

# Targeted Visual Prompting for Medical Visual Question Answering

Sergio Tascon-Morales ✉, Pablo Márquez-Neila, Raphael Sznitman

University of Bern, Bern, Switzerland

{sergio.tasconmorales, pablo.marquez, raphael.sznitman}@unibe.ch

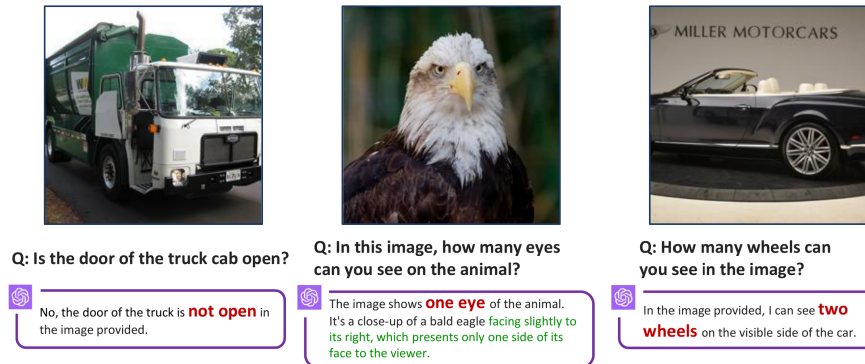
**Abstract.** With growing interest in recent years, medical visual question answering (Med-VQA) has rapidly evolved, with multimodal large language models (MLLMs) emerging as an alternative to classical model architectures. Specifically, their ability to add visual information to the input of pre-trained LLMs brings new capabilities for image interpretation. However, simple visual errors cast doubt on the actual visual understanding abilities of these models. To address this, region-based questions have been proposed as a means to assess and enhance actual visual understanding through compositional evaluation. To combine these two perspectives, this paper introduces targeted visual prompting to equip MLLMs with region-based questioning capabilities. By presenting the model with both the isolated region and the region in its context in a customized visual prompt, we show the effectiveness of our method across multiple datasets while comparing it to several baseline models. Our code and data are available at <https://github.com/sergiotasconmorales/locvqallm>.

**Keywords:** VQA · Localized Questions · Multimodal Large Language Model · Vision Transformer

## 1 Introduction

Visual Question Answering (VQA) is centered on developing models capable of answering questions about specific images [2]. This task is particularly challenging within the medical domain due to factors such as a scarcity of annotated data [11,7], the wide variety of imaging modalities and anatomical regions [5], as well as the unique characteristics of medical images and terminology, all of which necessitate specialized expertise [7,25]. Furthermore, approaches that leverage the detection of natural objects, which have significantly improved performance in the analysis of natural images [1], are less straightforward when applied to medical imagery [5].

Historically, models for Medical VQA (Med-VQA) treated visual and textual information independently, later merging these features through various fusion techniques. This composite data would then be input into a classifier to determine the most probable answer. However, recent developments in transformer-based models [21], including advancements in Large Language Models (LLMs),



**Fig. 1.** Examples of visual understanding failures using GPT-4V for the VQA task (Examples taken from [18]).

have led to a notable shift in VQA strategies. These advancements have paved the way for the adoption of multimodal LLMs (MLLMs) that integrate both visual and textual data more seamlessly, a trend that is emerging in both general [24,18,26] and Med-VQA [14,28] applications.

Despite the remarkable adoption of MLLMs, recent research has raised concerns about the quality of their visual capabilities (Fig. 1). This issue primarily arises from the pre-training process of the visual component, which typically relies on models like CLIP [12]. Surprisingly, MLLMs can perceive certain visually distinct images as similar, a phenomenon that human observers readily recognize as a visual error [18]. These visual understanding failures were also observed in VQA models before the widespread adoption of MLLMs [4,6,13,15].

