# MLLMS KNOW WHERE TO LOOK: TRAINING-FREE PERCEPTION OF SMALL VISUAL DETAILS WITH MULTIMODAL LLMS
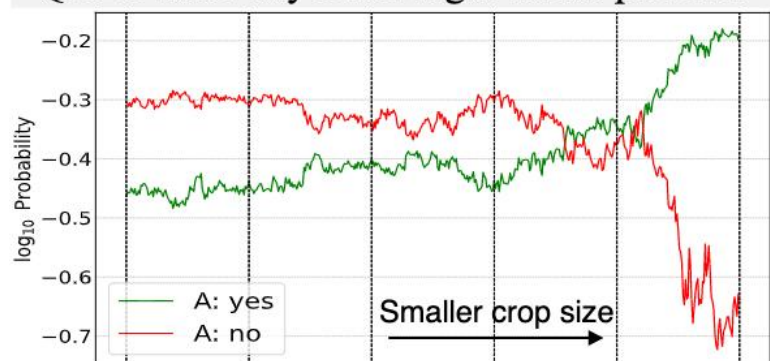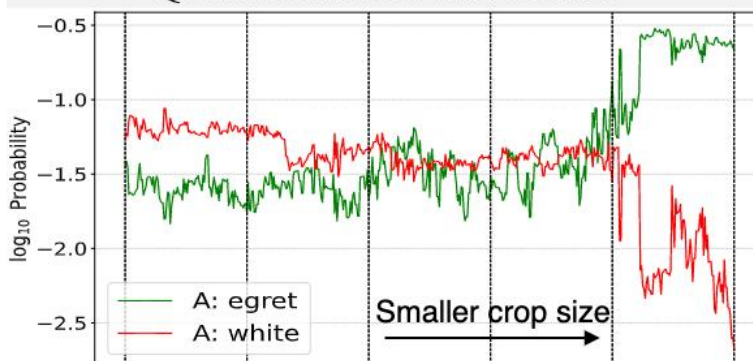
ICLR 2025

# Motivation

Given MLLMs' integration into many applications, it is important to understand the limitations of their **visual perception**.

- Model's performance is very sensitive to the size of the visual subject of the question.
- Model was not making a semantic error, rather it was unable to perceive sufficient details.
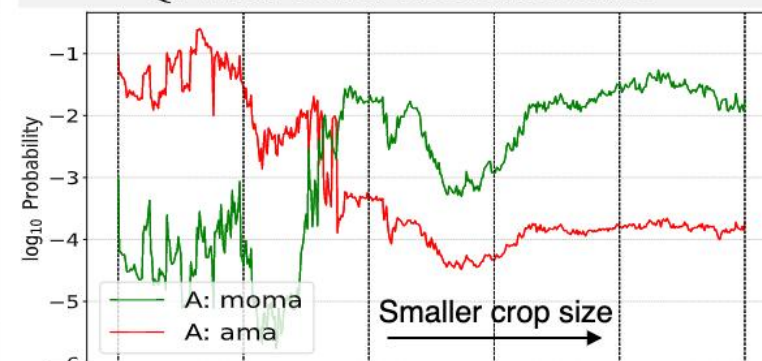- Model knows where to look based on the question.

# MLLMS' SENSITIVITY TO THE SIZE OF VISUAL CONCEPTS

question

what is the name of the company on the card?

