# Prompting Large Language Models with Answer Heuristics for Knowledge-based Visual Question Answering

Zhenwei Shao[1]    Zhou Yu[1*]    Meng Wang[2]    Jun Yu[1]

[1]Key Laboratory of Complex Systems Modeling and Simulation,
School of Computer Science and Technology, Hangzhou Dianzi University, China.
[2]School of Computer Science and Information Engineering, Hefei University of Technology, China

## A. Broader Impact

From a broader "VL+LLM" point of view, Prophet connects a vision-language (VL) model and a frozen LLM, aiming to endow the VL model with knowledge reasoning ability. Compared with Flamingo [1] which learns a large *general* VL model on top of a frozen LLM (Chinchilla-70B) in an end-to-end manner, Prophet introduces a decoupled VL+LLM learning paradigm to learn a small *task-specific* VL model independently and then feed its predictions to any external LLM for refinement. Consequently, Prophet significantly reduces the computational costs of model training and also benefits from the flexible choices of LLMs (both the offline and online LLMs are supported). We hope the decoupled learning paradigm in Prophet will inspire future research using LLMs.

## B. Implementation Details

### B.1. Improved MCAN Model

**Model architecture.** We propose an improved version of MCAN [11] based on its open-sourced MCAN-large implementation. Our modifications to the model architecture include: (i) we replace the original bottom-up-attention features with the grid-based features extracted from the CLIP's visual encoder with RN50×64 backbone [6]; (ii) we introduce the RoPE mechanism [7] to each image self-attention layer of MCAN to supplement the grid-based features with positional information; and (iii) we replace the original LSTM network with a pre-trained BERT-large model [2] as the text encoder before MCAN. Table 1 shows the accuracies of different model variants on the testing set of OK-VQA. By progressively adding the modifications to the original MCAN model, our improved MCAN model reports a 53.0% accuracy, which is on par with current state-of-the-art methods like KAT [4].

**Training recipe.** We first pretrain the model on the augmented *train+val+vg* dataset from VQAv2 [3] and Visual

---

| case | OK-VQA accuracy |
|---|---|
| original MCAN | 43.6 |
| + CLIP visual feats | 49.6 |
| + RoPE mechanism | 50.3 |
| + BERT as the text encoder | 53.0 |

Table 1. **Ablations for model architectures.** '+' denotes each modification is appended to the variant on the previous row.

Genome [5], with excluding the samples whose images are used in the testing split of OK-VQA to avoid data contamination. The settings for the pretraining stage are identical to the original implementation of MCAN. After that, the model is finetuned on the downstream OK-VQA and A-OKVQA datasets, respectively. For finetuning, the commonly used strategy is to replace the last linear layer (*i.e.*, the classification layer) with a new layer to adapt to the answer vocabulary of the downstream dataset. However, the answer vocabularies of the pretraining and finetuning datasets are *partially* overlapped. To maximally utilize the pretrained model parameters in the last layer, we inherit the parameters of existing answers and append new parameters for the new answers. After that, we freeze all the pretrained parameters and only update the new parameters for one epoch as a warm-up, and then train all model parameters for the rest training epochs. The detailed settings for the finetuning stage are shown in Table 2.

Table 3 shows the effects of different training strategies. Even without finetuning, the pretrained model (b) is superior to the model trained from scratch (a), implying the importance of pretraining. Moreover, our new finetuning strategy (d) leads to significantly better performance than the commonly used strategy (c), showing the effectiveness of inheriting model parameters for existing answers.

### B.2. Prompt Formats

We show an exemplar of Prophet's prompt in Table 7 and an exemplar of Prophet-MC's prompt (only used in the multiple-choice task of A-OKVQA) in Table 8.

| config | setting |
|---|---|
| optimizer | AdamW |
| weight decay | 0.01 |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.98$ |
| batch size | 64 |
| warm-up learning rate | 1e-3 |
| warm-up strategy | only update new parameters |
| warm-up epochs | 1 |
| base learning rate | 5e-5 |
| learning rate schedule | step decay |
| learning rate decay rate | 0.2 |
| learning rate decay epoch | 6 |
| total training epochs | 6 |

Table 2. **Training settings.** These hyper-parameters are used in both the OK-VQA and A-OKVQA experiments.

| training strategy | OK-VQA accuracy |
|---|---|
| (a) train from scratch | 35.6 |
| (b) pretrain, w/o finetune | 41.1 |
| (c) w/ finetune, replace last layer | 47.7 |
| (d) w/ finetune, append new answers | **53.0** |

Table 3. **Ablations for training strategies.** Four variants use the same architecture (last row in Table 1) except for the last layer.

