
UDKAG: Augmenting Large Vision-Language Models with Up-to-Date Knowledge



Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large vision-language models (LVLMs) are ignorant of the up-to-date knowledge,
2 such as LLaVA series, because they cannot be updated frequently due to the large
3 amount of resources required, and therefore fail in many cases. For example, if a
4 LVLM was released on January 2024, and it wouldn't know the detailed plot of the
5 new movie Dune 2, such as "the dead of an important duel", because this movie was
6 not yet released when the LVLM was trained. To solve the problem, a promising
7 solution is to provide LVLMs with up-to-date knowledge via internet search during
8 inference, *i.e.*, internet-augmented generation (IAG), which is already integrated in
9 some closed-source commercial LVLMs such as GPT-4V. However, the specific
10 mechanics underpinning them remain a mystery. In this paper, we propose a plug-
11 and-play framework, for augmenting existing LVLMs in handling visual question
12 answering (VQA) about up-to-date knowledge, dubbed UDKAG. A hierarchical
13 filtering model is trained to effectively and efficiently find the most helpful content
14 from the websites returned by a search engine to prompt LVLMs with up-to-date
15 knowledge. To train the model and evaluate our framework's performance, we propose
16 a pipeline to automatically generate news-related VQA samples to construct
17 a dataset, dubbed UDK-VQA. A multi-model voting mechanism is introduced to
18 label the usefulness of website/content for VQA samples to construct the training
19 set. For the test set, we perform manual screening to ensure the correctness of
20 test samples. Experimental results demonstrate significant improvements of our
21 framework over LVLMs, outperforming the self-contained IAG-capable GPT-4V
22 by ~25% in accuracy on UDK-VQA test set.

23 1 Introduction

24 Large vision-language models (LVLMs, *e.g.*, GPT-4V [1], Gemini Series [2], and Grok [3]) have
25 received much attention for their impressive generative capabilities. They require a large resource for
26 data collection, cleaning, and training, restricting them from frequently updating models. However,
27 new information and knowledge are created every time, making LVLMs ineffective in many scenarios.
28 For example, if we talk with LLaVA-1.6 [4] (released on January 30, 2024) about the detailed plot
29 of the new movie Dune 2, such as "the dead of an important duel", it performs very badly. It is
30 promising to augment LVLMs by retrieving up-to-date knowledge via internet search during inference,
31 *i.e.*, internet-augmented generation (IAG). Although commercial LVLMs such as GPT-4V [1] and
32 Claude3 [5] have the ability of IAG, the specific mechanics underpinning them remain undisclosed.
33 This paper proposes a plug-and-play framework to augment different LVLMs in handling visual
34 question answering (VQA) about up-to-date knowledge, named UDKAG.

35 We first introduce our overall framework applicable to different LVLMs for equipping them with
36 up-to-date knowledge during inference. It consists of four components: query generator, search

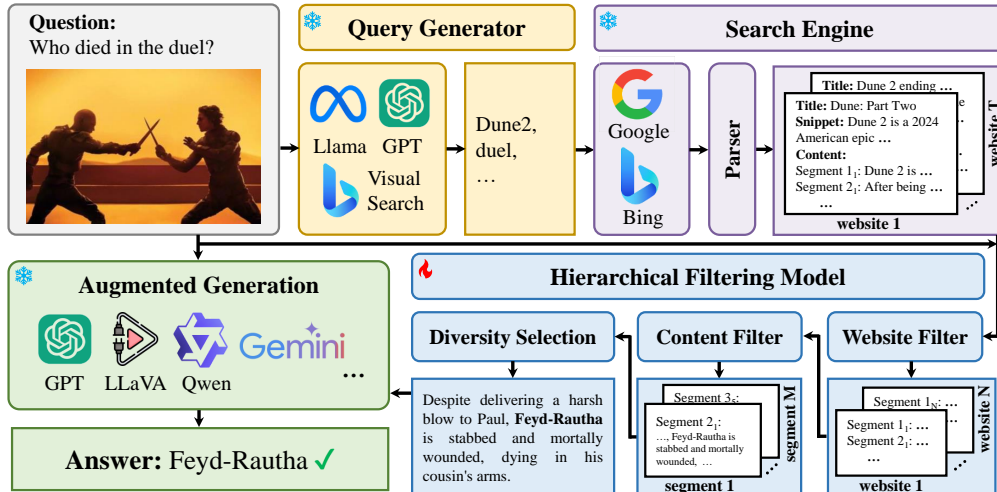


Figure 1: The proposed UDKAG, a framework for LVLMs to access up-to-date knowledge.

37 engine, hierarchical filtering model, and augmented generation, as shown in Figure 1. Specifically, we
 38 begin by extracting queries via Bing Visual Search and LLMs for an image-related question. Then, we
 39 acquire helpful websites through search engines and extract their contents by web scraping. However,
 40 it is impractical to augment LVLMs directly with the entire content of all websites, because: (1) Most
 41 LVLMs are poor at handling such long contexts. (2) Handling such long contexts is computationally
 42 intensive and time-consuming. To this end, a hierarchical filtering model is trained to find the most
 43 helpful content for answering the question, which first efficiently sifts the websites based on each
 44 website’s title and snippet, and then identifies the most helpful content from the filtered websites.
 45 Finally, the filtered content is fed to LVLMs to assist them in answering the question.

46 We then construct a dataset dubbed UDK-VQA about up-to-date news. It is used to train the
 47 hierarchical filtering model and also evaluate our overall framework’s performance. In particular,
 48 we propose a pipeline to automatically scrape the up-to-date news and generate news-related VQA
 49 samples. Specifically, we use search terms from Google Daily Search Trends and manually collected
 50 popular search terms as queries to search for hot news. For each piece of news, we divide its content
 51 into segments and ask GPT-3.5 to generate question-answer pairs based on each segment. Then, we
 52 extract an entity for each question and replace it with its hypernym. To compose a VQA sample,
 53 we use Bing to search images of the replaced entity and cluster them to reduce the outliers among
 54 them. In doing so, answering the generated VQA samples requires models to consider both visual
 55 and textual information. We use queries from different time periods to scrape news from different
 56 time periods to generate samples for constructing a training set and a test set, to avoid the test data
 57 being exposed in the training data. In the training set, we further use a multi-model voting mechanism
 58 to label website’s usefulness and content’s usefulness for VQA samples, and combine the samples
 59 with websites and their content based on the label for training the hierarchical filtering model. In the
 60 testing set, we conduct manual screening to ensure its correctness.

61 To validate the effectiveness and generalizability of the proposed framework, we incorporate 13
 62 state-of-the-art LVLMs into the framework, such as GPT-4V [1] and LLaVA-1.6 [4]. Notably,
 63 once the hierarchical filtering model is trained, our framework can adapt different LVLMs and
 64 improve their performance without any fine-tuning. Extensive experimental results demonstrate
 65 that our framework can significantly improve LVLm’s ability to answer questions about up-to-
 66 date knowledge. Incorporating the LLaVA-1.6 model of our framework even outperforms the
 67 self-contained IAG-capable GPT-4V by ~25% in accuracy on UDK-VQA test set.

68 The contributions of our work are summarized as follows. (1) We propose the first open-source
 69 framework which seamlessly incorporates existing LVLMs to equip them with up-to-date knowledge
 70 during inference. (2) We propose a pipeline which can automatically generate VQA samples related
 71 to up-to-date news, and construct the first test set for evaluating the ability of LVLMs in handling
 72 VQA about up-to-date knowledge. (3) Extensive experimental results on 13 state-of-the-art LVLMs
 73 demonstrate the effectiveness of our framework.