To detect such failures and enhance explainability in the visual component of Med-VQA, the work in [17] proposes a novel approach using the formulation of *localized questions* [17]. These questions allow fine-grained probing of images by focusing on user-defined regions rather than the entire image and facilitate a *compositional evaluation*. To enable such localized questions, the region to query is encoded and directly integrated into the attention mechanism of the model. Other proposed strategies include providing the model with a restricted region of the image [16] or relying on the language component of the VQA model to interpret region coordinates directly included in the question [22]. Yet, due to their design focused on traditional architectures, these methods fail to benefit MLLMs in Med-VQA. Other traditional [10] and MLLM-based methods [3,27] rely on object detectors, limiting their applicability to medical images.

To overcome these challenges and enable localized questions in MLLMs in Med-VQA, we introduce *Targeted Visual Prompting*. By carefully designing a prompt that provides both global and local visual tokens relative to the region of interest defined by the user, our method allows the full advantage of the MLLM to enhance the performance of the VQA model. To validate the effectiveness of

our method, we conduct exhaustive experiments across multiple datasets. Our results demonstrate clear performance benefits compared to previously proposed methods, all achieved without introducing additional parameters to the model.

## 2 Method

A VQA model with parameters  $\theta$  generates an answer  $\hat{a}$  when given an input image  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$  and a related question represented as a sequence of words,  $\mathbf{q}$ . In its most general form, this process can be described as a function  $\Psi_{\text{VQA}}$ , parameterized by  $\theta$ , that is applied on the image-question pair,

$$\hat{a} = \Psi_{\text{VQA}}(\mathbf{x}, \mathbf{q}; \theta). \quad (1)$$

In practice, this model’s output has traditionally been a distribution over a set of  $N$  candidate answers  $\{a_1, a_2, \dots, a_N\}$  set beforehand.

In this work, however, we choose the answer of the VQA to be generated by an LLM in an auto-regressive manner until the end-of-sentence (EOS) token is produced. To make the LLM multimodal, we adopt the widely used approach of projecting visual embeddings onto the input space of the LLM [8,20,23] and express this as,

$$\hat{a} = \Psi_{\text{LLM}}(\Psi_{\text{Vis}}(\mathbf{x}, \theta_{\text{Vis}}) \mathbf{W}^{\text{proj}}, \mathbf{q}; \theta_{\text{LLM}}), \quad (2)$$

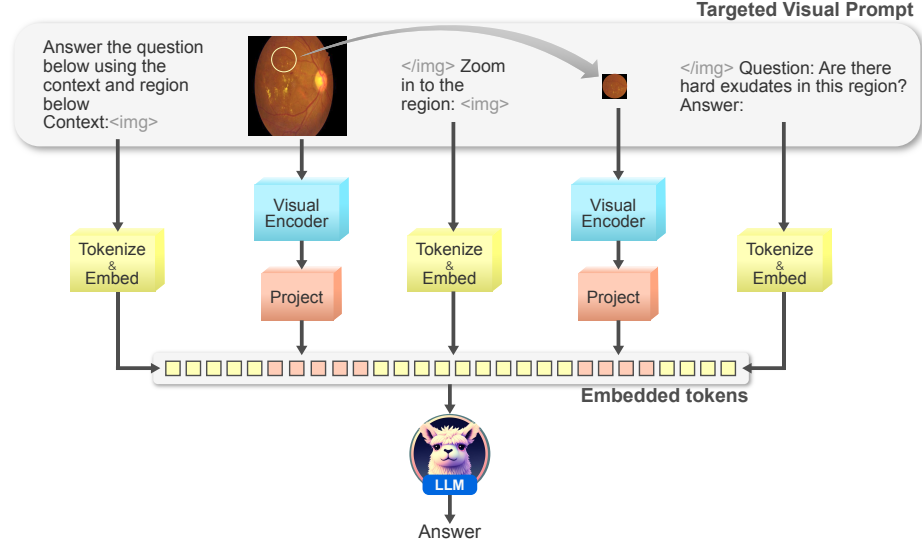
where  $\Psi_{\text{Vis}}$  refers to the visual encoder with parameters  $\theta_{\text{Vis}}$ , and  $\mathbf{W}^{\text{proj}}$  denotes the learnable parameters of the projection layer. Although not explicitly formalized, it is implied that the answer is generated in an auto-regressive fashion, meaning that the next word in the answer depends on the previously predicted words.