### B.3. Choice of Other Open-source LLMs

We have tried to replace GPT-3 with another open-source LLM GPT-J (6B) [8] and observed a distinct accuracy drop (53.6% *vs.* 61.1%). This can be explained that GPT-J's in-context learning ability is not strong enough to understand the answer candidates. As reported in [9], this advanced ability *emerges* only when a LLM reaches a certain model size (*e.g*, >100B parameters). As running such a LLM offline is far beyond the computational resource we have, using GPT-3 is our best choice at this moment. If a proper-size and open-source LLM with the similar in-context learning ability to GPT-3 is available in the future, we can use it to replace GPT-3 seamlessly.

### C. More Statistical Analyses

We provide more in-depth analyses of Prophet's performance on the testing set of OKVQA. All results are carried out using the default settings.

First, we show the per-type accuracies of MCAN (stage-1) and Prophet (stage-2) in Table 4. Prophet outperforms MCAN on all categories, indicating that generality of the knowledge in GPT-3. The improvement on the "Science and Technology" category is not as large as the rest categories. which can be explained that the required knowledge for this category is more specialized and professional. These questions are also challenging for humans.

Then, we display the distribution of four situations of Prophet's predictions before and after GPT-3 in Table 5.

| category | MCAN | Prophet |
|---|---|---|
| Plants and Animals | 52.58 | **63.67** |
| Science and Technology | 48.10 | **48.81** |
| Sports and Recreation | 59.08 | **66.00** |
| Geography, History, Language and Culture | 52.48 | **62.98** |
| Brands, Companies and Products | 51.98 | **54.77** |
| Vehicles and Transportation | 50.82 | **58.01** |
| Cooking and Food | 55.53 | **62.09** |
| Weather and Climate | 65.12 | **68.37** |
| People and Everyday life | 49.44 | **54.67** |
| Objects, Material and Clothing | 50.05 | **57.20** |

Table 4. **Per-category accuracies** of MCAN (stage-1) and Prophet (stage-2). This performance improvements of using GPT-3 are observed on all categories.

| before \ after | correct | wrong |
|---|---|---|
| correct | 54.4% | 4.2% |
| wrong | 12.0% | 29.4% |

Table 5. **The distribution of four situations** of Prophet's predictions before and after GPT-3. Prophet maintains the majority of correct predictions by MCAN, and the accuracy improvement by GPT-3 is mainly because the number of *wrong-to-correct* samples is larger than that of the *correct-to-wrong* samples.

| failure cause | proportion |
|---|---|
| (a) insufficient visual understanding | 27.3% |
| (b) incorrect knowledge reasoning | **44.1%** |
| (c) correct but differently expressed answer | 22.8% |
| (d) others | 5.8% |

Table 6. **The distribution of failure causes** by human studies.

From the results, we draw the following conclusions: (i) The proportion of the *correct-to-correct* samples (54.4%) is close to the accuracy of MCAN (53.0%) and proportion of the *correct-to-wrong* samples (4.2%) are relatively small. This means that Prophet maintains the majority of correct predictions; (ii) the accuracy improvement of Prophet is mainly due to the fact that the proportion of *wrong-to-correct* samples (12.4%) is larger than that of the *correct-to-wrong* samples (4.2%); (iii) there are still a considerable amount of samples (29.4%) that both MCAN and Prophet fail to give the correct answer, which leaves sufficient room for future improvement.

Finally, we perform human studies to analyze the causes of wrong predictions in Table 6. For each category, we randomly sample 10% testing samples that Prophet fails to get the correct answer. This results in 172 samples. We ask three annotators to categorize each sample into one of the following four failure causes: (a) insufficient visual understanding; (b) incorrect knowledge reasoning; (c) correct but differently expressed answer; (d) others (*e.g.*,

the failure is caused by the ambiguity of the question). From the results, we can see that the cause of "(b) incorrect knowledge reasoning" accounts for the highest proportion, which suggests that the bottleneck of Prophet still lies in the knowledge acquisition and reasoning. The cause of "(a) insufficient visual understanding" has the second highest proportion, showing the potential of devising more powerful VQA models. The cause of "(c) correct but differently expressed answer" also accounts for a considerable proportion. This reflects the limitation of the annotations and evaluation metric of OK-VQA.