74 2 Related Work

75 2.1 Retrieval-Augmented Generation

76 Recently retrieval-augmented generation (RAG) attracted increasing attention of both the natural
77 language processing [6, 7, 8, 9] and vision-and-language [10, 11, 12]. REALM [6] uses the query to
78 retrieve the top k most relevant article snippets, and uses large language models (LLMs) to generate
79 k responses, which are then combined to obtain a final output for question answering. Recently,
80 [13, 8, 14] explores the internet-augmented generation (IAG) of LLMs to enable language models
81 to access up-to-date information via search engines. Komeili *et.al.* [13] show that LLMs enhanced
82 via search engines can generate less factually incorrect information during dialogue with humans.
83 Lazaridou *et.al.* [8] uses few-shot prompting to enable LLMs to exploit knowledge returned from
84 Google search to answer questions about factual and up-to-date information. In vision-and-language,
85 REVEAL [10] builds a memory by encoding open-world knowledge including image-text pairs,
86 question-answering pairs, etc., and uses a retriever to find the most relevant knowledge entries in
87 the memory. The memory, encoder, retriever, and generator are pre-trained in an end-to-end manner.
88 Re-ViLM [11] augments Flamingo [15], by retrieving relevant image-text pairs from the external
89 image-text datasets [16, 17, 18] for zero and in-context few-shot image-to-text generations. RA-CM3
90 [19] performs retrieval from an external memory for generating images and text. Differently, we
91 focus on enabling LVLMs to retrieve up-to-date knowledge via Internet search during inference.

92 2.2 Large Models with Search Engine

93 Recent years have witnessed a growing interest in exploring external tools for LLMs [20, 21, 22, 23,
94 24]. Among them, some methods [24, 25, 26] can use search engines to access up-to-date knowledge.
95 Nonetheless, these methods usually focus on how to appropriately use different tools to enhance
96 LLMs, such as using Python interpreter to generate complex programs [21], incorporating more
97 external tools [24], or updating tools by acquiring new knowledge [26]. Although they can access
98 up-to-date knowledge, they usually directly use the website snippets for augmenting generation.
99 By contrast, this work focuses on internet-augmented generation and explores how to obtain more
100 relevant up-to-date knowledge and effectively use retrieved knowledge to augment LVLMs.

101 3 Framework

102 In this section, we introduce UDKAG, a framework that seamlessly incorporate existing LVLMs,
103 allowing these LVLMs to access up-to-date knowledge without fine-tuning. The whole framework is
104 illustrated in Figure 1. For a natural language question Q about an image V , we first extract queries
105 for both Q and V via the query generator. Then we enter the queries into search engines, and the
106 search engine would return related websites, each of which consist of a title and a snippet. To identify
107 the most helpful content within the websites, a website filter is used to filter the websites based on
108 their titles and snippets, and a content filter is further used to filter the content of the websites filtered
109 by the website filter. Finally, we stitch the filtered content together to prompt existing LVLMs.

110 3.1 Query Generator

111 **Question Query Generator.** To get queries that make search engines return websites containing
112 helpful content, we leverage large language models (LLMs) to extract queries for Q . Thanks to the
113 language understanding capability of LLMs, the role played by each word can be well inferred from
114 the grammatical information of Q , even if certain words are unknown for the LLMs. We use “*Do not*
115 *try to answer the question, just print the most informative no more than three entities in the question.*
116 *Put them on one line and separate them with comm.*” to prompt LLMs to generate queries.

117 **Image Query Generator.** For an image V , we leverage Bing Visual Search to analyze the image
118 entities of V as queries. The reason for using Bing Visual Search rather than a LVLM to extract
119 queries for V is that current LVLMs are inadequate in extracting image entities especially for emerging
120 entities. Notably, Bing Visual Search is a tool different from commonly used search engines, returning
121 image-related attributes, including image entity names, image-related search terms and image-related
122 websites. However, entity names are missing in most cases. To address this problem, we extract the
123 longest public ancestor of related search terms and related website titles as the queries for V .

124 3.2 Search Engine

125 The extracted queries are fed into a search engine, and the search engine returns relevant websites
126 with their titles and snippets. However, the returned titles and snippets often contain limited and
127 incomplete information. For example, for a website with title “*Pororo Dragon Castle Adventure*”, the
128 entire snippet returned by Bing is “*Pororo and his friends were having fun when a little red dragon*
129 *named Arthur appears above them Arthur who claims to be the king of dragons commands Pororo*
130 *and his friends to search for his Dragon ...*”, obviously there is more about “*Pororo Dragon Castle*
131 *Adventure*” contained in the website. Thus we parse the textual content of all websites. For a website,
132 not all of its content contributes to answering questions, we empirically divide the website content
133 into segments every third sentence for a more granular selection of content.

134 3.3 Hierarchical Filtering Model

135 Since most of the existing LVLMs cannot receive long context as inputs, and long contexts can be
136 computationally intensive and time consuming for them, it’s necessary to filter the website content
137 after obtaining the websites via the search engine. Towards this goal, we train a hierarchical filtering
138 model, which consists of a website filter and a content filter to perform a two-step filtering.

139 **Website Filter.** The aim of the website filter is to perform the filtering of websites based on their
140 titles and snippets. Specially, a website scoring model is trained via instruction tuning, to predict how
141 helpful a website will be in answering a question, and the N websites with higher scores would be
142 kept. The training samples are in the format (T, S, Q, V, R_w) , where R_w is a quantitative usefulness in
143 the interval $[0, 1]$ representing how helpful a website with title T and snippet S will be in answering
144 Q related to V . Based on the samples, we construct instructions like “*How helpful is an article with*
145 *such a title and snippet in answering the question based on the image? Choose the best option. Title:*
146 *<T> Snippet: <S> Question: <Q> Options: A. 1.0 B. 0.8 C. 0.6 D. 0.4 E. 0.2 F. 0.0*”. In doing so,
147 the score regression problem is converted into a classification problem, which is easier to learn.

148 **Content Filter.** The content filter is used to select the most helpful content segments from the
149 websites filtered by the website filter. For each content segment, we predict how helpful is it for
150 answering Q by a content scoring model. The content scoring model is trained by samples in the
151 format (C, Q, V, R_c) , where C is a content segment, and R_c is the quantitative usefulness of C in
152 answering Q . The instructions for training the content scoring model are in the format: “*How helpful*
153 *is this context in answering the question based on the image? Choose the best option. Context: <C>*
154 *Question: <Q> Options: A. 1.0 B. 0.8 C. 0.6 D. 0.4 E. 0.2 F. 0.0*”. We use the model to sort all
155 content segments and select the M highest scoring ones as the obtained segments.

156 **Diversity Selection.** To avoid LVLMs answer questions using bias from repetitive contexts, we
157 performed a quadratic selection on the obtained segments based on diversity. Specially, we extract
158 CLIP features [27] for all the segments and cluster them using k-means [28]. The segments closest to
159 the center of each cluster are stitched together as the final obtained content for prompting the LVLMs.

