

Token Activation Map to Visually Explain Multimodal LLMs

Yi Li, Hualiang Wang, Xinpeng Ding, Haonan Wang, Xiaomeng Li✉

Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology, Hong Kong SAR, China.
 {yilin,xmli}@ust.hk

Abstract

Multimodal large language models (MLLMs) are broadly empowering various fields. Despite their advancements, the explainability of MLLMs remains less explored, hindering deeper understanding, model credibility, and effective visualization. Unlike conventional vision models (e.g., CNNs, ViTs, CLIP) that produce a single output, MLLMs generate sequences of tokens progressively, where each generated token depends on the previous context. Therefore, earlier context tokens can introduce redundant activations that interfere with the explanation of later tokens beyond their original information. Existing studies often overlook this issue, but our observations reveal that these redundant correlations can significantly hurt the reliability of explanations. To address this, we propose an estimated causal inference method to mitigate the interference of context to achieve high-quality MLLM explanation, with a novel rank Gaussian filter to further reduce activation noises. We term this method Token Activation Map (TAM) to highlight the consideration of interactions between tokens. TAM also indicates that it excels at explaining multiple tokens of MLLM, which is different from the Class Activation Map (CAM) for a single prediction. Our TAM method significantly outperforms existing SoTA methods, showcasing high-quality visualization results that can be utilized for various scenarios, such as object localization, failure case analysis, video visualization, MLLMs visual comparison, and model understanding (e.g., color, shape, action, location, visual reasoning, multi-turn conversation, etc). The code is available at github.com/xmed-lab/TAM.

1. Introduction

Multimodal large language models (MLLMs or Multimodal LLMs) increasingly empower wide applications [17, 29, 30, 34], enabling multimodal inputs (e.g., images, videos, text) and human-like conversations. Although extensive efforts have been devoted to enhancing the performance of MLLMs [4, 16, 35, 47, 51], research on MLLM explainability remains less explored. It is crucial for user trust, model

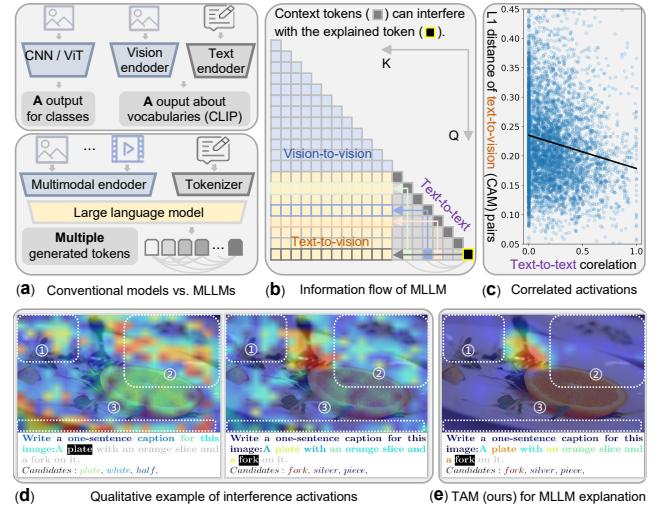


Figure 1. Illustration of Motivation. (a) MLLMs generate multiple tokens beyond a single output. (b) The information flow indicates that MLLM generates tokens progressively, where one generated token (each row) is correlated with the context (prompt + earlier answer tokens). (c) We randomly pair CAMs and count their L1 distance against text correlation. Higher text correlation corresponds to lower distance, indicating concurrent interferences. (d) The example shows that the context token “plate” introduces interference activations (marked by white boxes) to the later explained token “fork”. (e) Our TAM reveals the original information of the explained token “fork”, eliminating context interferences and beyond CAM for a single output in (d). Results in (c, d, e) derive from Qwen2-VL-2B on the COCO Caption dataset.

understanding, case analysis, and visualization. However, exploring the explainability of MLLMs is more challenging than conventional visual models (such as CNN [24], ViT [23], CLIP [41]), where MLLMs progressively generate multiple tokens beyond a single output as Fig. 1a. Thus, we aim to explain all tokens of MLLM, revealing the visual cues via token-level activation maps (e.g., Fig. 1e).

Current MLLM explanation works [7, 48, 55] are generally borrowed from techniques [38, 45] for conventional models. For example, LLaVA-CAM [55] uses Smooth Grad-CAM++ [38] to explain LLaVA [35]. The self-

attention weights [49] are applied to understand math problem-solving [48]. Besides, LVLM-Interpret [7] integrates some conventional methods as a tool. Unfortunately, existing visual explainability methods are designed for conventional models [23, 24, 41] that produce a single output, such as class activation maps (CAM) series [9, 26, 45, 56], attention-based approaches [1, 28], attention with relevance [2, 11], model-agnostic methods [36, 43], or multimodal explanation [3, 19, 31, 32]. *However, these methods overlook that MLLMs generate multiple tokens progressively, where earlier context tokens interfere with later tokens to be explained.* As shown in Fig. 1d, the conventional CAM [56] for the token “fork” is confounded by numerous redundant activations concurrent with the token “plate” in context (context refers to prompt and earlier answer tokens than the explained token). Quantitative analysis in Fig. 1c reveals that this issue is widespread, where correlated tokens lead to close and interference activations. These interferences obscure the original information conveyed by the explained token, significantly degrading the explainability of MLLMs. Notably, this activation interference across tokens is unique to MLLMs and has been overlooked.