answer

[ "blink", "intergrative nutrition", "blink",
"blink", "blink", "blink", "blink", "blink",
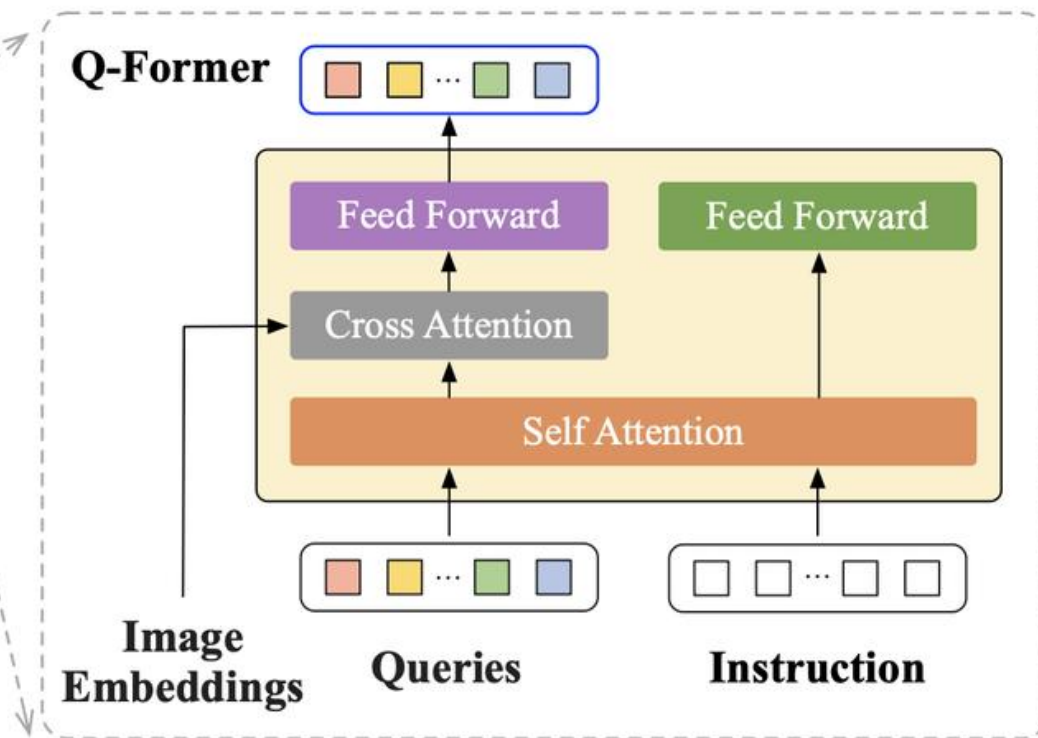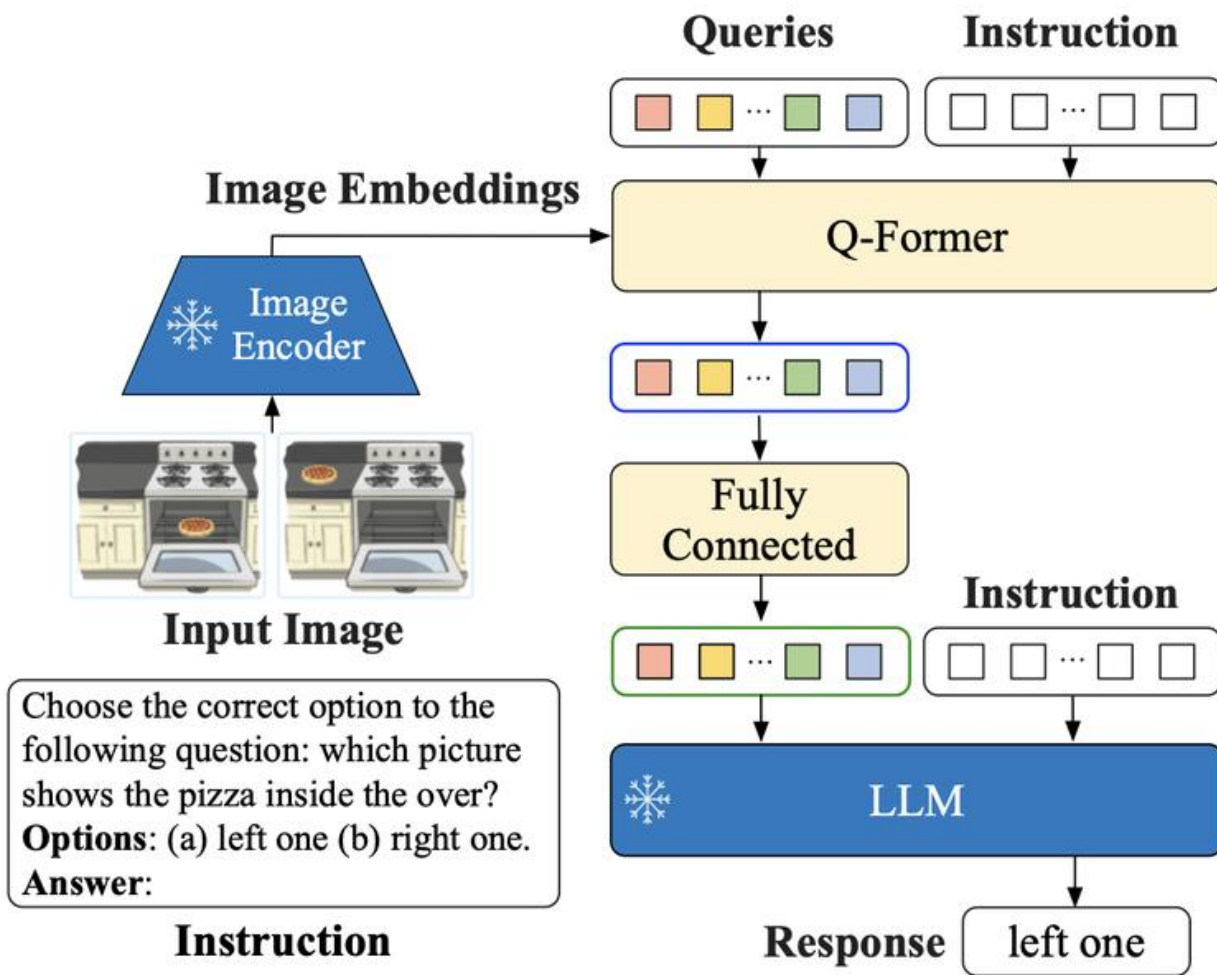"blink", "blink" ]