To expand the model’s capability to handle localized questions, we propose here a dedicated targeted visual prompt that allows two perspectives of the image to be encoded: one containing only the region of the image and the other containing the region in context.

The targeted visual prompt consists of five components: (1) comprises model instruction, denoted as  $\mathbf{w}_{\text{instr}}$ ; (2) the visual context represented by the image with the region drawn on it,  $\mathbf{x}_r$ ; (3)  $\mathbf{w}_{\text{det}}$  contains a textual prefix for the region; (4) the cropped region  $\mathbf{r}$ ; and (5)  $\mathbf{w}_q$  includes the question  $\mathbf{q}$ . Text-containing parts of the prompt undergo tokenization and embedding, while the visual components are processed by a visual encoder and then projected into the input space of the LLM. Subsequently, the results are concatenated and processed by the LLM, resulting in the generation of an answer. To handle global questions, the entire image is assigned to  $\mathbf{r}$ . We illustrate our model in Fig. 2 and summarize the computation of the answer as,

$$\hat{a} = \Psi_{\text{LLM}}(\mathbf{w}_{\text{instr}}, \Psi_{\text{Vis}}(\mathbf{x}_r, \theta_{\text{Vis}}) \mathbf{W}_{\mathbf{x}_r}^{\text{proj}}, \mathbf{w}_{\text{det}}, \Psi_{\text{Vis}}(\mathbf{r}, \theta_{\text{Vis}}) \mathbf{W}_{\mathbf{r}}^{\text{proj}}, \mathbf{w}_q; \theta_{\text{LLM}}). \quad (3)$$

To handle questions about the entire image, both  $\mathbf{x}_r$  and  $\mathbf{r}$  correspond to the original image.



**Fig. 2.** Our customized targeted visual prompt is created by providing the model with the region in context, as well as an isolated version of the region. Visual tokens are projected to the input space of the LLM and concatenated with the instruction tokens.

**Training.** As in [23], our model is trained using the original auto-regressive training loss of the LLM. The loss function is the standard negative log-likelihood accumulated over all time steps for predicting the correct next token. For a ground truth answer of length  $T$ , this loss is expressed as,

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \log p_{\theta}(a^t | \mathbf{x}, \mathbf{w}, a^{1:t-1}; \theta), \quad (4)$$

where  $\mathbf{x}$  and  $\mathbf{w}$  denote the visual and textual elements, respectively, and  $\mathbf{a} = \{a_1, a_2, \dots, a_T\}$  is the ground truth answer.

### 3 Experiments and results

**Datasets:** To evaluate our method, we make use of several publically available datasets [17]: (1) DME-VQA: contains questions on diabetic macular edema (DME) risk grade and about the presence of biomarkers in the entire image or specific regions. (2) RIS-VQA: contains images from the DaVinci robot during gastrointestinal surgery and questions related to surgical instruments. (3) INSEGCAT-VQA: contains frames from cataract surgery videos and questions about instruments used in this type of surgery. A summary of these is shown in Table 1. For all datasets, we use the same partitioning as in [17].



Dataset	Modality	# images	# QA-pairs
DME-VQA	Fundus	679	13470
RIS-VQA	Gastrointestinal	2978	32562
INSEGCAT-VQA	Cataract surgery	4647	39008

Table 1. Dataset parameters.




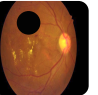


Baseline	No mask	Region in text	Draw region	Context only	Crop region	LocAtt
Input Image(s)						
Input Question	Are there hard exudates in this region?	Are there hard exudates in the ellipse contained in the bounding box (2 0 0, 8 1, 4 1 6, 2 2 4)	Are there hard exudates in this region?	Are there hard exudates in this region?	Are there hard exudates in this region?	Are there hard exudates in this region?

Fig. 3. Example input images and questions for evaluated baselines. In the baseline “Region in text,” the digits are separated to provide a fair scenario to the LLM.