We aim to reveal the original information for the explained token while minimizing interference from context, as discussed above. Herein, the original information reveals a causal relation [40] between the current token and the input, while the interferences represent the correlations that need to be excluded. Inspired by the potential outcome model (POM) of causal inference [27], we propose the estimated causal inference to explore this causal relation. Unlike raw causal inference by comparing ready-made results (yes vs. no), there is no single given output to represent the “no” result. Therefore, we estimated interferences from multiple context tokens that are highly related to the target, with a scale factor optimized by the least-squares method to generate the refined causal activation map. We also propose a novel rank Gaussian filter to reduce noises, thereby further enhancing the quality of activation maps. The above methods are collectively referred to as token activation map (TAM), emphasizing the consideration of interactions between tokens. TAM also indicates that it excels at explaining tokens in multiple rounds of MLLM, which differs from the Class Activation Map (CAM) for a single prediction.

Experimentally, the proposed TAM demonstrates quantitative improvement over existing SoTA explainability methods by more than 8.96% on the COCO Caption dataset [13] and 8.54% on the OpenPSG dataset [57] when applied to visually explain the Qwen2-VL-2B [51]. Meanwhile, it is complementary to existing methods to boost their performance. We further verify the scalability and applicability on 7 advanced MLLMs from the series of Qwen2-VL [51], InternVL2.5 [15] and LLaVA1.5 [35], where the performance improvements range from 5.45% to 11.0% on

COCO Caption dataset [13]. Besides quantitative results, we conduct qualitative experiments to understand MLLM attributes about color, shape, action, location, etc. Furthermore, we illustrate that TAM is capable of explaining failure cases on the QK-VQA dataset. [37], and it supports video visualization well on the STAR dataset [52]. In addition, it can be used as a tool to visually compare MLLMs qualitatively. Moreover, TAM shows wide applicability to explain multi-turn conversations, multi-image scenarios, and visual reasoning (see the catalog Table 5 in supplementary for extensive examples). Our work has three main contributions:

- We introduce TAM, a novel approach to explain MLLMs, which produces high-quality activation maps for multiple generative tokens by incorporating the estimated causal inference to alleviate interference from context tokens, with the rank Gaussian filter to reduce activation noises.
- TAM significantly outperforms current methods for MLLM explanation, while also complementing them.
- TAM demonstrates wide applicability and scalability, serving as a versatile tool for object localization, model understanding, failure case analysis, video visualization, MLLM comparison, and supporting diverse scenarios (e.g., multi-turn, multi-image, visual reasoning).

2. Related Work

Visual Explainability for Multimodal LLM. The emergence of MLLMs such as LLaVA [35], Qwen2-VL [51], InternVL [15], and closed-source GPT-4o [25] has significantly impacted various tasks involving multimodal inputs like text, images, and video. We mainly focus on the visual explanation of MLLMs via the activation maps for the input image or video frames. Regarding the modality perspective, the explainability of MLLMs is partially related to methods to explain multimodal models such as CLIP Surgery [32] and transformer register [19] designed for CLIP [41]; InterpreT [3] and Bi-Modal [10] for BERT [20]. From the technical aspect, conventional explanation methods are also related, including class activation map series (CAM [56], Grad-CAM [45], Grad-CAM++ [9], LayerCAM [26]), attention-based mechanisms (LRP [28], Rollout [1]), the combination (AttnLRP [2], Grad×AttnRoll [11]), or model-agnostic methods (LIME [43], SHAP [36]).

Although existing methods may be valid to explain MLLMs, they are not the optimal solution. Because they are usually designed for conventional models (e.g., CNN [24], ViT [23], CLIP [41]) that produce a single output rather than generating multiple tokens progressively in MLLMs. We have observed that these earlier context tokens interfere with later tokens by introducing redundant visual activations, as shown in Fig. 1. This phenomenon is first studied by us, solved by a novel estimated causal inference method, which is specially designed for MLLM. Notably, some workshop papers (LLaVA-CAM [55] and LVLM-

Interpret [7]) applied conventional methods [38, 44] to explain LLaVA [35], but they also overlooked the interference from context tokens and treated each token independently.

Causal Model. Causal inference [54] seeks to establish causal relation [40] between variables, focusing on how changes in one variable influence another. In contrast to statistical correlation analysis, it emphasizes causal relations. Another related concept, causal intervention [27], involves manipulating variables using a do operator. Model-agnostic explainability methods, such as LIME [43] and SHAP [36], can be viewed as causal interventions that explore the relationship between inputs and structured outputs by masking inputs or selecting subsets [12]. This concept is also applied to transformers through token masking (CLEANN [44]) or hidden state replacement [39]. However, causal intervention for MLLMs is impractical. Because there are too many input tokens are outputs in MLLM, requiring unbearable inference times. Besides, modified input leads to a changed context every inference, which is hard to evaluate. Our method is inspired by the potential outcome model (POM) to achieve causal inference without additional model inferences. Our contribution lies in estimating the unprocessed output from multiple tokens rather than a single ready-made output, considering the uniqueness of MLLMs.