dataset_id

36269

bbox

[ 712, 255, 64, 43 ]

$$S = \frac{A_{bb}}{A_{total}}$$

Small:      S < 0.005
Medium:  0.005 ⩽ S < 0.05
Large:      S ⩾ 0.05

# MLLMS' SENSITIVITY TO THE SIZE OF VISUAL CONCEPTS

| Model | Method | Answer Bbox Size ($S$) | | |
|---|---|---|---|---|
| | | small | medium | large |
| BLIP-2 (FlanT5$_{XL}$) | no cropping | 12.13 | 19.57 | 36.32 |
| | human-CROP | 55.76 | 52.02 | 45.73 |
| InstructBLIP (Vicuna-7B) | no cropping | 21.79 | 30.58 | 45.30 |
| | human-CROP | 69.60 | 61.56 | 53.39 |
| LLaVA-1.5 (Vicuna-7B) | no cropping | 39.38 | 47.74 | 50.65 |
| | human-CROP | 69.95 | 65.36 | 56.96 |
| Qwen-VL (Qwen-7B) | no cropping | 56.42 | 65.09 | 68.60 |
| | human-CROP | 70.35 | 75.49 | 71.05 |
| GPT-4o | no cropping | 65.76 | 72.81 | 69.17 |
| | human-CROP | 75.63 | 81.32 | 71.72 |

**Queries**

**Instruction**

**Image Embeddings**

Q-Former

**Image Encoder**

**Input Image**

Choose the correct option to the following question: which picture shows the pizza inside the over?
**Options**: (a) left one (b) right one.
**Answer**:

**Instruction**

Fully Connected

**Instruction**

LLM

**Response** left one

**Q-Former**

Feed Forward

Feed Forward

Cross Attention

Self Attention

**Image Embeddings**

**Queries**

**Instruction**

# DO MLLMS KNOW WHERE TO LOOK?

The limitation in perceiving small visual concepts:

(1) they are hard to locate in the larger image

(2) their small details are hard to perceive correctly

How to quantify?
**Attention Map**

Answer–to–token attention: $A_{st}(x, q) \in \mathbb{R}^{L \times H \times 1 \times T}$

$$\hat{A}_{st}(x, q) = \frac{1}{H} \sum_{h=1}^{H} A_{st}(x, q)$$

q': Write a general description of the image.

