# Harnessing LLMs for VQA: A Prompted Benchmark with Animate/Inanimate Keywords

Chanwoo Lee
*dept. of Software Convergence*
*Seoul Woman's University*
Seoul, South Korea
cksdn7479@gmail.com

Hyunjeong Lee
*dept. of Computer Engineering*
*Hanbat National University*
Daejeon, South Korea
20222039@edu.hanbat.ac.kr

Minsang Kim
*dept. of Computer Science and Engineering*
*Korea University*
Seoul, South Korea
kmswin1@korea.ac.kr

Hyun Kim
*Superintelligence Creative Research Lab.*
*Electronics and Telecommunications Research Institute*
Daejeon, South Korea
h.kim@etri.re.kr

Haneol Jang
*dept. of Computer Engineering*
*Hanbat National University*
Daejeon, South Korea
hejang@hanbat.ac.kr

Cheoneum Park*
*dept. of Computer Engineering*
*Hanbat National University*
Daejeon, South Korea
parkce@hanbat.ac.kr

*Abstract*—In the field of NLP, Large Language Models (LLMs) have recently achieved significant advancements, leading to the development of various benchmarks for their evaluation. Alongside NLP, Vision Language Models (VLMs) have also VLM have also significantly progressed, similar to LLMs. However, benchmarks for VLMs are still relatively underdeveloped compared to those for NLP, and their construction is often costly. In this work, we propose an automatically generated benchmark for evaluating VLMs based on LLMs and conduct a visual question answering task to assess this benchmark. The benchmark includes multiple-choice questions that not only distinguish between animate and inanimate objects but also generate these distinctions automatically, along with entity and object information within images. We evaluate the performance of open VLM using the generated multiple-choice questions, demonstrating the model's capabilities and the significance of the automatically generated benchmark. Finally, we discuss the necessity and future directions for benchmark research in this area.

*Index Terms*—Natural Language Processing, Prompt generation, Visual Question Answering benchmark, Multimodal synthetic data, Entity reasoning

## I. INTRODUCTION

Large language models (LLMs) are increasingly trained on vast amounts of high-quality raw text and instructions, leading to the generation of superior text outputs. According to the scaling law [1], as the number of parameters in these models increases to a large scale, the lower the test loss converges, and thus the quality of the generated text may improve. Consequently, the problem-solving capabilities of LLMs have also seen significant enhancement, prompting the development of various benchmarks to assess these capabilities. Problem-solving with LLMs extends beyond merely understanding context and predicting the next token [2]; it also encompasses reasoning tasks such as causal inference [3]–[5], domain- and language-specific knowledge [6]–[10], mathematics [5], dialogue [11], task decomposition [12], and planning [13].

Among these, [10] introduces auto annotation, which allows for periodic updates to the benchmark, thereby reducing the cost on human resources.

These benchmarks are crucial for evaluating and comparing the problem-solving abilities of various open LLMs. Traditionally, LLM evaluation has predominantly relied on text-based datasets. However, human problem-solving does not solely depend on textual information; it involves synthesizing various modalities such as visual and auditory data to draw conclusions. Consequently, there has been a growing interest in evaluating vision-language models (VLMs) trained on multimodal data [14]. Multimodal benchmarks focus on assessing an LLM's ability to process and solve problems using information that spans multiple modalities, with tasks commonly involving Visual Question Answering (VQA) [15]–[17], image captioning [18]–[21], and generating localized narratives [22], [23]. These benchmarks present more complex challenges compared to text-based benchmarks, as they allow for evaluating how LLMs comprehend visual information and integrate it with text to solve problems.

In this study, we generate VQA multiple choice questions based on a given image as part of a multimodal benchmark. The VQA multiple choice question format requires the model to understand the given image and question, and select the most appropriate answer from a set of choices. Unlike open-ended text generation, this format clearly defines the problem and assesses whether the model can accurately distinguish the correct answer within a specific context. A critical aspect of performing VQA is distinguishing between animate and inanimate objects in the image, as this plays a vital role in the model's problem-solving ability. For example, questions might involve identifying which objects in the image can move or understanding the relationships between objects. When animate and inanimate information is utilized in VQA multiple choice questions, it aids the model in achieving a deeper understanding and more accurate inference of the visual

---

* Corresponding author.

information presented by the given question [24], [25].