**Baselines:** We benchmark our method against multiple baselines, which are exemplified in Fig. 3. In **No mask**, the model receives no information about the location of the region; in **Region in text**, the region is specified in the question; in **Draw region**, the region is marked on top of the image. In **Context only**, the model only sees the context, but not the contents of the region; in **Crop region**, the model receives no context; finally, in **LocAtt** [17], the model has access to the image, as well as a binary image representing the region. For these baselines, the visual prompt given to the model is: “Answer the question below using the context below Context: <Img><Image></Img> Question:<Question> Answer:”

**Implementation details:** We use R2GenGPT [23] as base MLLM, adapting it from the task of radiology report generation to VQA. We use a pre-trained Swin Transformer [9] as visual encoder and Llama 2 7B [19] as LLM, initialized with its official weights. Different from R2GenGPT, we finetune all modules, including the LLM, end-to-end. We train all our models for 15 epochs, with a batch size of 8 and learning rate of 1e-4, with the AdamW optimizer and a cosine annealing scheduler with minimum learning rate 1e-6. For the text generation, we use a repetition penalty of 2.0 and a length penalty of -1.0. Our implementation uses PyTorch 2.0.1 and two Nvidia A100 cards with 80 GB of memory each.

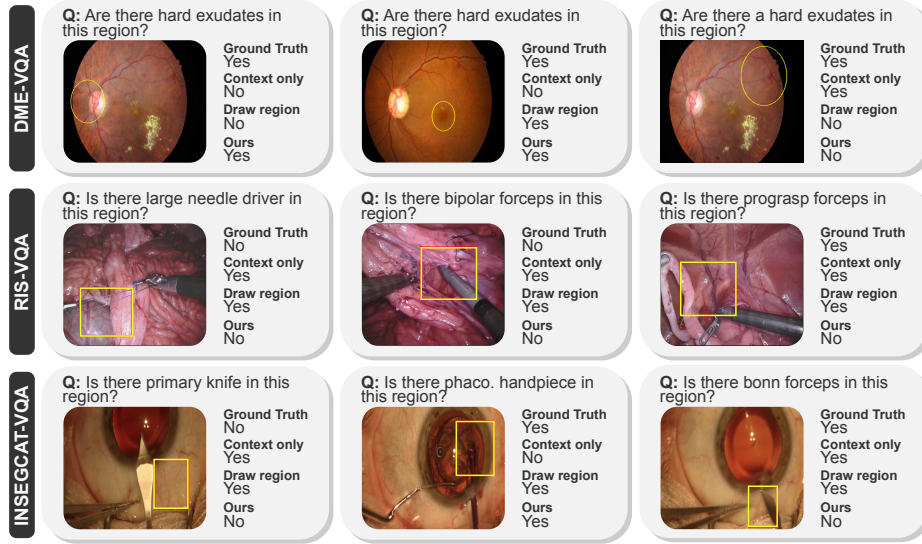
Dataset	Method	Accuracy (%)	F1 score (%)
DME-VQA	No Mask	57.32	57.32
	Region in Text [22]	62.12	63.59
	Crop Region [16]	86.52	87.26
	Draw Region	86.86	86.85
	Context Only	88.07	88.45
	<b>Ours</b>	<b>90.30</b>	<b>90.22</b>
	LocAtt [17]*	84.2	85.79
RIS-VQA	No Mask	50.00	50.00
	Region in Text [22]	64.81	65.39
	Crop Region [16]	85.50	85.64
	Draw Region	91.30	91.43
	Context Only	91.77	91.81
	<b>Ours</b>	<b>92.60</b>	<b>92.54</b>
	LocAtt [17]*	82.73	86.15
INSEGCAT-VQA	No Mask	50.00	50.00
	Region in Text [22]	73.51	74.55
	Crop Region [16]	90.91	90.93
	Draw Region	95.44	95.43
	Context Only	95.19	95.17
	<b>Ours</b>	<b>95.51</b>	<b>95.47</b>
	LocAtt [17]*	88.13	90.14

**Table 2.** Accuracy and F1 score comparison to SOTA approaches on the DME-VQA, RIS-VQA and INSEGCAT-VQA datasets. For the DME-VQA dataset only localized questions are considered. \*This result corresponds to a different architecture, but we include it for completeness.