Token–to–image attention: $A_{ti} \in \mathbb{R}^{L_c \times H_c \times T \times N^2}$

$$\hat{A}_{ti}(x) = \frac{1}{H_c} \sum_{h=1}^{H_c} A_{ti}(x)$$

$$A_{rel}(x, q) = \frac{A_{si}(x,q)}{A_{si}(x,q')}$$

Answer–to–image attention: $A_{si}(x, q) \in \mathbb{R}^{L \times L_c \times 1 \times N^2}$
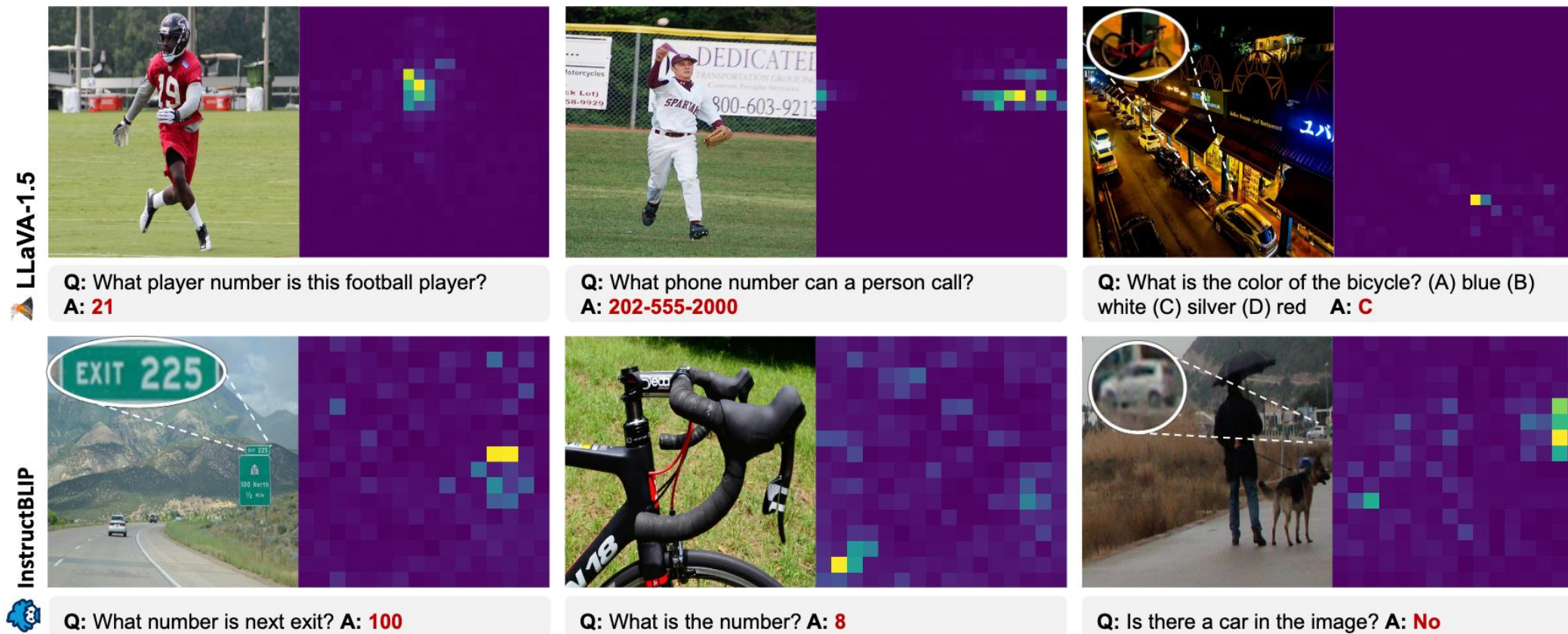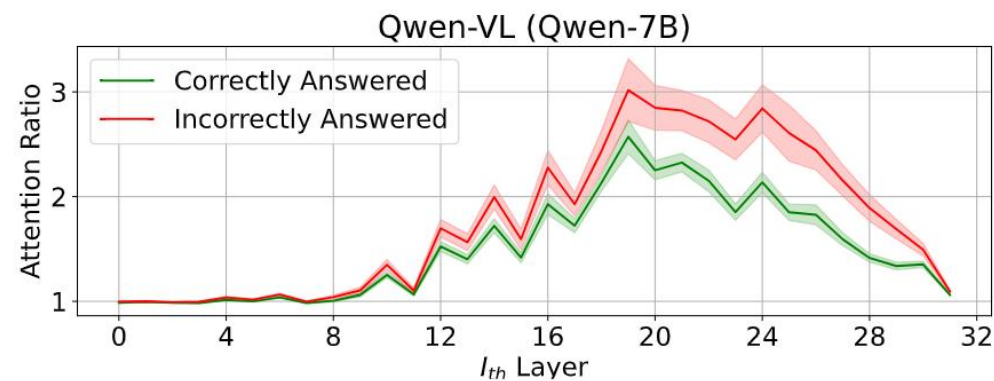
# DO MLLMS KNOW WHERE TO LOOK?