160 3.4 Augmented Generation

161 We augment existing LVLMs by prompting them with the final obtained content, to improve their
162 ability of answering questions about up-to-date knowledge. Taking answer the multiple choice
163 questions as an example, for a question Q with candidate answers A_1, A_2, A_3 and A_4 , we use the
164 prompt “*Given context: <X> Question: <Q> Answers: A.<A₁> B.<A₂> C.<A₃> D.<A₄> Answer*
165 *with the option’s letter from the given choices directly based on the context and the image.*”, where X
166 denotes the final content obtained by the hierarchical filtering model.

167 4 UDK-VQA Dataset

168 To evaluate the effectiveness of our framework, we propose a pipeline to automatically scrape the
169 up-to-date news and generate news-related VQA. The whole pipeline is demonstrated in Figure 2.
170 The pipeline is also used to collect training samples for the hierarchical filtering model. We first
171 collect hot search terms as queries to scrape relevant news returned by search engines. For each
172 piece of news, every third sentence is divided into a segment. We then employ GPT 3.5 to generate
173 a question-answer pair for each segment, and extract an entity in the question, replacing it with its

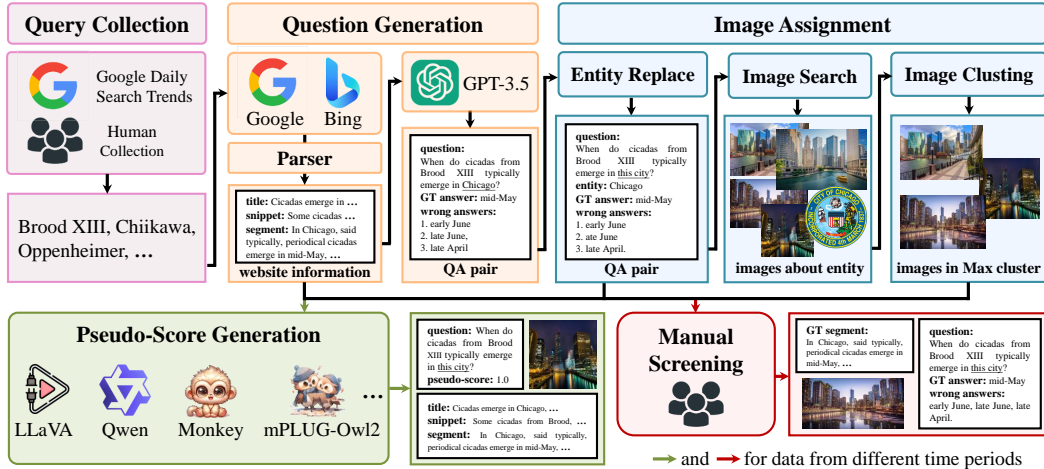


Figure 2: Overall pipeline of the sample generation for the UDK-VQA dataset. For brevity, we only show one output item at several steps, such as the content segment returned by the Parser. Notably, we use queries from different time periods to scrape news from different time periods to generate training samples and test samples, which is not reflected in this figure for brevity.

174 hypnym. Bing Image Search is used to find images for the replaced entity, and after removing
 175 outliers from the images using clustering, the images and the question-answer pair are composed
 176 into VQA samples. We combine the VQA samples and website information (*e.g.*, title, snippet and
 177 content), and introduce a multi-model voting mechanism to generate pseudo-score, constituting the
 178 training set. For the test set, manual screening is conducted to ensure the correctness of test samples.

179 4.1 Query Collection

180 Google daily search trends is an available data source that reflects what’s hot in real time, and is
 181 well suited as the query used to construct our dataset. However, we observe that most search terms
 182 of the Google daily search trends are related to politics and sports, which poses a great limitation.
 183 Therefore, we further manually collect popular search terms to improve the query diversity. The
 184 popular search terms are collected from many other domains including films, technological products,
 185 anime characters, places of interest, and so on. These human-collected queries were mixed with
 186 queries from Google daily search trends to be used for subsequent sample generation.

187 4.2 Question Generation

188 For each query, we use Bing to search for relevant and up-to-date news. For the scraped news content,
 189 we divide every third sentence into a segment, and use the following message to prompt GPT-3.5
 190 to generate a question-answer pair and several confused answers for each segment: “*Given context:*
 191 *<Con> Filling the blanks to generate a question about the most informative event of the context,*
 192 *generate an correct answer to the question in no more than three words based on context, and*
 193 *generate three incorrectly confused answers of no more than three words based on context. Question:*
 194 *___ Correct answer: ___ Incorrect answers: A. ___ B. ___ C. ___*”, where *<Con>* denotes a segment.

195 We design a simple but effective rule to ensure the correctness of the generated question-answer pairs,
 196 which requires a model can answer a question Q with A based on a segment C , if the model is able to
 197 generate a question-answer pair (Q, A) based on C .

198 4.3 Image Assignment

199 To generate VQA samples and avoid the model’s reliance on language priors for answering, we create
 200 samples that necessitate an understanding of the image for correct answers. Firstly, we extract an
 201 entity for each question via named an entity recognition (NER) model [29]. Images of the entity
 202 are then obtained by Bing Image Search. Since the images returned by the search engine are noisy

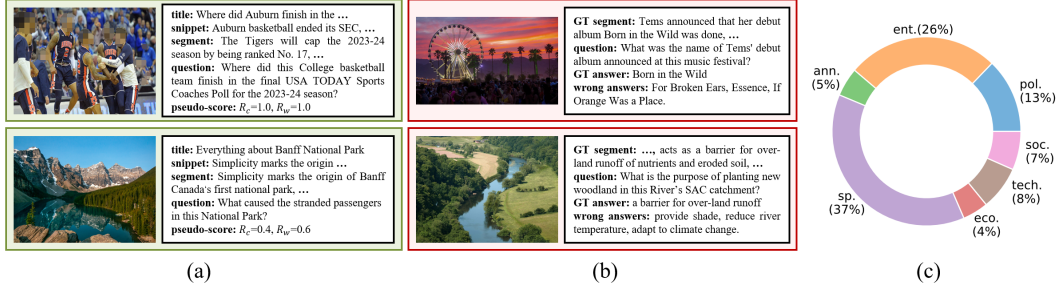


Figure 3: (a) Training samples. (b) Test samples. (c) Category statistics for the test set of UDK-VQA.

203 (outlier images), we cluster them based on the CLIP feature [27] of the images and keep only the
 204 images in the cluster with the highest number of images. Finally, the kept images are assigned to the
 205 new question-answer pairs where the entity is replaced by its hypernym, to compose VQA samples.
 206 An obtained VQA sample can be denoted as $(V, Q, A_{gt}, \{A_w^i\}_{i=1}^3)$, where Q is the generated question,
 207 V is the entity image, A represents the ground-truth answer, and A_w^i is i -th confused wrong answer.

208 4.4 Pseudo-Score Generation

209 For a VQA sample generated from the content segment C , which we denote as the ground-truth
 210 segment for the sample, it is certain that C is most helpful in answering this sample. Inevitably, we
 211 must consider to what extent do the other content segments contribute to answering the sample? We
 212 propose a pseudo-score generation method that uses five LVLMS for voting to quantify how helpful a
 213 content segment is to a VQA sample into six values: 1.0, 0.8, 0.6, 0.4, 0.2 and 0.0. Specially, for a
 214 VQA sample with the ground-truth segment C from a news fetched for a query, we first sample four
 215 content segments from the news for the query beyond its ground-truth segment. Then we use each
 216 sampled segment to prompt each of the five LVLMS to answer the VQA sample and count the rate of
 217 LVLMS that answer correctly as the pseudo-score for the segment.