### 3.1 Results

Table 2 summarizes our results on the DME-VQA, RIS-VQA, and INSEGCAT-VQA datasets. The accuracy and F1 score are reported for all datasets. Notably, our method consistently outperforms all evaluated baselines across all datasets, underscoring the efficacy of targeted visual prompting in enhancing MLLMs with localized question capabilities.

In the case of the DME-VQA and RIS-VQA datasets, we observe that the performance of *context only* surpasses that of *crop region*. At first glance, this suggests that the context holds more relevance than the specific contents of the region. However, this behavior is likely influenced by spurious correlations between region sizes/locations, and the corresponding answers. For instance, in DME-VQA, images with a high amount of biomarkers often feature smaller regions associated with negative answers. Another reason for this behavior is that in many cases the context provides more evidence, in terms of pixel count, to answer the question, as compared to the region. For instance, in RIS-VQA, the tool can often be determined from its body without considering the tip.

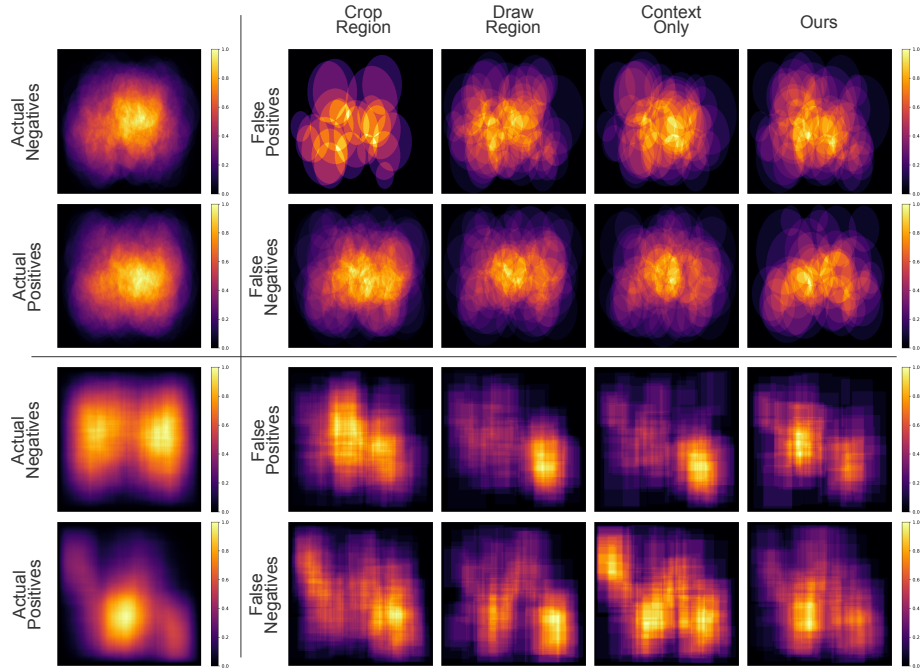


**Fig. 4.** Qualitative examples on the DME-VQA (first row), RIS-VQA (second row), and INSEGCAT-VQA (third row) datasets.

Notably, the *region in text* baseline exhibits poor performance. Given the use of a powerful LLM in the pipeline, higher performance might be expected. Different variations were explored for this baseline, including not separating the coordinate digits or replacing coordinate digits with words, but performance did not improve. We hypothesize that the model fails to correctly map location information from the text to the image, which can be at least partly attributed to using a ViT to embed the image.

We provide qualitative example results in Fig. 4. The first column exemplifies cases where our method demonstrates robustness to subtle evidence (small biomarkers), correlations (surgical suture is usually close to the needle driver), and borderline cases (evidence close to the region border). The second column highlights the weaknesses of *context only* when the context fails to provide enough evidence for the answer. Finally, the third column shows errors made by our model in tricky cases (subtle or ambiguous evidence in the region).