Transformer Denoising. Activation maps of transformers [49] usually present many noise activations, impeding high-quality visualization. Recent studies give their explanations about it, including attention sink [53], lack of register [19], and redundant features across classes [32]. In MLLM, these noises still exist obversely, even system tokens already play the role of registers, or deploying the feature surgery to mitigate redundant features (See Table 2). To solve this problem practically, we aim to introduce denoising filters to mitigate it as post-processing. Since these noises belong to salt-and-pepper noise morphologically, we apply the median filter, adaptive median filter [8], and the Gaussian filter. However, these filters are not the optimal solution, where the Gaussian filter keeps too much noise signal and the median filter overlooks smaller responses. Therefore, we propose the rank Gaussian filter, a novel, simple, and effective filter to denoise transformer activation. It merges ranked values within a sliding window, weighted summed by a 1-d Gaussian kernel with the center at the median rank. Besides, we improve the Gaussian kernel by the coefficient of variation for more robust results. Overall, our method addresses the activation noises in a practical aspect through a new filter.

3. Method

In this section, we first provide an overview of the proposed Token Activation Map (TAM), as illustrated in Fig. 2. Next, we elaborate on two key modules within TAM: the estimated causal inference and the rank Gaussian filter. Finally, we introduce three metrics designed to enable fine-grained

evaluations of the explanations for MLLMs.

3.1. Token Activation Map

TAM aims to dipect explainable activation maps for multiple tokens generated from MLLMs by leveraging causal inference to mitigate inter-token interference. To be more specific, given the visual data and prompt text as the input, the MLLM is employed to progressively generate the visual features $\mathbf{F}^v \in \mathbb{R}^{n_v \times c}$, prompt features $\mathbf{F}^p \in \mathbb{R}^{n_p \times c}$ and answer features $\mathbf{F}^a \in \mathbb{R}^{n_a \times c}$. n_v , n_p , n_a and c indicate the number of three type of features, and the feature dimension, respectively. Next, a token classifier, instantiated as a fully connected layer, is used to generate answer tokens $t^a = \{t_1^a \dots t_{n_a}^a\}$ from these features. Prompt tokens are defined as $t^p = \{t_1^p \dots t_{n_p}^p\}$ similarly. Based on the variables described above, the visual activation maps for the prompt tokens and answer tokens are calculated as follows:

$$\begin{aligned}\mathbf{A}_i^p &= [\mathbf{F}^v \mathbf{w}_{t_i^p}]_+, i \in [1, n_p] \\ \mathbf{A}_i^a &= [\mathbf{F}^v \mathbf{w}_{t_i^a}]_+, i \in [1, n_a],\end{aligned}\quad (1)$$

where $\mathbf{A}_i^p \in \mathbb{R}^{n_v \times 1}$ and $\mathbf{A}_i^a \in \mathbb{R}^{n_v \times 1}$ present the activation map for the i -th prompt token t_i^p and answer token t_i^a . $\mathbf{w}_{t_i^p} \in \mathbb{R}^{c \times 1}$ and $\mathbf{w}_{t_i^a} \in \mathbb{R}^{c \times 1}$ are the corresponding weight vectors within the token classifier. $\lfloor \cdot \rfloor_+$ is the function to keep positive activations.

Then, the activation maps for prompt and answer tokens are concatenated (\perp) as $\mathbf{A}_{:n_p+n_a} = \mathbf{A}^p \perp \mathbf{A}^a$. For the i -th answer token, its context token is formulated as $\mathbf{A}_{:n_p+i-1}$, including all earlier context maps (prompt tokens + answer tokens earlier than i). On top of it, the refined activation maps for answer tokens are formulated as follows:

$$\bar{\mathbf{A}}_i^a = [\mathcal{D}(\mathbf{A}_i^a - s\mathcal{E}(\mathbf{A}_{:n_p+i-1}))]_+, i \in [1, n_a], \quad (2)$$

where $\bar{\mathbf{A}}_i^a$ indicates the refined activation map for the i -th answer token. The number of i increases sequentially from 1 to n_a . $\mathbf{A}_i^a - s\mathcal{E}(\mathbf{A}_{:n_p+i-1})$ is the proposed estimated causal inference that produces the causal activation map with the estimation function \mathcal{E} described in Sec. 3.2. \mathcal{D} is the proposed rank Gaussian filter module in Sec. 3.3.

In the last, we concatenate (\perp) the refined visual activation map $\bar{\mathbf{A}}_i^a$ of the i -th answer token with the raw textual relevance $\mathbf{r}_i \in \mathbb{R}^{n_p+i-1}$ for a multimodal activation map (image or video with text) $\mathbf{M}_i \in \mathbb{R}^{n_v+n_p+i-1}$ as follows:

$$\mathbf{M}_i = \mathcal{N}(\bar{\mathbf{A}}_i^a \perp \mathbf{r}_i), i \in [1, n_a], \quad (3)$$

where $\mathcal{N}()$ is a min-max normalization function for visualizing visual and textual activations at an aligned level.