218 In doing so, we obtain training samples for the content filter, in the format (C, Q, V, R_c) , where C is
 219 a content segment, Q denotes a question related to the image V , and R_c is the pseudo-score of how
 220 helpful C is to answer Q . Moreover, we count the maximum pseudo-score of all content segments in
 221 a news for a VQA sample as the pseudo-score for the news website, dubbed R_w , to build training
 222 samples for the website filter. The training sample format for the website filter is (T, S, Q, V, R_w) ,
 223 where T is the website title, S is the website snippet. By merging these training samples into the
 224 training instructions mentioned in Section 3.3, the hierarchical filtering model can be implemented.

225 4.5 Manual Screening

226 For constructing the test set, we do not use the pseudo-score generation method. A test sample
 227 $(C, V, Q, A_{gt}, \{A_w^i\}_{i=1}^3)$ can be seen as a VQA sample with its ground-truth content segment C . It is
 228 worth noting that C is only provided when testing the upper bound of performance. For each test
 229 sample, we randomly mix A_{gt} and $\{A_w^i\}_{i=1}^3$, then assign them the options (*i.e.* A, B, C and D), and
 230 add a complementary option E . *No Correct Answers*, to evaluate LVLMS in a multiple choice format.
 231 Moreover, we manually review all test samples to ensure that they are correct.

232 4.6 Dataset Analysis

233 To avoid the test data being exposed in the training set, we use queries from different time periods
 234 to scrape news from different time periods for constructing the training set and the test set. For the
 235 training set, we use the queries from February 17, 2024 to March 31, 2024 to scrape news before
 236 April 10, 2024. The training sample number for the website filter and the content filter are 599, 700
 237 and 850, 267, respectively. For the test set, we use the queries from April 1, 2024 to April 31, 2024
 238 to scrape news after April 10, 2024. The number of test samples is 1,000. We manually divide the
 239 test sample into seven categories, including politics, entertainment, announcement, sports, economic,
 240 technology and society, based on their required knowledge. We visualize some samples in UDK-VQA
 241 and the statistics for test samples in each category in Figure 3.

Table 1: Comparison with SOTA LVLMS on UDK-VQA, where Raw represents the model without IAG ability (e.g., official API version), IAG represents the model with self-contained IAG-capable ability (official web version), LC represents the model with long context input. “Ours” stands for incorporating the Raw baseline into our framework. The value outside/in () indicates the accuracy over samples that do not violate the content management policy of current/all model(s).

Model Variant		pol.	ent.	ann.	sp.	eco.	tech.	soc.	overall
Gemini 1.5 Pro	Raw	6.2 (5.7)	15.8 (16.3)	10.2 (11.9)	7.4 (8.1)	2.3 (2.5)	8.0 (6.2)	3.0 (3.7)	9.1 (9.5)
	LC	61.7 (65.7)	71.5 (77.3)	73.5 (76.2)	77.2 (79.2)	72.7 (72.5)	81.3 (83.1)	62.1 (66.7)	76.4 (76.1)
	Ours	82.8 (82.9)	79.6 (79.0)	91.8 (92.9)	81.5 (80.1)	97.7 (97.5)	84.0 (83.1)	90.9 (88.9)	83.3 (82.3)
GPT 4V	Raw	21.1 (23.8)	31.5 (30.9)	16.3 (19.0)	16.7 (17.5)	15.9 (17.5)	41.3 (38.5)	21.2 (22.2)	24.2 (23.8)
	IAG	62.5 (68.0)	61.9 (63.6)	63.3 (66.7)	62.4 (63.3)	70.5 (67.5)	80.0 (78.5)	69.7 (70.4)	64.5 (65.9)
	Ours	76.6 (83.8)	85.8 (85.8)	89.8 (90.5)	86.2 (86.4)	97.7 (97.5)	92.0 (90.8)	87.9 (92.6)	87.2 (87.4)
GPT 4o	Raw	36.7 (40.0)	34.2 (36.1)	42.9 (47.6)	28.6 (30.4)	40.9 (42.5)	65.3 (67.7)	36.4 (38.9)	37.2 (37.8)
	IAG	61.7 (63.1)	57.3 (58.4)	69.4 (64.3)	48.4 (47.9)	70.5 (75.0)	81.3 (81.5)	62.1 (61.1)	57.8 (57.9)
	Ours	86.7 (92.4)	89.6 (91.4)	98.0 (100)	83.9 (88.0)	97.7 (100)	90.7 (96.9)	89.4 (94.4)	91.8 (91.6)
LLaVA 1.6	Raw	43.8 (45.7)	32.3 (31.8)	22.4 (23.8)	24.9 (23.2)	25.0 (25.0)	53.3 (55.4)	33.3 (31.5)	31.8 (31.2)
	Ours	86.7 (87.6)	91.9 (91.8)	93.9 (97.6)	88.1 (88.0)	90.9 (90.0)	93.3 (93.8)	95.5 (100)	90.2 (90.7)

242 5 Experiments

243 5.1 Settings

244 **Training.** We implement two versions of the hierarchical filtering model, one using LLaVA-1.5-
 245 vicuna-7b [30] and the other using Qwen-VL-Chat [31]. In each version, we use same hyper-
 246 parameters to fine-tune two same LVLMS with LoRA [32] as the website filter and the content
 247 filter, respectively. Whether fine-tuning LLaVA-1.5-vicuna-7b or Qwen-VL-Chat, the entire training
 248 process is facilitated on two Nvidia A100 GPUs, using a batch size of 128 over 3 epochs.

249 **Baselines.** We incorporate 13 representative LVLMS into the proposed framework including Gemini
 250 1.5 Pro [2], GPT-4V [1], GPT-4o, LLaVA-1.6 [4], XComposer2 [33], Monkey [34], CogVLM [35],
 251 MiniCPM-V2 [36], mPLUG-Owl2 [37], Qwen-VL [31], MMAIaya [38], Xtuner [39] and VisualGLM
 252 [40]. We implement Gemini 1.5 Pro, GPT-4V and GPT-4o via their official webs and APIs. For other
 253 LVLMS, we implement them based on the VLMEvalKit toolkit [41].

254 **Evaluation.** We evaluate LVLMS on the test set of our UKD-VQA dataset. In addition to evaluating
 255 LVLMS via VLMEvalKit, we design additional matching patterns for each LVLMS with respect to its
 256 answer format. For example, we additionally use the pattern “The answer is XXX.” for XComposer2
 257 as it often answers in this format. All evaluations are conducted with a single Nvidia A100 GPU.

258 5.2 Quantitative Comparison with SOTA LVLMS

259 We compare with state-of-the-art LVLMS on the UDK-VQA test set, including Gemini 1.5 Pro [2],
 260 GPT-4V [1], GPT-4o and LLaVA-1.6 [4]. For Gemini 1.5 Pro, GPT-4V and GPT-4o, we implement
 261 their **Raw** version via official APIs, which do not have the ability of IAG. Since Gemini 1.5 Pro
 262 is famous for its ability to receive long contexts, we use all website content returned by the search
 263 engine of our framework to prompt it directly, dubbed **LC**. For GPT-4V and GPT-4o, we test their
 264 self-contained IAG-capable ability via prompting their official web versions with “*Retrieve relevant*
 265 *news and answer the question directly from the given options using the option letters based on the*
 266 *image.*”, dubbed **IAG**. We incorporate each Raw baseline into our framework as **Ours**.