Fig. 5 shows error maps by region location for the DME-VQA and INSEGCAT-VQA datasets and for the four strongest baselines. On the left side of the plot, the locations of actual positives and negatives are illustrated. For the INSEGCAT-VQA dataset, this visualization reveals a location bias that other baselines without access to the region or the context may be exploiting. Due to the nature of the images (cataract surgery) and questions, regions with positive answers tend to cluster in a specific area. This, coupled with the dissimilarity of objects mentioned in the questions, explains why a baseline like *crop region* achieves relatively high performance on this dataset compared to the other two datasets



**Fig. 5.** Error analysis by region location for the four strongest baselines. The maps are obtained by adding binary masks representing the regions for all QA pairs in each category and then normalizing. **Top:** DME-VQA dataset. **Bottom:** INSEGCAT-VQA dataset.

(see Table 2). Similarly, in the case of DME-VQA, it becomes evident that the lack of context in *crop region* results in lower sensitivity, highlighting the significance of context even when the isolated region should theoretically provide sufficient evidence. Fig. 5 also demonstrates that *draw region* and *context only* exhibit marked clusters of false positives and negatives in INSEGCAT-VQA, potentially indicating the utilization of location biases. In contrast, our method produces a more evenly distributed location for both types of errors.

## 4 Conclusions

In this work, we introduced a novel approach to enable localized questions in multimodal LLMs for the tasks of VQA. Our proposed approach involves the utilization of targeted visual prompting, granting the model access not only to the region and its context within the image but also to an isolated version of the region. Doing so allows two perspectives to be encoded in the prompt and more fine-grained information to be leveraged. Our approach demonstrates enhanced performance across all evaluated datasets compared to various baselines. Future

works include extending the methodology to accommodate multiple images and enabling the use of comparison questions.

**Prospect of application:** This approach aims to be useful for medical assistants/chatbots that can help doctors assess specific parts of an image that look suspicious. By providing a second opinion, it can improve the accuracy of diagnoses. Additionally, this technology could help medical students learn and reinforce medical concepts by enabling a more modular analysis of medical images.

**Acknowledgments.** This work was partially funded by the Swiss National Science Foundation through grant 191983.

**Disclosure of Interests.** The authors have no competing interests to declare that are relevant to the content of this article.

## References

1. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang, L.: Bottom-up and top-down attention for image captioning and visual question answering. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 6077–6086 (2018)
2. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: Vqa: Visual question answering. In: *Proceedings of the IEEE international conference on computer vision*. pp. 2425–2433 (2015)
3. Chen, C., Qin, R., Luo, F., Mi, X., Li, P., Sun, M., Liu, Y.: Position-enhanced visual instruction tuning for multimodal large language models. *arXiv preprint arXiv:2308.13437* (2023)
4. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 6904–6913 (2017)
5. Gupta, D., Suman, S., Ekbal, A.: Hierarchical deep multi-modal network for medical visual question answering. *Expert Systems with Applications* **164**, 113993 (2021)
6. Hudson, D.A., Manning, C.D.: Gqa: a new dataset for compositional question answering over real-world images. *arXiv preprint arXiv:1902.09506* **3**(8) (2019)
7. Liu, F., Peng, Y., Rosen, M.P.: An effective deep transfer learning and information fusion framework for medical visual question answering. In: *International Conference of the Cross-Language Evaluation Forum for European Languages*. pp. 238–247. Springer (2019)
8. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. *arXiv preprint arXiv:2304.08485* (2023)
9. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 10012–10022 (2021)
10. Mani, A., Yoo, N., Hinthorn, W., Russakovsky, O.: Point and ask: Incorporating pointing into visual question answering. *arXiv preprint arXiv:2011.13681* (2020)