We construct an automated benchmark through prompting, which involves two key steps: 1) generating animate and inanimate keywords from the given image, and 2) using this information contextually to generate VQA multiple choice questions. Finally, to verify the validity of our hypothesis, we conduct a comparative experiment by creating two types of prompts for solving VQA multiple-choice questions: one that uses animate and inanimate information as context to select the correct answer, and another that relies solely on the question without additional context.

## II. METHOD

We build a benchmark by prompting the LLM to generate VQA multiple choice questions and animate and inanimate keywords. The benchmark is entirely composed of AI-generated outputs, and we design prompts to ensure the creation of high-quality results.

### A. Problem Description

The method of generating VQA multiple choice question and animate and inanimate keywords from an LLM uses a simple prompt $p(i)$ = "Input image $i$,generate animate and inanimate keywords and multiple choice question." consisting of the input image $I$ and a generation request token. For the generation of the keywords and the multiple choice question, we use LLMs, denoted by $\mathcal{M}$, such as ChatGPT [26], by prompting without additional training. The method can be formulated as follows:

$$\hat{y} = \mathcal{M}(p(i)) \tag{1}$$

In the above equation, $p(.)$ is the prompting function, $i$ is the input image, $\mathcal{M}$ is the LLM, and the output result is $\hat{y}$.

The evaluation is conducted by extracting multiple-choice questions from the output $\hat{y}$. The model is executed according to the following equation, where the input consists of an image $i$ and a multiple-choice question $q$, and the answer $a$ corresponding to $q$ is generated as the output.

$$a = \mathcal{M}(p_{vqa}(d)) \\ d = concat(i, q) \tag{2}$$

### B. Prompt Template

We design two types of prompts: 1) The base prompt is constructed to generate VQA multiple choice questions and animate and inanimate keywords (Figure 1) and 2) contrastive prompt is a prompt utilizing the counterfactual summary (Figure 3).

We design three types of prompts: 1) we generate keywords that distinguish between animate and inanimate objects in the given image, and then create VQA multiple-choice questions based on the generated context (1). 2) we input the previously generated VQA validation data along with the corresponding paired images to generate results appropriate

for the prompt(2, 3). Here, {QUESTION} and {CHOICE[n]} represent the questions and multiple-choice options in the generated benchmark.

To create an evaluation set, we need to generate the following three types of questions and answers: {Animate, Inanimate, Multiple Choice Question}
[order]:
1. Generate three requirements for the given image.
2. Evaluate whether the generated questions align with the given image.
3. Evaluate whether the generated answers are correct based on the questions.
4. Evaluate whether the generated answers align with the given image.
5. Provide the question that is most similar to the Topic.
6. Write all text in English

[Format]:
[Animate]
List all animate entities such as $[w_1, w_2, ..., w_n]$. If not, just say 'None'.

[Inanimate]
List all inanimate objects such as $[w_1, w_2, ..., w_n]$. If not, just say 'None'.

[Multiple Choice Question]
(Q)  Write a question.
A)  Write an option.
B)  Write an option.
C)  Write an option.
D)  Write an option.
(A)  Write the correct answer.

Fig. 1.  Prompt template design for generating VQA and animate and inanimate keywords.

Generate appropriate answers by looking at the given image and questions. There are one question types: (Multiple Choice)

[Multiple Choice]
Look at the question and answer candidates below.
Write the correct answer down after (A) with alphabet number.
(Q) {QUESTION}
{CHOICE[0]}
{CHOICE[1]}
{CHOICE[2]}
{CHOICE[3]}
(A)

Fig. 2.  Prompt template base design for selecting VQA answers.

```
Generate appropriate answers by looking at the given
image and questions. There are three question types:
(Animate, Inanimate, Multiple Choice)

[Animate]
List all animate entities such as [w_1, w_2, ..., w_n]. If not,
just say 'None'.

[Inanimate]
List all inanimate objects such as [w_1, w_2, ..., w_n]. If not,
just say 'None'.

[Multiple Choice]
Look at the question and answer candidates below.
Choose the correct answer considering the generated
animate and inanimate information.
Write the correct answer down after (A) with alphabet
number.
(Q) {QUESTION}
{CHOICE[0]}
{CHOICE[1]}
{CHOICE[2]}
{CHOICE[3]}
(A)
```

Fig. 3. Prompt template design for selecting VQA answers involves generating animate and inanimate keywords.