3.2. Estimated Causal Inference

Overall, Eq. 2 optimizes later visual activation maps based on visual activation maps of all earlier context tokens. To

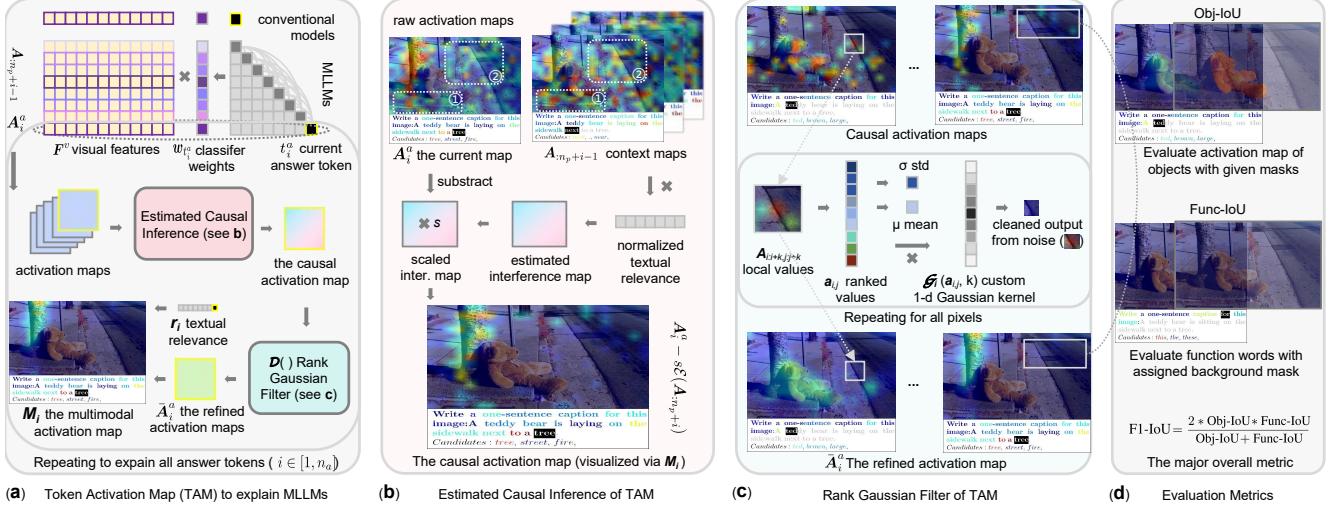


Figure 2. Illustration of method. (a, Eq. 1 - Eq. 3) The overall framework of TAM. (b, Eq. 4 - Eq. 5) Details of the estimated causal inference module. (d, Eq. 6 - Eq. 7) Details of the rank Gaussian filter module. (d, Eq. 8 - Eq. 10) Fine-grained evaluation metrics.

be more specific, causal inference is employed to investigate the correct causal relation between the current answer token and visual tokens, eliminating interferences from prompt tokens and preceding answer tokens. To this end, we quantify the interference map from the context tokens for the i -th answer token as follows:

$$\mathcal{E}(\mathbf{A}_{:n_p+i-1}) = \sum_{k=1}^{n_p+i-1} \frac{r_i^k}{\sum r_i + \epsilon} \mathbf{A}_k \quad (4)$$

$$s.t. r_i^k = 0, \text{ if } t_i^a = t_k,$$

where r_i^k is the textual relevance between the k -th context token t_k and the current answer token t_i^a . The relevance is obtained from activations among textual tokens $\lfloor (\mathbf{F}^p \perp \mathbf{F}^a) \mathbf{w}_{t_i^a} \rfloor_+$ following Eq. 1. If the current context token t_k is the same as the answer token t_i^a , we set their relevance as 0, i.e., $r_i^k = 0$, to avoid eliminating visual activations for tokens identical to the current answer token.

Overall, Eq. 4 aims for each answer token to encapsulate more visual information, which implies minimizing similarity to the visual semantics of previously context tokens. Meanwhile, we use the scale factor s in Eq. 2 to control the scale of the interference map, ensuring that the influence of the interference map is removed from the activation map at a comparable degree. Specifically, we employ the least squares method to minimize the residual between the two maps, for the optimal value of s , as follows:

$$s = \arg \min_s \sum_{j=1}^{n_v} (\mathbf{A}_{j,i}^a - s \mathcal{E}(\mathbf{A}_{j,:n_p+i-1}))^2. \quad (5)$$

After processing by the estimated causal inference module, there are still many salt-and-pepper noises in activation

maps. To address this issue, we will elaborate on the solution in the following section.

3.3. Rank Gaussian Filter

As described in Sec. 2 the activation noise exists obviously and belongs to the salt-and-pepper type. Different from existing methods [19, 32] with suboptimal performances on MLLM, we address this problem in a new aspect: image denoising. Compared with existing filters like the Gaussian filter that keeps too much noise signal and the median filter which overlooks smaller responses, our rank Gaussian filter is more robust beyond existing methods. The overall equation of the rank Gaussian filter for any activation map $\mathbf{A}_{i,j} \in \mathbb{R}^{h \times w}$ is written as:

$$\mathcal{D}(\mathbf{A}_{i,j}, k) = \sum \mathbf{A}_{i,j} \mathcal{G}(\mathbf{a}_{i,j}, k), \forall i \in H, \forall j \in W \quad (6)$$

$$s.t. \mathbf{a}_{i,j} = \text{rank}(\mathbf{A}_{i:i+k,j:j+k}).$$

Herein, $\mathbf{A}_{i:i+k,j:j+k}$ indicates local values within a sliding window at kernel size k . We sort them to get a ranked array to multiply with a custom 1-d Gaussian kernel $\mathcal{G}(\cdot)$. This custom kernel allows the median to have the largest weight, meanwhile aggregating signals ranked close to the median. It can be regarded as a kind of smoothed median filter to enhance the robustness weighted by the Gaussian kernel.