11. Nguyen, B.D., Do, T.T., Nguyen, B.X., Do, T., Tjiputra, E., Tran, Q.D.: Overcoming data limitation in medical visual question answering. In: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (eds.) Medical Image Computing and Computer Assisted Intervention – MICCAI 2019. pp. 522–530. Springer International Publishing, Cham (2019)
12. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)
13. Ribeiro, M.T., Guestrin, C., Singh, S.: Are red roses red? evaluating consistency of question-answering models. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. pp. 6174–6184 (2019)
14. Seenivasan, L., Islam, M., Kannan, G., Ren, H.: Surgicalgpt: End-to-end language-vision gpt for visual question answering in surgery. arXiv preprint arXiv:2304.09974 (2023)
15. Selvaraju, R.R., Tendulkar, P., Parikh, D., Horvitz, E., Ribeiro, M.T., Nushi, B., Kamar, E.: Squinting at vqa models: Introspecting vqa models with sub-questions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10003–10011 (2020)
16. Tascon-Morales, S., Márquez-Neila, P., Sznitman, R.: Consistency-preserving visual question answering in medical imaging. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part VIII. pp. 386–395. Springer (2022)
17. Tascon-Morales, S., Márquez-Neila, P., Sznitman, R.: Localized questions in medical visual question answering. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 361–370. Springer (2023)
18. Tong, S., Liu, Z., Zhai, Y., Ma, Y., LeCun, Y., Xie, S.: Eyes wide shut? exploring the visual shortcomings of multimodal llms. arXiv preprint arXiv:2401.06209 (2024)
19. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al.: Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023)
20. Tsimpoukelli, M., Menick, J.L., Cabi, S., Eslami, S., Vinyals, O., Hill, F.: Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems* **34**, 200–212 (2021)
21. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
22. Vu, M.H., Löfstedt, T., Nyholm, T., Sznitman, R.: A question-centric model for visual question answering in medical imaging. *IEEE transactions on medical imaging* **39**(9), 2856–2868 (2020)
23. Wang, Z., Liu, L., Wang, L., Zhou, L.: R2gengpt: Radiology report generation with frozen llms. *Meta-Radiology* **1**(3), 100033 (2023)
24. Yin, S., Fu, C., Zhao, S., Li, K., Sun, X., Xu, T., Chen, E.: A survey on multimodal large language models. arXiv preprint arXiv:2306.13549 (2023)
25. Zhan, L.M., Liu, B., Fan, L., Chen, J., Wu, X.M.: Medical visual question answering via conditional reasoning. In: Proceedings of the 28th ACM International Conference on Multimedia. pp. 2345–2354 (2020)
26. Zhang, D., Yu, Y., Li, C., Dong, J., Su, D., Chu, C., Yu, D.: Mm-llms: Recent advances in multimodal large language models. arXiv preprint arXiv:2401.13601 (2024)

27. Zhang, S., Sun, P., Chen, S., Xiao, M., Shao, W., Zhang, W., Chen, K., Luo, P.: Gpt4roi: Instruction tuning large language model on region-of-interest. arXiv preprint arXiv:2307.03601 (2023)
28. Zhang, X., Wu, C., Zhao, Z., Lin, W., Zhang, Y., Wang, Y., Xie, W.: Pmc-vqa: Visual instruction tuning for medical visual question answering. arXiv preprint arXiv:2305.10415 (2023)

# Supplementary Material

## Targeted Visual Prompting for Medical Visual Question Answering

Anonymous

Anonymous Organization  
email

Method	Accuracy (%)			
	Overall	Grade	Whole	Macula
No Mask	60.50	81.13	76.42	85.85
Region in Text	64.75	79.25	83.96	82.08
Crop Region	86.05	80.19	83.96	84.91
Draw Region	86.18	79.25	83.02	83.02
Context Only	82.61	76.42	87.74	90.57
<b>Ours</b>	89.29	79.25	83.96	84.91

Table 1. Accuracy for the DME-VQA dataset by question type.