## D. Case Study

Figure 1 provides some testing samples along with their in-context examples to illustrate how the answer heuristics work. The results show that both the answer candidates and the answer-aware examples are helpful for GPT-3 to generate high-quality answers. It is worth noting that even though some predicted answers do not hit the ground-truth answers, they are still reasonable under human evaluation.

Figure 2 demonstrates some testing samples from different knowledge categories. In the 1st-3rd columns, we show the correctly answered samples with different prediction behaviors (*i.e.*, keep top-1, in top 2-$K$, and beyond top-$K$). The visualized results indicate that Prophet can adaptively choose suitable answers from candidates. In the last column, we show some failure samples, implying that there is still room for future improvement.

## References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022. 1

[2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, pages 4171–4186, 2019. 1

[3] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *CVPR*, 2017. 1

[4] Liangke Gui, Borui Wang, Qiuyuan Huang, Alex Hauptmann, Yonatan Bisk, and Jianfeng Gao. Kat: A knowledge augmented transformer for vision-and-language. In *NAACL*, 2021. 1, 4

[5] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 123(1):32–73, 2017. 1

[6] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021. 1

[7] Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021. 1

[8] Ben Wang and Aran Komatsuzaki. Gpt-j-6b: A 6 billion parameter autoregressive language model, 2021. 2

[9] Jason Wei, Yi Tay, Rishi Bommasani, et al. Emergent abilities of large language models. *TMLR*, 2023. 2

[10] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. An empirical study of gpt-3 for few-shot knowledge-based vqa. In *AAAI*, pages 3081–3089, 2022. 4

[11] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *CVPR*, pages 6281–6290, 2019. 1

```
Please answer the question according to the context and the answer candidates. Each answer
candidate is associated with a confidence score within a bracket. The true answer may not be
included in the candidates.

===
Context: The motorcycle racers are getting ready for a race.
===
Question: What sport are these guys doing?
===
Candidates: motorcross(0.94), motocross(0.79), bike(0.35), dirt bike(0.28), motorcycle(0.03),
bmx(0.03), cycling(0.02), motorbike(0.02), race(0.02), bicycle(0.02)
===
Answer: motorcross

===
Context: A man riding on the back of a motorcycle.
===
Question: Can you name a sport this person could be a
part of?
===
Candidates: race(0.94), motocross(0.80), dirt bike(0.70), motorcross(0.25), motorcycle(0.08),
bike(0.03), cycling(0.03), motorbike(0.01), ride(0.01), bicycle(0.01)
===
Answer: race

===
Context: A man on a motorcycle flying through the air.
===
Question: Which sport is this?
===
Candidates: motorcross(0.91), motocross(0.91), dirt bike(0.25), bike(0.05), bmx(0.03),
motorbike(0.01), motorcycle(0.01), bicycling(0.01), cycling(0.01), dirt(0.01)
===
Answer: motocross

===
Context: a black motorcycle parked in a parking lot.
===
Question: What sport can you use this for?
===
Candidates: race(0.53), motorcycle(0.41), motocross(0.19), bike(0.17), motorcross(0.15),
cycling(0.11), dirt bike(0.10), ride(0.08), bicycling(0.01), bicycle(0.01)
===
Answer:
```