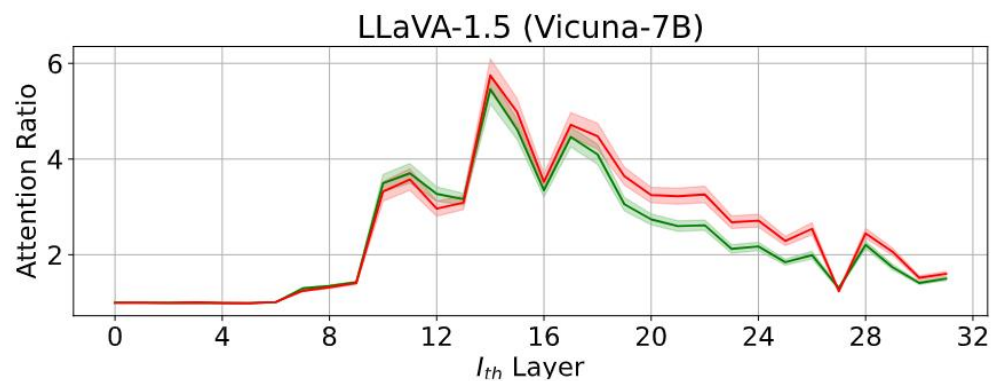Figure 2: Examples of MLLMs knowing where to look despite answering incorrectly. The right panel in each example displays relative attention to the image (defined in Sec. 4) of one layer in the MLLM.

# DO MLLMS KNOW WHERE TO LOOK?

Attention Ratio: the ratio of the total (sum) relative attention inside the answer ground–truth bounding box to its average across all bounding boxes of the same size.

# AUTOMATIC VISUAL CROPPING (VICROP)



Figure 4: Illustration of the proposed visual cropping approach applied to two MLLMs.

# AUTOMATIC VISUAL CROPPING (VICROP)

**Relative Attention ViCrop (`rel-att`).** In this method, we directly compute the relative attention $A_{rel}(x, q)$ defined in Sec. 4 for each image-question pair $(x, q)$. We then select a target layer (in LLM and connector) based on a small held-out set of samples in TextVQA and use its relative attention as the importance map for visual cropping (discussed below). We ablate on the choice of layer in Sec. 6.