For the custom 1-d Gaussian kernel, we formulate it as Eq. 7. Specifically, $(i - k^2//2)^2$ indicates the distance from the median rank ($k^2//2$, size exact division 2), and $\sigma(\mathbf{a})/\mu(\mathbf{a})$ is the coefficient of variation which is more stable than the standard division $\sigma(\cdot)$ in general Gaussian kernel using the mean $\mu(\cdot)$ only.

$$\mathcal{G}_i(\mathbf{a}, k) = e^{-\frac{(i-k^2//2)^2}{2(\sigma(\mathbf{a})/\mu(\mathbf{a}))^2}}, \forall i \in k^2 \quad (7)$$

3.4. Evaluation for MLLM Explanation

Unlike conventional models with a single output, MLLMs generate multiple tokens. To evaluate how well activation maps correspond to objects, we propose two fine-grained IoU-based metrics: Obj-IoU for objects with manual masks, and Func-IoU for backgrounds. These metrics focus on plausibility instead of faithfulness (see concept differences in Supp. E). Plausibility better suits MLLMs since perturbation-based faithfulness tests [2, 11] which cause inconsistent text generation, making evaluation invalid (see detailed reasons in Supp. E).

To be specific, we formulate Obj-IoU in Eq. 8, counting the IoU between activation map (\mathbf{A}_i) after binarization ($\mathcal{B}()$ using opencv OTSU auto threshold) and the ground-truth mask \mathbf{G}_i . Since the object classes may be unfixed, we count the IoU for all tokens at size o over the dataset without an average on classes. We obtain the mask by matching the token to mask names. For a word or phrase containing multiple tokens, we record their max IoU for a single object.

$$\text{Obj-IoU} = \frac{1}{o} \sum_{i=1}^o \frac{\mathcal{B}(\mathbf{A}_i) \cap \mathbf{G}_i}{\mathcal{B}(\mathbf{A}_i) \cup \mathbf{G}_i} \quad (8)$$

Func-IoU measures the activation degree of tokens that are unreadable from the image (e.g., “is”, “and”, “so”, “the”). Higher Func-IoU indicates fewer false positives and unnecessary responses. Specifically, we record the mean binarization thresholds b_i (opencv OTSU auto threshold) for all noun tokens in the same image to divide foreground and background on function words. Activation map lower than it ($\mathbf{A}_i < b_i$) is the background prediction used to count IoU with the all-one matrix \mathbf{J} (ground-truth is all background class). It is applied to total u tokens over the whole dataset.

$$\text{Func-IoU} = \frac{1}{u} \sum_{i=1}^u \frac{(\mathbf{A}_i < b_i) \cap \mathbf{J}}{(\mathbf{A}_i < b_i) \cup \mathbf{J}} \quad (9)$$

We try to merge the above IoUs as the overall major metrics. We find their average is insufficient to measure biased and low-quality activation maps. For example, if one method predicts almost all tokens to be background, the Obj-IoU is very low but Func-IoU is high, forming a certain inverse ratio (e.g., 5.74% vs. 96.5% of Attention-Rollout [1] in Table 3). One suitable solution is to use the F1 value to merge metrics with a certain inverse ratio. Thus, we have the F1-IoU to evaluate the overall plausibility stably:

$$\text{F1-IoU} = \frac{2 * \text{Obj-IoU} * \text{Func-IoU}}{\text{Obj-IoU} + \text{Func-IoU}} \quad (10)$$

4. Experiments

4.1. Setup

Datasets. For quantitative experiments, we use datasets with texts and pixel-level annotations. Our main dataset is the COCO Caption dataset [13], which draws images and

masks from COCO2014 [33]. Since explainability methods require testing only, we use its 5K-image minival split without training. Additional datasets include GranDf [42] (1K images) and OpenPSG [57] (3,176-image validation set). Masks in COCO Caption and GranDf are manually annotated, while OpenPSG masks are integrated by Rasheed et al. [42]. For datasets without masks, we perform qualitative visualization tests (Sec. 4.3), including attribute-only images and the QK-VQA [37] validation set (5,046 VQA samples with textual answers). We also qualitatively evaluate the STAR video dataset [52] (914 videos).

Implementations. We use the SciPy minimize function with the BFGS method to optimize the scale factor s in Eq. 5. The only hyperparameter, the kernel size k in Eq. 7, is set to 3, as our observations indicate that salt-and-pepper noises typically cluster between 1 to 5 tokens. A kernel size of 3, consisting of 9 elements, is deemed suitable. Eq. 9 incorporates part-of-speech identification. For Func-IoU and Obj-IoU, we tag each word’s part of speech using the pos_tag function from the NLTK Python package (version 3.8.1, see tag details in Supp. E). For word matching to mask names, we employ the lemmatize function from nltk.stem. Note that the thresholds in Eq. 8 and Eq. 9 are automatically derived from OpenCV’s OTSU method, requiring no manual operation. Details of the baselines are provided in Supp. D.