267 Experimental results are listed in Table 1, we can observe that: (1) GPT-4o with our framework
 268 achieves the best performance on almost categories of UDK-VQA. (2) For all four baselines, our
 269 framework consistently improves their accuracy (e.g., 22.7% and 34.0% absolute performance gains
 270 in overall accuracy for GPT-4V and GPT-4o, respectively). (3) Our framework uses shorter contexts
 271 but has higher accuracy (e.g., 76.4% vs 83.3% in accuracy for LC and Ours variants of Gemini,
 272 respectively). The observations suggest that our framework is generalizable and effective in enhancing
 273 the ability of LVLMS to answer questions about up-to-date knowledge.

Table 2: Ablation studies of our framework on UDK-VQA.

Model	Variant	Hierarchical Filtering Model		Query Generator (Q)			Query Generator (V)	Acc. (%)	
		LLaVA-1.5	QWen-VL	NER	LLaMA3	GPT-3.5	Bing Visual Search		
LLaVA-1.6	Raw	-	-	-	-	-	-	31.8	
	IAG (SIM Q)	-	-	-	-	-	-	46.1	
	IAG (SIM V)	-	-	-	-	-	-	47.1	
	IAG (SIM QV)	-	-	-	-	-	-	47.7	
	Ours		✓	-	-	-	-	✓	49.3
			✓	-	-	-	✓	-	65.9
			✓	-	✓	-	-	✓	81.4
			✓	-	-	✓	-	✓	86.6
			✓	-	✓	✓	✓	✓	87.6
			✓	-	✓	✓	✓	✓	90.2
	-	✓	✓	✓	✓	✓	89.6		

274 **5.3 Ablation Studies**

275 The results of ablation studies on UDK-VQA are shown in Table 2. Firstly, we investigate simple
 276 IAG methods, including using the similarity between questions and segments to select segments, *i.e.*,
 277 **IAG (SIM Q)**, using the similarity between images and segments to select segments *i.e.*, **IAG (SIM**
 278 **V)**, using the averaged similarity of the the above two similarities to select segments, *i.e.*, **IAG (SIM**
 279 **QV)**. These methods show limited improvements and achieve unsatisfactory accuracy.

280 Then, we study the influences of different components of our framework on the performance. For
 281 the hierarchical filtering model, we study two popular LVLMS, LLaVA-1.5 [30] and Qwen-VL [31].
 282 For the query generator, we conduct experiments with NER [29], LLaMA3 [42], GPT-3.5 and Bing
 283 Visual Search. We observe that: (1) Using different backbone for the hierarchical filtering model has
 284 little effect on performance. (2) Using multiple question query generators at the same time can result
 285 in better performance than using only one. (3) Using both the question query generator and the image
 286 query generator gives the best performance. These observations suggest that all components of our
 287 framework are effective in improving the baseline, and components are complementary to each other.

288 **5.4 Analysis of Pseudo-Score Generation**

289 We analyze the influences of using different LVLMS to gener-
 290 ate pseudo-scores on the performance. We categorize 10
 291 LVLMS into two groups based on their released date, the first
 292 group contains LLaVA-1.6, XComposer2, Monkey, CogVLM
 293 and MiniCPM-V2, the second group contains mPLUG-Owl2,
 294 Qwen-VL, MMAlaya, Xtuner and VisualGLM. Using these
 295 two groups to generate pseudo-scores are dubbed **PSG with G1**
 296 and **PSG with G2**. Experimental results are shown in Figure 4,
 297 which reveal that: (1) The proposed framework can be directly
 298 used to boost LVLMS that are not used for generating pseudo-
 299 scores, which show the transferability of our framework. (2)
 300 The use of more recent LVLMS for generating pseudo-scores al-
 301 lows for greater improvements in general. (3) Different LVLMS
 302 have different performance upper bound, some of them achieve
 303 limited accuracy (*e.g.*, $\sim 70\%$ in accuracy for VisualGLM) even
 304 are augmented with ground-truth segments (**GT Segment**).

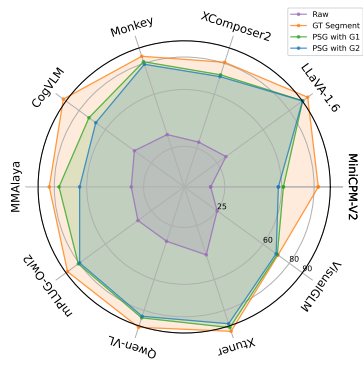


Figure 4: Accuracy using different LVLMS to generate pseudo-scores.

305 **5.5 Analysis of Diversity Selection**

306 In this section, we investigate the necessity of diversity selection. We compare our diversity selection
 307 (Div- K) with Top- K selection, and the experimental results of 10 LVLMS are shown in Figure 5.
 308 The Top- K selection means stitching K content segments with the highest scores together to prompt
 309 the LVLMS. For Div- K , K denotes the number of clusters. Experimental results demonstrate that:

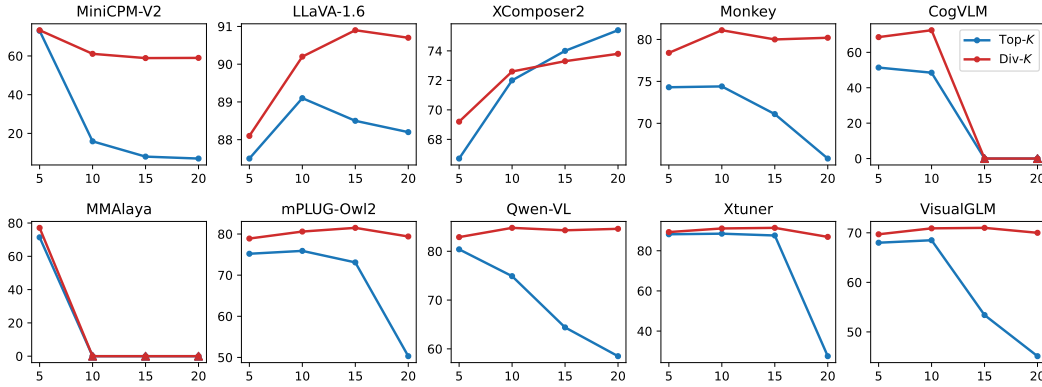


Figure 5: Comparison between Top- K selection and diversity selection (Div- K), where K denotes the number of stitched content segments for prompting LVLMs. For each sub-figure, the horizontal coordinate is K and the vertical coordinate is the accuracy. Note that an accuracy of 0 means that the model fails at the context length under the current setting of K , and is labeled as a triangle.

310 (1) Our diversity selection outperforms the Top- K selection regardless of the setting of K for most
 311 LVLMs. (2) As K increases, the performance using the Top- K selection plummets. This is because
 312 content with high scores is similar, and if a LVLm receives too many duplicate content as inputs, it
 313 will misinterpret the instruction and thus repeat the inputs instead of answering the question. These
 314 experimental results prove the necessity and effectiveness of the diversity selection.

315 5.6 Analysis of Website Filter

316 An important capability of the website filter is the trade-off
 317 between the content filter efficiency and the LVLms' accuracy.
 318 Adjusting the filtered website number N can control the token
 319 number that the content filter needs to process as a percent-
 320 age of the total token number returned by the search engine,
 321 dubbed θ . The variation in accuracy of LVLms as θ increases
 322 is shown in Figure 6, we can observe that: (1) The accuracy of
 323 LVLms increases with θ , especially when $\theta \leq 40\%$. (2) The
 324 increase in accuracy of LVLms slows down after $\theta \geq 40\%$.
 325 Therefore, setting $\theta = 40\%$ achieves a better trade-off, because
 326 the accuracy obtained by processing 40% tokens is close to
 327 98% of the accuracy obtained when processing 100% tokens.