**Gradient-Weighted Attention ViCrop (`grad-att`).** The relative attention runs an additional generic instruction through the MLLM to normalize the answer-to-image attention and emphasize semantically relevant attention. As an alternative that does not require a second forward pass, we consider using the gradients to normalize attention, because the gradient of the model's decision with respect to an attention score shows how sensitive the decision is to changes in that attention, hence how semantically relevant the attention is for answering the question. To get a differentiable representation of the model's decision, we consider the logarithm of the maximum output probability at the starting answer token, $v = \log \text{softmax}(z(x, q))_{t^*} \in \mathbb{R}$, where $z \in \mathbb{R}^D$ is the output logit of the LLM at the starting answer position, $D$ the vocabulary size, and $t^* = \arg\max_t z_t$. Then, following our notation in Sec. 4, we can compute the gradient-weighted versions of answer-to-token attention $\tilde{A}_{st}(x, q) = A_{st}(x, q) \odot \sigma(\nabla_{A_{st}} v(x, q))$ and token-to-image attention $\tilde{A}_{ti}(x, q) = A_{ti}(x) \odot \sigma(\nabla_{A_{ti}} v(x, q))$, where $\odot$ is element-wise product and $\sigma(w) = \max(0, w)$. We remove negative gradients because they correspond to tokens that if attended to will reduce the model certainty, hence often distracting locations Selvaraju et al. (2017). Finally, we compute the gradient-weighted answer-to-image attention as the following tensor product $\tilde{A}_{si}(x, q) = \tilde{A}_{st}(x, q) \otimes \tilde{A}_{ti}(x, q) \in \mat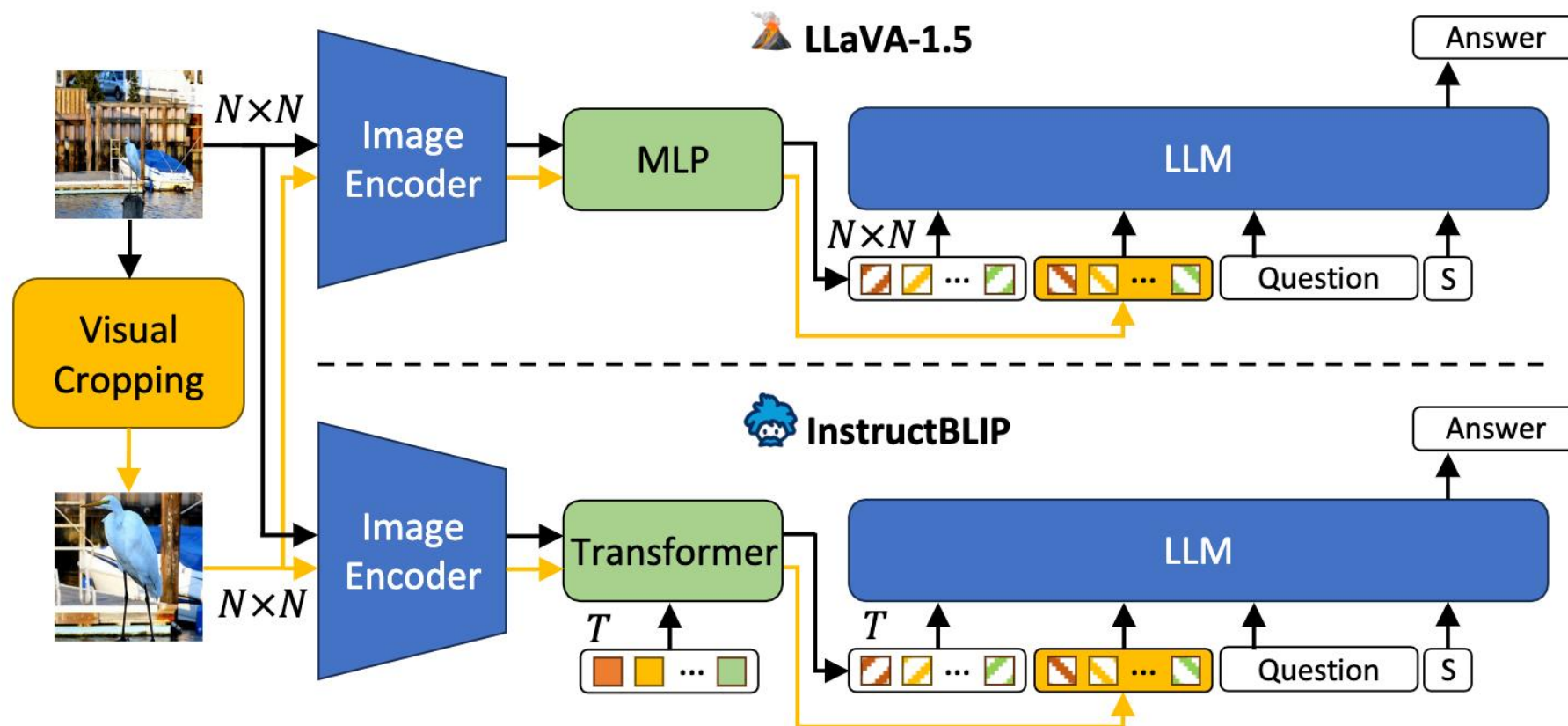hbb{R}^{L \times L_c \times 1 \times N^2}$. We will select the same target layer identified in `rel-att` from $\tilde{A}_{si}(x, q)$ as the importance map for visual cropping.

