> .	602
564	quire fine-tuning of large language models. We	Matthew Honnibal, Ines Montani, Sofie Van Lan-	603
565	strive to ensure that our approach is not only	deghem, and Adriane Boyd. 2020. <i>spaCy:</i>	604
566	accurate and efficient but also fair and unbiased.	<i>Industrial-strength Natural Language Processing</i>	605
567	We recognize the potential of significant impact	<i>in Python</i> .	606
568	of visual question answering technology on so-	Ilija Ilievski and Jiashi Feng. 2017. Multimodal	607
569	ciety and pledge to maintain transparency in	learning and reasoning for visual question answer-	608
570	sharing our findings and progress with relevant	ing. <i>Advances in neural information processing</i>	609
571	users and stakeholders.	<i>systems</i> , 30.	610
572	References	Mahmoud Khademi, Ziyi Yang, Felipe Frujeri,	611
573	Mayank Agrawal, Anand Singh Jalal, and Himanshu	and Chenguang Zhu. 2023. Mm-reasoner: A	612
574	Sharma. 2023. A review on vqa: Methods, tools	multi-modal knowledge-aware framework for	613
575	and datasets. In <i>2023 International Conference</i>	knowledge-based visual question answering. In	614
576	<i>on Computer Science and Emerging Technologies</i>	<i>Findings of the Association for Computational</i>	615
577	<i>(CSET)</i> , pages 1–6. IEEE.	<i>Linguistics: EMNLP 2023</i> , pages 6571–6581.	616
578	Ali Furkan Biten, Lluís Gomez, Marçal Rusinol,	Kuang-Huei Lee, Xi Chen, Gang Hua, Houdong Hu,	617
579	and Dimosthenis Karatzas. 2019. Good news,	and Xiaodong He. 2018. Stacked cross attention	618
580	everyone! context driven entity-aware caption-	for image-text matching. In <i>Proceedings of the</i>	619
581	ing for news images. In <i>Proceedings of the</i>	<i>European conference on computer vision (ECCV)</i> ,	620
582	<i>IEEE/CVF Conference on Computer Vision and</i>	pages 201–216.	621
583	<i>Pattern Recognition</i> , pages 12466–12475.	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	622
584	Dongping Chen, Ruoxi Chen, Shilin Zhang, Yinuo	Petroni, Vladimir Karpukhin, Naman Goyal,	623
585	Liu, Yaochen Wang, Huichi Zhou, Qihui Zhang,	Heinrich Küttler, Mike Lewis, Wen tau Yih,	624
		Tim Rocktäschel, Sebastian Riedel, and Douwe	625
		Kiela. 2021. <i>Retrieval-augmented generation for</i>	626
		<i>knowledge-intensive nlp tasks</i> .	627

628	Jiaqi Li, Miaozen Du, Chuanyi Zhang, Yongrui Chen, Nan Hu, Guilin Qi, Haiyun Jiang, Siyuan Cheng, and Bozhong Tian. 2024. Mike: A new benchmark for fine-grained multimodal entity knowledge editing. <i>arXiv preprint arXiv:2402.14835</i> .	671
629		672
630		673
631		
632		
633		
634	Kunpeng Li, Yulun Zhang, Kai Li, Yuanyuan Li, and Yun Fu. 2019. Visual semantic reasoning for image-text matching. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 4654–4662.	
635		
636		
637		
638		
639	Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, and Lingpeng Kong. 2023a. Silk: Preference distillation for large visual language models. <i>arXiv preprint arXiv:2312.10665</i> .	
640		
641		
642		
643		
644	Yunxin Li, Longyue Wang, Baotian Hu, Xinyu Chen, Wanqi Zhong, Chenyang Lyu, and Min Zhang. 2023b. A comprehensive evaluation of gpt-4v on knowledge-intensive visual question answering. <i>arXiv preprint arXiv:2311.07536</i> .	
645		
646		
647		
648		
649	Yuanze Lin, Yujia Xie, Dongdong Chen, Yichong Xu, Chenguang Zhu, and Lu Yuan. 2022. Re- vive: Regional visual representation matters in knowledge-based visual question answering. <i>Advances in Neural Information Processing Systems</i> , 35:10560–10571.	
650		
651		
652		
653		
654		
655	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning .	
656		
657	Ziyu Liu, Zeyi Sun, Yuhang Zang, Wei Li, Pan Zhang, Xiaoyi Dong, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. 2024. Rar: Retrieving and ranking augmented mllms for visual recognition. <i>arXiv preprint arXiv:2403.13805</i> .	
658		
659		
660		
661		
662	Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. 2021. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. <i>arXiv preprint arXiv:2110.13214</i> .	
663		
664		
665		
666		
667		
668	TR Mahesh, T Rajan, K Vanitha, HK Shashikala, et al. 2023. Intelligent systems for medical diag- nostics with the detection of diabetic retinopathy at reduced entropy. In <i>2023 International Con- ference on Network, Multimedia and Information Technology (NMITCON)</i> , pages 1–8. IEEE.	671
669		672
670		673
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable vi- sual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	674
		675
		676
		677
		678
		679
		680
	Elias Stengel-Eskin, Jimena Guallar-Blasco, Yi Zhou, and Benjamin Van Durme. 2022. Why did the chicken cross the road? rephrasing and analyzing ambiguous questions in vqa. <i>arXiv preprint arXiv:2211.07516</i> .	681
		682
		683
		684
		685
	Zhou Su, Chen Zhu, Yinpeng Dong, Dongqi Cai, Yurong Chen, and Jianguo Li. 2018. Learning visual knowledge memory networks for visual question answering. In <i>Proceedings of the IEEE conference on computer vision and pattern recog- nition</i> , pages 7736–7745.	686
		687
		688
		689
		690
		691
	Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. 2014. Deepface: Closing the gap to human-level performance in face verification . In <i>2014 IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 1701–1708.	692
		693
		694
		695
		696
	Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordani, Philip Bachman, and Kaheer Suleman. 2016. Newsqa: A ma- chine comprehension dataset. <i>arXiv preprint arXiv:1611.09830</i> .	697
		698
		699
		700
		701
	Qifan Yu, Juncheng Li, Longhui Wei, Liang Pang, Wentao Ye, Bosheng Qin, Siliang Tang, Qi Tian, and Yueting Zhuang. 2023. Hallucidoctor: Miti- gating hallucinatory toxicity in visual instruction data. <i>arXiv preprint arXiv:2311.13614</i> .	702
		703
		704
		705
		706
	Dongxiang Zhang, Rui Cao, and Sai Wu. 2019. In- formation fusion in visual question answering: A survey. <i>Information Fusion</i> , 52:268–280.	707
		708
		709
	Xinlu Zhang, Yujie Lu, Weizhi Wang, An Yan, Jun Yan, Lianke Qin, Heng Wang, Xifeng Yan, William Yang Wang, and Linda Ruth Petzold. 2023. Gpt-4v(ision) as a generalist evaluator for vision-language tasks .	710
		711
		712
		713
		714

- 715 Wenfeng Zheng, Yu Zhou, Shan Liu, Jiawei Tian,
716 Bo Yang, and Lirong Yin. 2022. A deep fusion
717 matching network semantic reasoning model. *Ap-
718 plied Sciences*, 12(7):3416.
- 719 Jie Zhu, Shufang Wu, Xizhao Wang, Guoqing Yang,
720 and Liyan Ma. 2018. Multi-image matching
721 for object recognition. *IET Computer Vision*,
722 12(3):350–356.