TABLE I
DATA STATISTICS, MCQ IS MULTIPLE CHOICE QUESTION

| Language | Num. of Image | Num. of MCQ |
|----------|---------------|-------------|
| English  | 2186          | 2186        |
| Korean   |               | 2100        |

## III. EXPERIMENTS

### A. Datasets

We present the data statistics in Table I for images sampled from the MS COCO [18] validation set, ensuring there is no overlap. The generated VQA dataset is constructed as bilingual, including both Korean and English.

We utilize Figure 1 to first prompt and generate VQA multiple-choice questions and animate/inanimate evaluation sets for each language. Based on the generated data, we follow the process outlined in Figure 3 to generate and evaluate the animate/inanimate multiple-choice results for the given evaluation images. We only evaluate the VQA multiple choice question and compute accuracy metrics separately for each the generated VQA dataset.

### B. Setups

a) Implementation details: For the experiments, we use GPT-4o (gpt-4o-mini) [26] with a temperature of 0.0 when calling the GPT-4o API.

TABLE II
MULTIPLE CHOICE RESULTS FOR THE VQA BENCHMARK GENERATED BY GPT-4O

| Language | w/ animate (acc) | w/o animate (acc) |
|----------|------------------|-------------------|
| English  | 96.47            | 94.64             |
| Korean   | 88.47            | 87.09             |

### C. Main Results

The animate and inanimate distinctions, as well as the VQA multiple-choice questions, are automatically generated by GPT-4o using specific prompts. In this paper, we evaluate the performance of the VQA multiple-choice questions based on the automatically generated data using GPT-4o-mini. The initial experimental results are presented in Table II. We analyze the experiment by comparing a method (w/ animate) that solves the VQA by prompting with a conjugation of animate and inanimate information and a method (w/o animate) that only processes the VQA.

The experimental results demonstrate significant performance overall, despite being based on benchmarks generated by LLM. In particular, a clear distinction is observed when using the animate information proposed in this study compared to when it is not used. This distinction is likely due to the method's ability to identify objects in the image and compare the semantic meaning of the VQA questions and candidates to find the correct answer. This feature is especially prominent in the Korean experiment, where the prompts are written in English, but the questions and candidates use Korean data.

In the case of the w/o animate prompt, the performance is relatively low at 87.09. However, when using both animate and inanimate information, the system appears to leverage semantic information sufficiently at the natural language level, enabling it to compare the relationships between the question, candidates, and image objects to find the correct answer. Therefore, it can be concluded that the use of animate and inanimate information in VQA, as hypothesized in this paper, is meaningful.

### D. Qualitative Analysis

a) Benchmark Sample: In this paper, the benchmark we generated uses GPT-4o and GPT-4o-mini to automatically create and evaluate questions through auto annotation. Although the automatically generated benchmark is somewhat less reliable compared to human-annotated data, advancements in LLM performance have demonstrated that auto annotation can produce high-quality benchmarks, as evidenced by studies such as [10]. Figure4 shows a VQA question generated using the prompts we employed. The generated animate, inanimate, and VQA multiple-choice questions demonstrate that the necessary information is accurately extracted from the images without hallucination, and that the VQA questions are well-constructed. Based on the previous experimental results, we confirm that it is possible to generate VQA questions using GPT-4o-mini and to provide appropriate answers to the automatically generated questions.

[Image]

[Animate]
[woman]

[Inanimate]
[phone, bag, wall, sweater]

[Multiple Choice Question]
(Q) What is the woman holding?
A) A book
B) A phone
C) A camera
D) A drink
(A) B) A phone

Fig. 4. Generated Benchmark Data Sample.

[Image]

[Animate]
[Player 1, Player 2, Player 3]

[Inanimate]
[Baseball Bat]

[Multiple Choice Question]
(Q) What is the primary object being held by one of the players?
A) Football
B) Baseball Bat
C) Tennis Racket
D) Golf Club
(A) B) Baseball Bat

Fig. 5. The Results Depending on The Use of Animate and Inanimate Information.