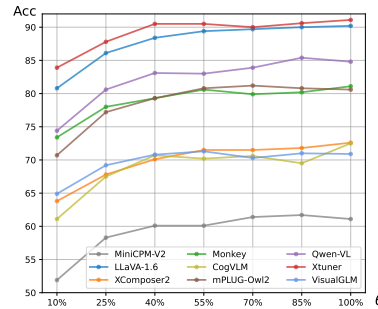


Figure 6: Accuracy under the content filter processing different percentages of website content.

328 6 Limitations

329 In our implementation, the hierarchical filtering model is trained apart from the LVLms, which may
 330 affect the performance considering the differences in training data, architecture design, and abilities
 331 of LVLms. In the future, we will consider training the filter and the LVLms in an end-to-end manner.

332 7 Conclusion

333 In this work, we have presented UDKAG, a plug-and-play framework to augment LVLms in handling
 334 visual question answering about up-to-date knowledge. By introducing a hierarchical filtering model,
 335 the framework enables LVLms to access up-to-date knowledge. A UDK-VQA dataset is further
 336 curated by scraping up-to-date news and generating news-related VQA samples. The dataset enables
 337 quantitatively evaluate the ability of LVLms to respond to questions about up-to-date knowledge.
 338 Experimental results on UDK-VQA demonstrate that our framework can significantly boost the
 339 performance of LVLms for answering questions requiring up-to-date knowledge.

References

- 340
- 341 [1] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar,
342 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early
343 experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 344 [2] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
345 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
346 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 347 [3] Grok Contributors. Grok. <https://github.com/xai-org/grok-1>, 2024.
- 348 [4] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next:
349 Improved reasoning, ocr, and world knowledge, January 2024.
- 350 [5] Claude Contributors. Claude. <https://claude.ai/>, 2024.
- 351 [6] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented
352 language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR,
353 2020.
- 354 [7] Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Heung-Yeung
355 Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of large language model with
356 knowledge graph. *arXiv preprint arXiv:2307.07697*, 2023.
- 357 [8] Angeliki Lazaridou, Elena Gribovskaya, Wojciech Stokowiec, and Nikolai Grigorev. Internet-augmented
358 language models through few-shot prompting for open-domain question answering. *arXiv preprint*
359 *arXiv:2203.05115*, 2022.
- 360 [9] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve,
361 generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- 362 [10] Ziniu Hu, Ahmet Iscen, Chen Sun, Zirui Wang, Kai-Wei Chang, Yizhou Sun, Cordelia Schmid, David A
363 Ross, and Alireza Fathi. Reveal: Retrieval-augmented visual-language pre-training with multi-source
364 multimodal knowledge memory. In *Proceedings of the IEEE/CVF conference on computer vision and*
365 *pattern recognition*, pages 23369–23379, 2023.
- 366 [11] Zhuolin Yang, Wei Ping, Zihan Liu, Vijay Korthikanti, Weili Nie, De-An Huang, Linxi Fan, Zhiding Yu,
367 Shiyi Lan, Bo Li, et al. Re-vilm: Retrieval-augmented visual language model for zero and few-shot image
368 captioning. *arXiv preprint arXiv:2302.04858*, 2023.
- 369 [12] Chenchen Jing, Yukun Li, Hao Chen, and Chunhua Shen. Retrieval-augmented primitive representations
370 for compositional zero-shot learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
371 volume 38, pages 2652–2660, 2024.
- 372 [13] Mojtaba Komeili, Kurt Shuster, and Jason Weston. Internet-augmented dialogue generation. *arXiv preprint*
373 *arXiv:2107.07566*, 2021.
- 374 [14] Junfeng Tian, Hehong Chen, Guohai Xu, Ming Yan, Xing Gao, Jianhai Zhang, Chenliang Li, Jiayi Liu,
375 Wenshen Xu, Haiyang Xu, et al. Chatplug: Open-domain generative dialogue system with internet-
376 augmented instruction tuning for digital human. *arXiv preprint arXiv:2304.07849*, 2023.
- 377 [15] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc,
378 Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for
379 few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- 380 [16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
381 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014:*
382 *13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages
383 740–755. Springer, 2014.
- 384 [17] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,
385 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual*
386 *Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565,
387 2018.
- 388 [18] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale
389 image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF conference*
390 *on computer vision and pattern recognition*, pages 3558–3568, 2021.

- 391 [19] Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Richard James, Jure Leskovec, Percy Liang, Mike
392 Lewis, Luke Zettlemoyer, and Wen-Tau Yih. Retrieval-augmented multimodal language modeling. In
393 *International Conference on Machine Learning*, pages 39755–39769. PMLR, 2023.
- 394 [20] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without
395 training. In *CVPR*, pages 14953–14962, 2023.
- 396 [21] Dídac Surfís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for
397 reasoning. In *ICCV*, pages 11888–11898, 2023.
- 398 [22] Difei Gao, Lei Ji, Luowei Zhou, Kevin Qinghong Lin, Joya Chen, Zihan Fan, and Mike Zheng Shou.
399 Assistgpt: A general multi-modal assistant that can plan, execute, inspect, and learn. *arXiv preprint*
400 *arXiv:2306.08640*, 2023.
- 401 [23] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt:
402 Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- 403 [24] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and
404 Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. In *NeurIPS*,
405 pages 43447–43478, 2023.
- 406 [25] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu,
407 Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and
408 action. *arXiv preprint arXiv:2303.11381*, 2023.
- 409 [26] Zhi Gao, Yuntao Du, Xintong Zhang, Xiaojian Ma, Wenjuan Han, Song-Chun Zhu, and Qing Li. Clova: A
410 closed-loop visual assistant with tool usage and update. In *Proceedings of the IEEE/CVF conference on*
411 *computer vision and pattern recognition (CVPR)*, 2024.
- 412 [27] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
413 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from
414 natural language supervision. In *Proceedings of the International Conference on Machine Learning*
415 (*ICML*), pages 8748–8763, 2021.
- 416 [28] Anil K Jain and Richard C Dubes. *Algorithms for clustering data*. Prentice-Hall, Inc., 1988.
- 417 [29] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep
418 bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- 419 [30] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
420 tuning, 2023.
- 421 [31] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and
422 Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint*
423 *arXiv:2308.12966*, 2023.
- 424 [32] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
425 Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on*
426 *Learning Representations*, 2022.
- 427 [33] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang
428 Zhang, Haodong Duan, Maosong Cao, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Xinyue Zhang,
429 Wei Li, Jingwen Li, Kai Chen, Conghui He, Xingcheng Zhang, Yu Qiao, Dahua Lin, and Jiaqi Wang.
430 Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language
431 large model. *arXiv preprint arXiv:2401.16420*, 2024.
- 432 [34] Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang
433 Bai. Monkey: Image resolution and text label are important things for large multi-modal models. *arXiv*
434 *preprint arXiv:2311.06607*, 2023.
- 435 [35] Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,
436 Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint*
437 *arXiv:2311.03079*, 2023.
- 438 [36] OpenBMB. Minicpm-v. <https://github.com/OpenBMB/OmniLMM>, 2024.
- 439 [37] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and
440 Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration.
441 2023.