# AUTOMATIC VISUAL CROPPING (VICROP)

**Input Gradient ViCrop (`pure-grad`).** In this method, we seek to find the relevant regions on the image directly using the gradient of the MLLM's decision with respect to the input image. Compared to the previous attention-based methods, `pure-grad` is more versatile because it does not rely on the Transformer-based architecture. Specifically, for each image-question pair $(x, q)$, we will compute $G(x, q) = \|\nabla_x v(x, q)\|_2$, where $v(x, q)$ is the logarithm of the maximum output probability of the LLM at the starting answer token (as defined in `grad-att` above), and the L2-norm is taken over the image channel dimension. However, gradients sometimes show high values in entirely constant-color regions (*e.g.*, blue skies). Given that these non-edge regions do not contain any visual details, we will explicitly diminish them in $G$. To that end, we first apply a $3 \times 3$-size Gaussian high-pass filter to the image, followed by a median filter of the same size to reduce salt-and-pepper noise. The resulting filtered image is then thresholded at its spatial median value to become a binary mask and is element-wise multiplied by $G$. Finally, the edge-emphasized $G$ is spatially average-pooled into the $N \times N$ patches of the MLLM to become an importance map for visual cropping.

**Bounding Box Selection for Visual Cropping.** To convert the importance map (from each of the methods described above) to a bounding box, we use sliding windows of different sizes inspired by object detection literature Redmon et al. (2016). Specifically, for each MLLM, we define a set of windows with sizes equal to a multiple of the input image resolution of the MLLM. The multiples are in $\{1, 1.2, \ldots 2\}$. We slide each window over the image with a stride of 1 and find its best position where the sum of the importance map inside the window is maximized. After selecting the best position per window, we choose the window whose internal sum has the largest difference from the average internal sum of its adjacent positions. This latter step is a heuristic to avoid choosing too small or too large windows (notice that in both cases, moving the window slightly left/right or up/down will not change its internal sum significantly). The chosen window is then cropped from the image, resized to the input image resolution of the MLLM, and provided to the MLLM in addition to the image-question pair.

# AUTOMATIC VISUAL CROPPING (VICROP)

**High-Resolution Visual Cropping.** In one of the benchmarks we consider in this work, V* Wu and Xie (2023), the images are of very high resolution (always more than 1K) and consequently, the resized input image provided to the MLLM might completely lose the visual concept of interest for a question. To mitigate this, on this benchmark, we employ a two-stage strategy. In the first stage, we divide images into non-overlapping blocks of smaller than $1024 \times 1024$ with an aspect ratio close to 1, compute the importance map separately for each block according to the ViCrop methods, and then re-attach the blocks back together. In the second stage, we find the bounding box for visual cropping on this re-attached importance map exactly as described before and provide the original image-question pair with the resized cropped image to the MLLM.

# AUTOMATIC VISUAL CROPPING (VICROP)



Figure 5: Examples of `rel-att` helping MLLMs correct their mistakes (cyan-colored bounding box shows cropped region by `rel-att`; zoom-in insets are displayed for better readability).