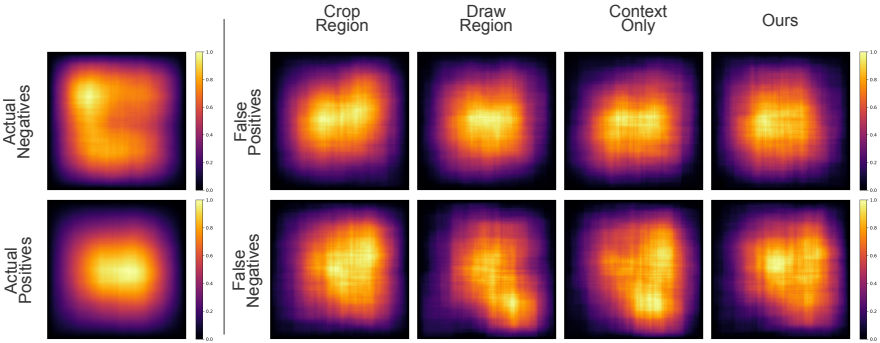
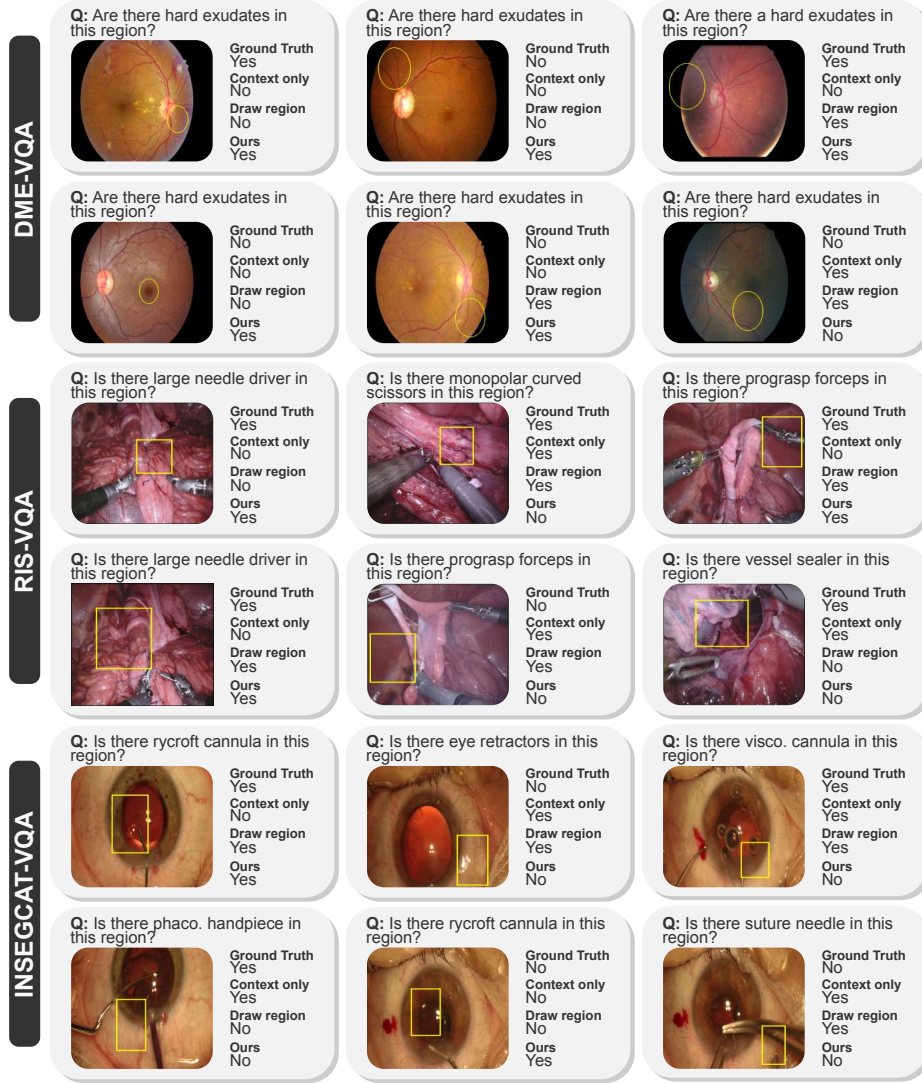


Fig. 1. Error analysis by region location for the four strongest baselines for the RIS-VQA dataset. The maps are obtained by adding binary masks representing the regions for all QA pairs in each category and then normalizing.





**Fig. 2.** Additional examples for DME-VQA (rows 1 and 2), RIS-VQA (rows 3 and 4) and Insecat-vQA (rows 5 and 6).