Table 7. **A prompt exemplar for Prophet.** We only show 3 in-context examples here for better visualization. Following the implementations in PICa [10] and KAT [4], we use a special symbol '===' to separate each two lines.

```
Please choose the correct answer in the choices according to the context, the question and the
answer candidates. Each answer candidate is associated with a confidence score within a bracket.
The true answer may not be included in the candidates.

===
Context: A young man riding a skateboard on a sidewalk.
===
Question: What part of his body will be most harmed by the item in his mouth?
===
Candidates: skateboard(0.02), nothing(0.02), table(0.01), leg(0.01), helmet(0.00), knees(0.00),
skateboarding(0.00), head(0.00), teeth(0.00), falling(0.00)
===
Choices: (A) back, (B) lungs, (C) feet, (D) eyes
===
Answer: (B)

===
Context: A guy jumping in the air on a skateboard.
===
Question: What is touching the skateboard?
===
Candidates: shoe(0.89), foot(0.67), shoes(0.51), feet(0.45), sneakers(0.17), boy(0.03),
road(0.03), man(0.03), skateboard(0.02), nothing(0.02)
===
Choices: (A) sneakers, (B) dress shoes, (C) pogo stick, (D) hands
===
Answer: (A)

===
Context: A man flying through the air over a skateboard.
===
Question: Before going aloft what did the man ride?
===
Candidates: skateboard(1.00), skating(0.03), skate board(0.01), skate(0.01), board(0.01),
skateboarding(0.00), scooter(0.00), car(0.00), trick(0.00), longboard(0.00)
===
Choices: (A) unicycle, (B) skateboard, (C) plane, (D) car
===
Answer: (B)

===
Context: a young boy kneeling on a skateboard on the street.
===
Question: What did this lad likely injure here?
===
Candidates: skateboard(0.18), shoes(0.02), shoe(0.02), skateboarding(0.01), street(0.01),
flowers(0.01), skating(0.01), boy(0.01), head(0.00), skateboarder(0.00)
===
Choices: (A) knee, (B) elbow, (C) rear, (D) board
===
Answer:
```

Table 8. **A prompt exemplar for Prophet-MC.** Compared to the prompt in Table 7, we add one extra line of choices for the example and testing input, and change the output format to adapt to the multiple-choice task. All the differences are marked in red.

Figure 1. **Answer candidates and answer-aware examples.** We show some typical samples of the testing inputs and their in-context examples. The predicted answers of Prophet have a high probability to appear in the answer candidates and answer-aware examples, showing the effectiveness of answer heuristics in enhancing GPT-3's ability to predict the correct answer. Although the samples in the last column are not answered correctly, their predictions are still reasonable under human evaluation.
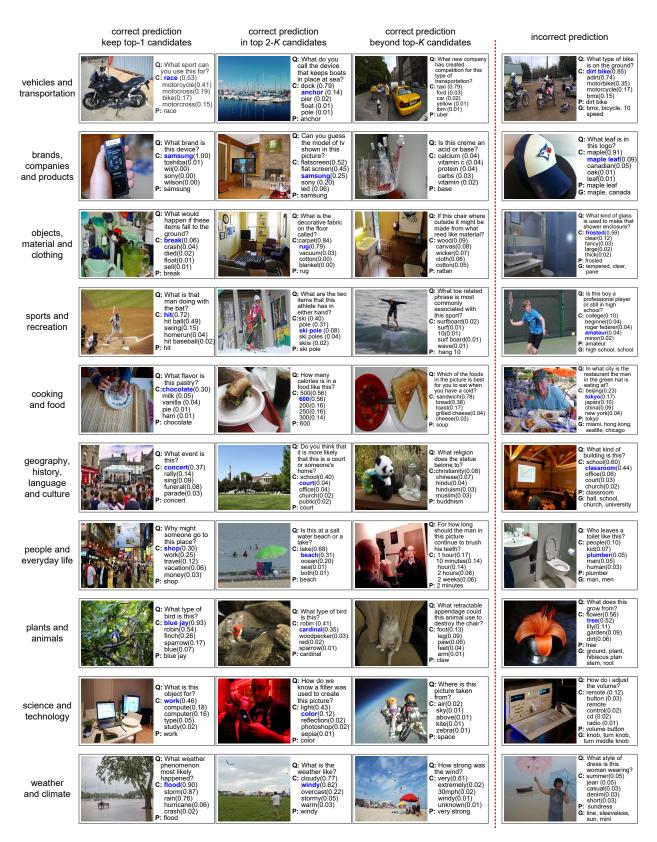
Figure 2. **Different categories and prediction behaviors.** Each row contains four testing samples from a specific knowledge category. The first to the third columns correspond to the correctly answered samples of different prediction behaviors (*i.e.*, keep top-1, in top 2-$K$, and beyond top-$K$). The last column contains failure samples.