- 442 [38] DataCanvas Ltd. mmalaya. <https://github.com/DataCanvasIO/MMAlaya>, 2024.
- 443 [39] XTuner Contributors. Xtuner: A toolkit for efficiently fine-tuning llm. [https://github.com/InternLM/](https://github.com/InternLM/xtuner)
444 [xtuner](https://github.com/InternLM/xtuner), 2023.
- 445 [40] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm: General
446 language model pretraining with autoregressive blank infilling. In *Proceedings of the Annual Meeting of*
447 *the Association for Computational Linguistics (ACL)*, pages 320–335, 2022.
- 448 [41] OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models.
449 <https://github.com/open-compass/opencompass>, 2023.
- 450 [42] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix,
451 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation
452 language models. *arXiv preprint arXiv:2302.13971*, 2023.

453 **NeurIPS Paper Checklist**

454 **1. Claims**

455 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's
456 contributions and scope?

457 Answer: [Yes]

458 Justification: The abstract contains our main claims including motivation, the up-to-date knowledge
459 retrieval-augmented framework, the pipeline for generating news-related VQA samples and the
460 curated dataset UDK-VQA.

461 Guidelines:

- 462 • The answer NA means that the abstract and introduction do not include the claims made in the
463 paper.
- 464 • The abstract and/or introduction should clearly state the claims made, including the contributions
465 made in the paper and important assumptions and limitations. A No or NA answer to this
466 question will not be perceived well by the reviewers.
- 467 • The claims made should match theoretical and experimental results, and reflect how much the
468 results can be expected to generalize to other settings.
- 469 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not
470 attained by the paper.

471 **2. Limitations**

472 Question: Does the paper discuss the limitations of the work performed by the authors?

473 Answer: [Yes]

474 Justification: We have discussed the limitations in a separate section.

475 Guidelines:

- 476 • The answer NA means that the paper has no limitation while the answer No means that the paper
477 has limitations, but those are not discussed in the paper.
- 478 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 479 • The paper should point out any strong assumptions and how robust the results are to violations of
480 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,
481 asymptotic approximations only holding locally). The authors should reflect on how these
482 assumptions might be violated in practice and what the implications would be.
- 483 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested
484 on a few datasets or with a few runs. In general, empirical results often depend on implicit
485 assumptions, which should be articulated.
- 486 • The authors should reflect on the factors that influence the performance of the approach. For
487 example, a facial recognition algorithm may perform poorly when image resolution is low or
488 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide
489 closed captions for online lectures because it fails to handle technical jargon.
- 490 • The authors should discuss the computational efficiency of the proposed algorithms and how
491 they scale with dataset size.
- 492 • If applicable, the authors should discuss possible limitations of their approach to address problems
493 of privacy and fairness.
- 494 • While the authors might fear that complete honesty about limitations might be used by reviewers
495 as grounds for rejection, a worse outcome might be that reviewers discover limitations that
496 aren't acknowledged in the paper. The authors should use their best judgment and recognize
497 that individual actions in favor of transparency play an important role in developing norms that
498 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize
499 honesty concerning limitations.

500 **3. Theory Assumptions and Proofs**

501 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete
502 (and correct) proof?

503 Answer: [NA]

504 Justification: Our work is not related to theorems.

505 Guidelines:

- 506 • The answer NA means that the paper does not include theoretical results.
- 507 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.

- 508 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 509 • The proofs can either appear in the main paper or the supplemental material, but if they appear in
- 510 the supplemental material, the authors are encouraged to provide a short proof sketch to provide
- 511 intuition.
- 512 • Inversely, any informal proof provided in the core of the paper should be complemented by
- 513 formal proofs provided in appendix or supplemental material.
- 514 • Theorems and Lemmas that the proof relies upon should be properly referenced.

515 4. Experimental Result Reproducibility

516 Question: Does the paper fully disclose all the information needed to reproduce the main experimental
 517 results of the paper to the extent that it affects the main claims and/or conclusions of the paper
 518 (regardless of whether the code and data are provided or not)?

519 Answer: [Yes]

520 Justification: We provide the implementation details including hyperparameter settings, baseline
 521 selection and evaluation details.

522 Guidelines:

- 523 • The answer NA means that the paper does not include experiments.
- 524 • If the paper includes experiments, a No answer to this question will not be perceived well by the
 525 reviewers: Making the paper reproducible is important, regardless of whether the code and data
 526 are provided or not.
- 527 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
 528 their results reproducible or verifiable.
- 529 • Depending on the contribution, reproducibility can be accomplished in various ways. For
 530 example, if the contribution is a novel architecture, describing the architecture fully might suffice,
 531 or if the contribution is a specific model and empirical evaluation, it may be necessary to either
 532 make it possible for others to replicate the model with the same dataset, or provide access to
 533 the model. In general, releasing code and data is often one good way to accomplish this, but
 534 reproducibility can also be provided via detailed instructions for how to replicate the results,
 535 access to a hosted model (e.g., in the case of a large language model), releasing of a model
 536 checkpoint, or other means that are appropriate to the research performed.
- 537 • While NeurIPS does not require releasing code, the conference does require all submissions
 538 to provide some reasonable avenue for reproducibility, which may depend on the nature of the
 539 contribution. For example
 - 540 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
 541 reproduce that algorithm.
 - 542 (b) If the contribution is primarily a new model architecture, the paper should describe the
 543 architecture clearly and fully.
 - 544 (c) If the contribution is a new model (e.g., a large language model), then there should either be
 545 a way to access this model for reproducing the results or a way to reproduce the model (e.g.,
 546 with an open-source dataset or instructions for how to construct the dataset).
 - 547 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are
 548 welcome to describe the particular way they provide for reproducibility. In the case of
 549 closed-source models, it may be that access to the model is limited in some way (e.g.,
 550 to registered users), but it should be possible for other researchers to have some path to
 551 reproducing or verifying the results.

552 5. Open access to data and code

553 Question: Does the paper provide open access to the data and code, with sufficient instructions to
 554 faithfully reproduce the main experimental results, as described in supplemental material?

555 Answer: [No]

556 Justification: We do not provide open access to the data and code at this time, but can publish part of
 557 them at the rebuttal stage if the reviewers need it. The complete data and code will be published after
 558 the paper is accepted.

559 Guidelines:

- 560 • The answer NA means that paper does not include experiments requiring code.
- 561 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/public/
 562 guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 563 • While we encourage the release of code and data, we understand that this might not be possible,
 564 so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless
 565 this is central to the contribution (e.g., for a new open-source benchmark).

- 566 • The instructions should contain the exact command and environment needed to run to reproduce
567 the results. See the NeurIPS code and data submission guidelines ([https://nips.cc/public/
568 guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 569 • The authors should provide instructions on data access and preparation, including how to access
570 the raw data, preprocessed data, intermediate data, and generated data, etc.
- 571 • The authors should provide scripts to reproduce all experimental results for the new proposed
572 method and baselines. If only a subset of experiments are reproducible, they should state which
573 ones are omitted from the script and why.
- 574 • At submission time, to preserve anonymity, the authors should release anonymized versions (if
575 applicable).
- 576 • Providing as much information as possible in supplemental material (appended to the paper) is
577 recommended, but including URLs to data and code is permitted.

578 6. Experimental Setting/Details

579 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,
580 how they were chosen, type of optimizer, etc.) necessary to understand the results?