# AUTOMATIC VISUAL CROPPING (VICROP)

Table 2: Accuracy of the proposed ViCrop methods on visual question answering benchmarks.

| Model | | Smaller Visual Concepts | | | | Larger Visual Concepts | | |
|---|---|---|---|---|---|---|---|---|
| | | TextVQA[†] | V* | POPE | DocVQA | AOKVQA | GQA | VQAv2 |
| LLaVA-1.5 | no cropping | 47.80 | 42.41 | 85.27 | 15.97 | 59.01 | 60.48 | 75.57 |
| | `rel-att` | 55.17 | **62.30** | **87.25** | 19.63 | **60.66** | 60.97 | **76.51** |
| | `grad-att` | **56.06** | 57.07 | 87.03 | **19.84** | 59.94 | **60.98** | 76.06 |
| | `pure-grad` | 51.67 | 46.07 | 86.06 | 17.70 | 59.92 | 60.54 | 75.94 |
| InstructBLIP | no cropping | 33.48 | 35.60 | 84.89 | 9.20 | 60.06 | 49.41 | 76.25 |
| | `rel-att` | 45.44 | **42.41** | 86.64 | 9.95 | 61.28 | 49.75 | **76.84** |
| | `grad-att` | **45.71** | 37.70 | **86.99** | **10.81** | **61.77** | **50.33** | 76.08 |
| | `pure-grad` | 42.23 | 37.17 | 86.84 | 8.99 | 61.60 | 50.08 | 76.71 |

# AUTOMATIC VISUAL CROPPING (VICROP)

Table 3: Ablation study on the choice of layer and the use of high-resolution visual cropping.

| Model | | Choice of Layer | | | High-Resolution ViCrop | | |
|---|---|---|---|---|---|---|---|
| | | Selective | Average | Δ | w/ High-Res | w/o High-Res | Δ |
| LLaVA-1.5 | no cropping | 47.80 | – | – | 42.41 | 42.41 | – |
| | rel-att | 55.17 | 55.45 | +0.28 | 62.30 | 47.64 | -14.66 |
| | grad-att | 56.06 | 56.26 | +0.20 | 57.07 | 49.74 | -7.33 |
| | pure-grad | 51.67 | – | – | 46.07 | 45.03 | -1.04 |
| InstructBLIP | no cropping | 33.48 | – | – | 35.60 | 35.60 | – |
| | rel-att | 45.44 | 44.40 | -1.04 | 42.41 | 38.74 | -3.67 |
| | grad-att | 45.71 | 44.98 | -0.73 | 37.70 | 42.41 | +4.71 |
| | pure-grad | 42.23 | – | – | 37.17 | 42.41 | +5.24 |

# AUTOMATIC VISUAL CROPPING (VICROP)

| | Model | Original | SAM | YOLO | CLIP | rel-att | grad-att | pure-grad |
|---|---|---|---|---|---|---|---|---|
| Accuracy (TextVQA) | LLaVA-1.5 | 47.80 | 49.42 | 48.84 | 48.55 | 55.17 | 56.06 | 51.67 |
| | InstructBLIP | 33.48 | 39.23 | 36.49 | 39.61 | 45.44 | 45.71 | 42.23 |
| CPU Time | LLaVA-1.5 | 2.26 | 91.53 | 0.97 | 5.46 | 14.43 | 11.33 | 29.86 |
| | InstructBLIP | 0.66 | | | | 4.35 | 3.78 | 7.04 |
| GPU Time | LLaVA-1.5 | 0.17 | 3.33 | 0.35 | 1.07 | 1.16 | 0.89 | 2.36 |
| | InstructBLIP | 0.06 | | | | 0.28 | 0.29 | 0.60 |

**Limitations**

1. Questions concerning relations and counting are particularly difficult for ViCrop methods to help answer.

2. ViCrop can only focus on one region in the image.

3. Time overhead and the addition of visual tokens.