*b) Analysis with animate/inanimate information:* Moreover, we hypothesize that using object information, such as animate and inanimate attributes, is beneficial for reasoning about the meaning and relationships between image objects and textual entities in VQA tasks. In the case of Figure 5 performed in Figure 3, *Player 1, Player 2, Player 3* are generated as animate entities, and *Baseball bat* is generated as an inanimate entity. The question in the VQA includes the animate entity *players*, and one of the candidates includes the inanimate entity *Baseball bat*. The language model generates animate and inanimate entities based on the prompt, understands this information as context, and then performs the VQA task. The output for the VQA question is *B) Baseball bat*, which is the correct answer. On the other hand, when Figure 2 is executed, a different answer is generated. The animate and inanimate information generated in this process aids in solving the VQA task, thereby supporting our hypothesis. This suggests that accurately understanding the relationships between image objects and text plays a crucial role in deriving the correct answer.

## IV. CONCLUSION

In this paper, we propose a novel approach to automatically generating multimodal benchmark with LLM prompting. The proposed data is VQA multiple choice questions and animate/inainmate keywords for the input image. We conducted VQA multiple-choice experiments using the generated data, and the results showed significant performance, with accuracy ranging from approximately 87% to 88% for Korean and 94% to 96% for English. Additionally, we validated our hypothesis through experiments, demonstrating that generating animate and inanimate entities from the image and incorporating them into the prompt can aid the language model's reasoning when solving VQA multiple-choice problems.

For future work, we plan to conduct human annotation to verify the quality of the generated questions. Additionally, we aim to develop methods to filter high-quality questions and enhance the diversity of the generated questions.

## REFERENCES

[1] J. Kaplan et al. 2020. Scaling laws for neural language models. [Online]. Available: https://arxiv.org/abs/2001.08361
[2] G. Son et al., "KMMLU: measuring massive multitask language understanding in Korean", arXiv preprint arXiv:2402.11548, 2024. [Online]. Available: https://arxiv.org/abs/2402.11548.
[3] J. Ham, Y. J. Choe, K. Park, I. Choi, and H. Soh, "KorNLI and KorSTS: new benchmark datasets for Korean natural language understanding," T. Cohn, Y. He, and Y. Liu, Eds., Association for Computational Linguistics, Nov. 2020, in *Find. Assoc. Comput.*

*Linguistics: EMNLP 2020*, pp. 422–430. [Online]. Available: https://aclanthology.org/2020.findings-emnlp.39.

[4] J. Jin, J. Kim, N. Lee, H. Yoo, A. Oh, and H. Lee, "KoBBQ: Korean bias benchmark for question answering," *Trans. Assoc. Comput. Linguistics*, vol. 12, pp. 507–524, May 2024. [Online]. Available: https://doi.org/10.1162/tacl_a_00661.

[5] D. Hendrycks et al., "Measuring massive multitask language understanding," in *Proc. Int. Conf. Learn. Represent.*, 2021. [Online]. Available: https://openreview.net/forum?id=d7KBjmI3GmQ.

[6] E. Kim, J. Suk, P. Oh, H. Yoo, J. Thorne, and A. Oh, "CLIcK: A benchmark dataset of cultural and linguistic intelligence in Korean," N. Calzolari, M.-Y. Kan, V. Hoste, A. Lenci, S. Sakti, and N. Xue, Eds., Turin, Italy: ELRA and ICCL, May 2024, in *Proc. 2024 Joint Int. Conf. Comput. Linguistics, Lang. Resources Eval. (LREC-COLING 2024)*, pp. 3335–3346. [Online]. Available: https://aclanthology.org/2024.lrec-main.296.

[7] W. Hwang, D. Lee, K. Cho, H. Lee, and M. Seo, "A multi-task benchmark for Korean legal language understanding and judgement prediction," in *Proc. Thirty-Sixth Conf. Neural Inf. Process. Syst. Datasets and Benchmarks Track*, 2022.

[8] J. Lee, M. Kim, S. Kim, J. Kim, S. Won, H. Lee, and E. Choi, "KorNAT: LLM alignment benchmark for Korean social Values and common knowledge", arXiv preprint arXiv:2402.13605, 2024. [Online]. Available: https://arxiv.org/abs/2402.13605.