4.2. Quantitative Results

ECI	RGF	Obj-IoU (%)	Func-IoU (%)	F1-IoU (%)
✗	✗	21.23	51.93	30.14
✓	✗	22.41	69.03	33.84 ^{+3.7}
✗	✓	24.82	43.34	31.57 ^{+1.33}
✓	✓	27.37	68.44	39.1 ^{+8.96}

Table 1. Ablation study on the COCO Caption [13] dataset using Qwen2-VL-2B [51]. These two modules are mutually beneficial, the combination exceeds the sum of their individual gains. ECI indicates the proposed estimated causal inference, and RGF is the rank Gaussian filter. Metrics are IoU for object words, IoU for function words, and their F1-score-like combination, respectively.

Setting	Obj-IoU (%)	Func-IoU (%)	F1-IoU (%)
Baseline	21.23	51.93	30.14
Replace rank Gaussian filter			
Adaptive Median [8]	25.48	68.2	37.1
Median Filter	26.01	68.26	37.67
Gaussian Filter	26.56	67.78	38.16
Replace estimated causal inference			
Feature surgery [32]	18.5	48.66	26.81
ECI-mean	27.84	49.85	35.72
ECI-attnWeights	27.11	54.61	36.23
Ours	27.37	68.44	39.1

Table 2. Effectiveness study on the COCO Caption [13] dataset using Qwen2-VL-2B [51]. ECI-mean and ECI-attnWeights are candidate implementations of ECI described in Supp. D.

Ablation and Effectiveness. We present the ablation study results in Table 1, which indicates that each module of the TAM is valuable. Specifically, the estimated causal infer-

Method	Type	COCO Caption			GranDf			OpenPSG		
		Obj-IoU	Func-IoU	F1-IoU	Obj-IoU	Func-IoU	F1-IoU	Obj-IoU	Func-IoU	F1-IoU
Grad-CAM [45]	Gradient	21.23	51.93	<u>30.14</u>	17.85	62.15	27.74	22.93	48.57	31.15
Grad-CAM++ [9]		19.52	62.83	29.78	17.3	73.42	<u>28.01</u>	22.21	59.95	<u>32.41</u>
Grad-Rollout [1]		1.27	99.51	2.51	1.4	99.61	2.77	1.57	99.58	13.08
Layer-CAM [26]		11.43	84.88	20.15	13.11	82.09	22.62	14.12	85.29	24.22
Attention		8.2	92.87	15.07	9.6	93.56	17.42	10.58	94.28	19.03
Attention-Rollout [1]		5.74	96.5	10.83	7.21	96.65	13.42	7.94	97.04	14.68
CP-LRP [5]		9.9	53.97	16.73	12.61	53.24	20.39	13.3	53.36	21.3
Attn-LRP [2]		9.92	52.41	16.69	12.15	52.19	19.72	12.78	52.26	20.54
CAM [56]	Combination	21.23	51.93	<u>30.14</u>	17.85	62.15	27.74	22.93	48.57	31.15
Archi.-Surgery [32]		15.69	63.82	25.19	16.59	62.28	26.2	19.83	58.77	29.65
TAM (ours)		27.37	68.44	39.1	18.65	88.97	30.83	26.26	92.99	40.95
Archi.-Surgery [32]		+4.13	+10.68	+6.12	+1.41	+26.48	+3.74	+1.59	+31.92	+5.0
Grad-CAM++ [9]	+ Ours	+2.83	+3.91	+3.7	+1.26	+5.39	+2.04	+2.75	+9.44	+4.3
Layer-CAM [26]		+5.04	-7.95	+6.98	+2.79	+0.52	+4.05	+6.62	-2.49	+8.95
CAM [56]		+6.14	+16.52	+8.96	+0.79	+26.82	+3.09	+3.33	+44.42	+9.8

Table 3. Comparison with SoTA methods using Qwen2-VL-2B on diverse datasets. We adopt the “Logit” type that uses the classifier weights in Eq. 1, without back-propagation in “Gradient” and “Combination” (Gradient + Attention). Besides, we support FlashAttention [18] and SdpaAttention, which do not return attention weights, while “Attention” and “Combination” rely on it. Our TAM is also complementary to existing methods with gains marked by “+”. Note, CAM and Grad-CAM are equivalent as discussed in Supp. D. The major metric is the F1-IoU (%) to reflect the overall result, merged from IoUs for Object words (Obj-IoU) and function words (Func-IoU).

CAM [56] +Ours	COCO Caption			GranDf			OpenPSG		
	Obj-IoU	Func-IoU	F1-IoU	Obj-IoU	Func-IoU	F1-IoU	Obj-IoU	Func-IoU	F1-IoU
LLava1.5-7B [35]	+4.47	+18.27	+7.97	+0.65	+11.67	+2.47	+3.47	+9.51	+5.16
LLava1.5-13B [35]	+4.3	+7.31	+5.45	+0.76	+7.03	+2.11	+4.22	+11.51	+6.37
InternVL2.5-2B [16]	+5.44	+19.48	+8.57	+2.2	+48.29	+8.47	+3.24	+40.45	+8.65
InternVL2.5-4B [16]	+3.53	+22.17	+7.14	+1.62	+45.19	+7.56	+2.2	+55.05	+10.85
InternVL2.5-8B [16]	+5.38	+2.13	+6.93	+3.53	+28.53	+7.02	+3.27	+26.54	+6.44
Qwen2-VL-2B [51]	+6.14	+16.52	+8.96	+0.79	+26.82	+3.09	+3.33	+44.49	+9.8
Qwen2-VL-7B [51]	+5.62	+29.4	+11.01	+1.29	+22.54	+3.4	+3.53	+46.94	+11.15
Mean Improvements	+4.98	+16.47	+8.01	+1.55	+27.15	+4.88	+3.32	+33.5	+8.35