581 Answer: [Yes]

582 Justification: We provide the implementation details including hyperparameter settings, baseline
583 selection and evaluation details.

584 Guidelines:

- 585 • The answer NA means that the paper does not include experiments.
- 586 • The experimental setting should be presented in the core of the paper to a level of detail that is
587 necessary to appreciate the results and make sense of them.
- 588 • The full details can be provided either with the code, in appendix, or as supplemental material.

589 7. Experiment Statistical Significance

590 Question: Does the paper report error bars suitably and correctly defined or other appropriate informa-
591 tion about the statistical significance of the experiments?

592 Answer: [No]

593 Justification: We follow existing work in the areas we work in and do not provide statistical significance
594 for fair comparisons.

595 Guidelines:

- 596 • The answer NA means that the paper does not include experiments.
- 597 • The authors should answer "Yes" if the results are accompanied by error bars, confidence
598 intervals, or statistical significance tests, at least for the experiments that support the main claims
599 of the paper.
- 600 • The factors of variability that the error bars are capturing should be clearly stated (for example,
601 train/test split, initialization, random drawing of some parameter, or overall run with given
602 experimental conditions).
- 603 • The method for calculating the error bars should be explained (closed form formula, call to a
604 library function, bootstrap, etc.)
- 605 • The assumptions made should be given (e.g., Normally distributed errors).
- 606 • It should be clear whether the error bar is the standard deviation or the standard error of the
607 mean.
- 608 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report
609 a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is
610 not verified.
- 611 • For asymmetric distributions, the authors should be careful not to show in tables or figures
612 symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 613 • If error bars are reported in tables or plots, The authors should explain in the text how they were
614 calculated and reference the corresponding figures or tables in the text.

615 8. Experiments Compute Resources

616 Question: For each experiment, does the paper provide sufficient information on the computer
617 resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

618 Answer: [Yes]

619 Justification: We provide the computer resources for reproducing the experiments.

620 Guidelines:

- 621 • The answer NA means that the paper does not include experiments.
- 622 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud
- 623 provider, including relevant memory and storage.
- 624 • The paper should provide the amount of compute required for each of the individual experimental
- 625 runs as well as estimate the total compute.
- 626 • The paper should disclose whether the full research project required more compute than the
- 627 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into
- 628 the paper).

629 9. Code Of Ethics

630 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code
631 of Ethics <https://neurips.cc/public/EthicsGuidelines>?

632 Answer: [Yes]

633 Justification: Our work conforms with the NeurIPS Code of Ethics.

634 Guidelines:

- 635 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 636 • If the authors answer No, they should explain the special circumstances that require a deviation
- 637 from the Code of Ethics.
- 638 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due
- 639 to laws or regulations in their jurisdiction).

640 10. Broader Impacts

641 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts
642 of the work performed?

643 Answer: [NA]

644 Justification: There is no societal impact of our work performed.

645 Guidelines:

- 646 • The answer NA means that there is no societal impact of the work performed.
- 647 • If the authors answer NA or No, they should explain why their work has no societal impact or
- 648 why the paper does not address societal impact.
- 649 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,
- 650 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-
- 651 ment of technologies that could make decisions that unfairly impact specific groups), privacy
- 652 considerations, and security considerations.
- 653 • The conference expects that many papers will be foundational research and not tied to particular
- 654 applications, let alone deployments. However, if there is a direct path to any negative applications,
- 655 the authors should point it out. For example, it is legitimate to point out that an improvement in
- 656 the quality of generative models could be used to generate deepfakes for disinformation. On the
- 657 other hand, it is not needed to point out that a generic algorithm for optimizing neural networks
- 658 could enable people to train models that generate Deepfakes faster.
- 659 • The authors should consider possible harms that could arise when the technology is being used
- 660 as intended and functioning correctly, harms that could arise when the technology is being used
- 661 as intended but gives incorrect results, and harms following from (intentional or unintentional)
- 662 misuse of the technology.
- 663 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies
- 664 (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitor-
- 665 ing misuse, mechanisms to monitor how a system learns from feedback over time, improving the
- 666 efficiency and accessibility of ML).

667 11. Safeguards

668 Question: Does the paper describe safeguards that have been put in place for responsible release of
669 data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or
670 scraped datasets)?

671 Answer: [Yes]

672 Justification: We use search engines to access Internet data, and search engines have their own methods
673 to avoid security safety risks. Moreover, samples in the test set we curated have been reviewed case by
674 case.

675 Guidelines:

- 676 • The answer NA means that the paper poses no such risks.

- 677 • Released models that have a high risk for misuse or dual-use should be released with necessary
678 safeguards to allow for controlled use of the model, for example by requiring that users adhere to
679 usage guidelines or restrictions to access the model or implementing safety filters.
- 680 • Datasets that have been scraped from the Internet could pose safety risks. The authors should
681 describe how they avoided releasing unsafe images.
- 682 • We recognize that providing effective safeguards is challenging, and many papers do not require
683 this, but we encourage authors to take this into account and make a best faith effort.

684 12. Licenses for existing assets

685 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,
686 properly credited and are the license and terms of use explicitly mentioned and properly respected?

687 Answer: [Yes]

688 Justification: We've cited the original paper of the code and model we used.

689 Guidelines:

- 690 • The answer NA means that the paper does not use existing assets.
- 691 • The authors should cite the original paper that produced the code package or dataset.
- 692 • The authors should state which version of the asset is used and, if possible, include a URL.
- 693 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 694 • For scraped data from a particular source (e.g., website), the copyright and terms of service of
695 that source should be provided.
- 696 • If assets are released, the license, copyright information, and terms of use in the package should
697 be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for
698 some datasets. Their licensing guide can help determine the license of a dataset.
- 699 • For existing datasets that are re-packaged, both the original license and the license of the derived
700 asset (if it has changed) should be provided.
- 701 • If this information is not available online, the authors are encouraged to reach out to the asset's
702 creators.

703 13. New Assets

704 Question: Are new assets introduced in the paper well documented and is the documentation provided
705 alongside the assets?

706 Answer: [No]

707 Justification: We will provide open access to part of the new assets at the rebuttal stage if the reviewers
708 need it. The complete assets will be published after the paper is accepted.

709 Guidelines:

- 710 • The answer NA means that the paper does not release new assets.
- 711 • Researchers should communicate the details of the dataset/code/model as part of their sub-
712 missions via structured templates. This includes details about training, license, limitations,
713 etc.
- 714 • The paper should discuss whether and how consent was obtained from people whose asset is
715 used.
- 716 • At submission time, remember to anonymize your assets (if applicable). You can either create an
717 anonymized URL or include an anonymized zip file.

718 14. Crowdsourcing and Research with Human Subjects

719 Question: For crowdsourcing experiments and research with human subjects, does the paper include
720 the full text of instructions given to participants and screenshots, if applicable, as well as details about
721 compensation (if any)?

722 Answer: [NA]

723 Justification: The test samples of our curated UDK-VQA dataset are checked by co-authors.

724 Guidelines:

- 725 • The answer NA means that the paper does not involve crowdsourcing nor research with human
726 subjects.
- 727 • Including this information in the supplemental material is fine, but if the main contribution of the
728 paper involves human subjects, then as much detail as possible should be included in the main
729 paper.
- 730 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other
731 labor should be paid at least the minimum wage in the country of the data collector.

732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The test samples of our curated UDK-VQA dataset are checked by co-authors.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.