[9] Y. Huang et al., "C-EVAL: a multi-level multi-discipline Chinese evaluation suite for foundation models," in *Proc. 37th Int. Conf. Neural Inf. Process. Syst.* Red Hook, NY, USA: Curran Associates Inc., 2024.

[10] Z. Gu et al., "Xiezhi: an ever-updating benchmark for holistic domain knowledge evaluation," *Proc. AAAI Conf. Artif. Intell.*, vol. 38, no. 16, pp. 18 099–18 107, Mar 2024. [Online]. Available: https://ojs.aaai.org/index.php/AAAI/article/view/29767.

[11] L. Zheng et al., "Judging LLM-as-a-Judge with MT-Bench and chatbot arena," in *Proc. 37th Int. Conf. Neural Inf. Process. Syst.* Red Hook, NY, USA: Curran Associates Inc., 2024.

[12] Y. Shen et al., "TaskBench: benchmarking large language models for task automation", arXiv preprint arXiv:2311.18760, 2023. [Online]. Available: https://arxiv.org/abs/2311.18760.

[13] K. Valmeekam, M. Marquez, A. Olmo, S. Sreedharan, and S. Kambhampati, "Planbench: an extensible benchmark for evaluating large language models on planning and reasoning about change," in *Proc. 37th Int. Conf. Neural Inf. Process. Syst.* Red Hook, NY, USA: Curran Associates Inc., 2024.

[14] K. Ying et al., "MMT-Bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask AGI", arXiv preprint arXiv:2404.16006, 2024. [Online]. Available: https://arxiv.org/abs/2404.16006.

[15] S. Antol et al., "VQA: Visual Question Answering," in *Proc. Int. Conf. Comput. Vis. (ICCV)*, Santiago, Chile, 2015.

[16] X. Fu et al., "Generate then select: Open-ended visual question answering guided by world knowledge," A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds., Toronto, Canada, Jul. 2023, in *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 2333–2346. [Online]. Available: https://aclanthology.org/2023.findings-acl.147.

[17] S. Subramanian et al., "Modular visual question answering via code generation", arXiv preprint arXiv:2306.05392, 2023. [Online]. Available: https://arxiv.org/abs/2306.05392.

[18] T.-Y. Lin et al., "Microsoft COCO: common objects in context", arXiv preprint arXiv:1405.0312, 2015. [Online]. Available: https://arxiv.org/abs/1405.0312.

[19] R. Krishna et al., "Visual genome: connecting language and vision using crowdsourced dense image annotations," *Int. J. Comput. Vis.*, vol. 123, no. 1, pp. 32–73, 2017. [Online]. Available: https://doi.org/10.1007/s11263-016-0981-7.

[20] T. Thrush et al., "Winoground: Probing vision and language models for visio-linguistic compositionality," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2022, pp. 5238–5248.

[21] C. Schuhmann et al., "LAION-5B: An open large-scale dataset for training next generation image-text models," *Adv. NeurIPS.*, vol. 35, pp. 25 278–25 294, 2022.

[22] J. Pont-Tuset, J. Uijlings, S. Changpinyo, R. Soricut, and V. Ferrari, "Connecting vision and language with localized narratives," in *Proc. Eur. Conf. Comput. Vis. (ECCV).* Glasgow, UK: Springer, 2020.

[23] Z. Gan et al., "Vision-language pre-training: basics, recent advances, and future trends," *Found. Trends Comput. Graph. Vis.*, vol. 14, no. 3–4, pp. 163–352, 2022.

[24] B. Zhou, A. Khosla, A. Lapedriza, A. Torralba, and A. Oliva, "Places: An image database for deep scene understanding", arXiv preprint arXiv:1610.02055, 2016. [Online]. Available: https://arxiv.org/abs/1610.02055.

[25] D. Teney, P. Anderson, X. He, and A. van den Hengel, "Tips and tricks for visual question answering: Learnings from the 2017 challenge," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 4223–4232.

[26] OpenAI, "GPT-4 technical report", arXiv preprint arXiv:2303.08774, 2023. [Online]. Available: https://arxiv.org/abs/2303.08774.