Table 4. Explainability improvements of TAM on diverse MLLMs and datasets. Specific results are listed in Table 6 of Supp. K.

ence significantly improves Func-IoU by reducing context interference on function words, while the rank Gaussian filter performs better for object words. Importantly, these modules are mutually beneficial; their combination leads to an F1-IoU improvement of 8.96%, exceeding the total of their individual contributions (3.7% and 1.33%).

We further assess the effectiveness of these modules by comparing them with other methods, as shown in Table 2. Replacing our proposed rank Gaussian filter with alternatives yields lower performance, such as the adaptive median filter [8], the Gaussian filter, and the median filter. Although the differences are less pronounced than those seen with the causal inference module, denoising remains a challenging area. Our method outperforms feature surgery [32] (which operates on the class dimension) by 12.29% and also surpasses other ECI implementations in Supp. D.

Comparison with SoTA Methods. Table 3 presents a comparison between TAM and SoTA methods. Our TAM outperforms all others in the major F1-IoU and Obj-IoU metrics. While some methods, such as Grad-Rollout [1] and Attention, achieve high Func-IoU by producing a limited number of activations, their performance on other metrics re-

mains low. The second-best method is Grad-CAM++ [9] on the COCO Caption [13] and OpenPSG [57] datasets, while CAM [56] and Grad-CAM [45] rank second on GranDf [42]. Our method surpasses these approaches by 8.96%, 2.82%, and 8.54% on these datasets, respectively. Furthermore, TAM is complementary to existing methods and enhances them significantly, with maximum gains of 6.62%, 44.42%, and 9.8% across the three metrics, respectively.

Applicability and Scalability. We evaluate the applicability on three prominent MLLMs, as shown in Table 4, and assess its scalability across various model sizes. The results in Table 3 indicate that TAM consistently outperforms the baseline across all MLLMs and model sizes by significant margins. This demonstrates that TAM can be effectively applied to a range of MLLMs while providing superior explainability compared to existing methods. Moreover, TAM supports MLLMs comparison, complementing quantitative results. It introduces a new evaluation perspective by highlighting visual alignment quality, as illustrated in Fig. 8.

4.3. Qualitative Results

Comparison of Explainability Methods.

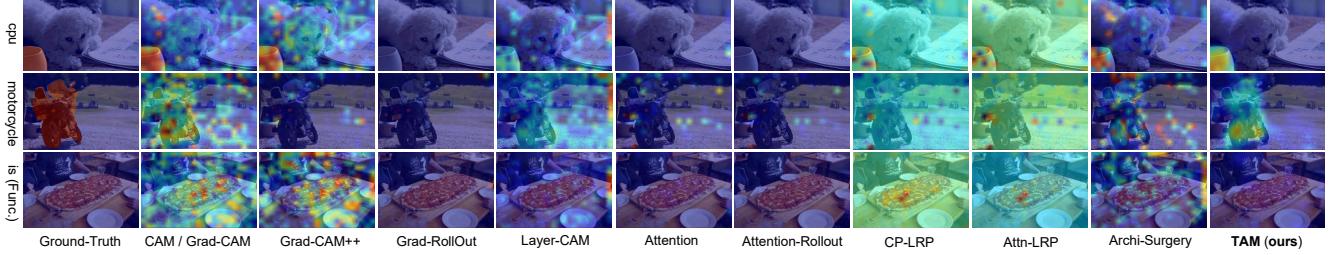


Figure 3. Visual comparison between TAM and SoTA methods on the COCO Caption dataset [13] using the Qwen2-VL-2B [51] model. The “is” in the last row is a functional word with background ground-truth. More examples are shown in Supp. G including complex cases.

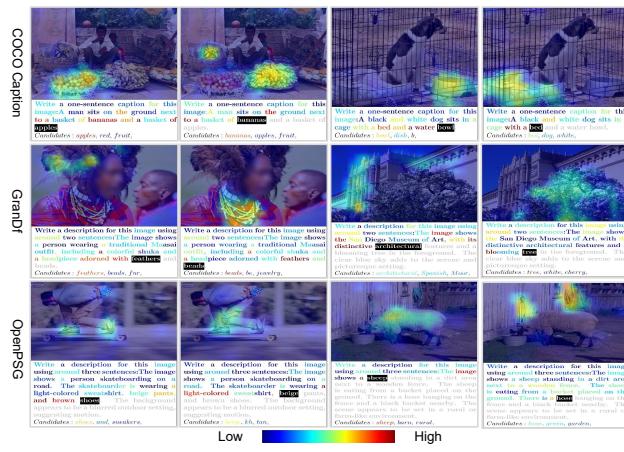


Figure 4. TAM presents high-quality localization results on diverse datasets [13, 42, 57] under the Qwen2-VL-7B [51] model.

We perform comparisons between the proposed TAM and SoTA explanation methods, as illustrated in Fig. 3. The visualization results clearly show that TAM more accurately highlights target objects and renders fewer false activations on function words. In contrast, attention-based methods [1, 2, 5] produce only a limited number of activations, while gradient-based methods [9, 45] present too many correlated responses and noise, performing significantly worse than our approach. Additionally, we present more complex examples in Supp. G, where TAM notably outperforms existing SoTA methods. TAM excels at explaining multiple tokens of MLLMs. In Supp. C, results show how TAM focuses on important tokens, whereas the baseline methods generate many redundant activations.

High-quality Localization Results. We visualized the results using the multimodal activation map in Eq. 3 (see detailed examples in Supp. B). The results demonstrate that TAM effectively locates objects across diverse datasets, as shown in Fig. 4. These findings indicate that TAM enhances MLLMs with localization capabilities, which could be potentially beneficial for various downstream tasks, including segmentation [32], object counting [46], anomaly detection [14], image editing [6], autonomous driving [22], and ap-

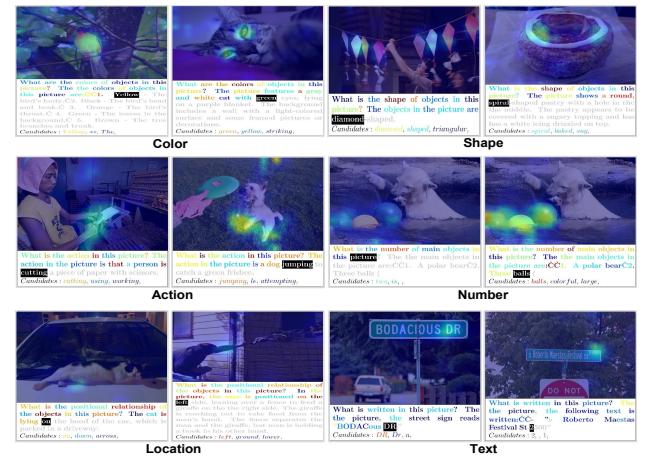


Figure 5. TAM supports to explain and analyze diverse attributes of MLLM (Qwen2-VL-7B). See extensive examples in Supp. M.

plications in the medical field [21, 50], which require pixel-level activations.

Visualization for Various Attributes. We apply TAM to explain various attributes of MLLMs, specifically focusing on fine-grained token types as Fig. 5. The results indicate that TAM effectively supports the explanation of attributes such as color, shape, action, number, location, and text. Notably, the representation of numbers is less pronounced compared to object localization, suggesting that some attributes may be encoded within related object tokens. In contrast, other attributes are more prominently featured.

Explaining Failure Cases. We apply TAM to analyze failure cases in the QK-VQA dataset [37] as Fig. 7. Our findings indicate that while the model successfully locates objects, it struggles to align them with specific knowledge. For instance, the model can identify the desert and focus on the key prompt (“year”) when outputting “don’t know.” It suggests a misalignment between the desert and the year, rather than recognition or prompt issues. Besides, some responses consist of synonyms, hypernyms, or hyponyms of the correct answers, resulting in mismatches. For more case studies, please refer to the catalog in Supp. Table 5.



Figure 6. TAM significantly improves the quality of video visualization, with much fewer redundant activation and noises on STAR [52] dataset using Qwen2-VL-2B [51]. Please see extensive examples in Supp. 27 and failure cases analysis for videos in Fig. 25.



Figure 7. TAM supports failure cases analysis on QK-VQA [37]. See more cases in Supp. about VQA (Fig. 24), videos (Fig. 25), visual reasoning (Fig. 28), and multi-turn conversation (Fig. 31).

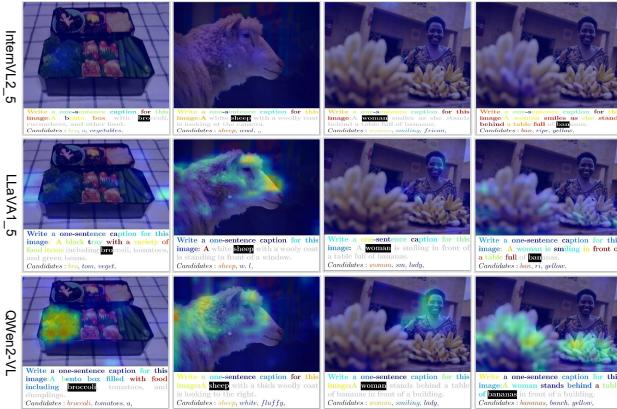


Figure 8. TAM enables qualitative comparison among MLLMs. Extensive examples are given in Fig. 17 and Fig. 18 in Supp. K.

Video Visualization. The proposed TAM excels at video visualization, where there are much fewer redundant activations and noises as shown in Fig. 6. Our method provides very clear results, while the conventional CAM makes it hard to see the video content without obvious highlights.

A Tool to Visually Compare MLLMs. TAM can be used as a tool to compare MLLMs qualitatively. If an MLLM exhibits superior activation maps, it can be inferred that its explainability and degree of vision alignment are also better, offering valuable pixel-level predictions at no

additional cost. As shown in Fig. 8, Qwen2-VL-7B [51] shows greater explainability than LLaVA1_5-7B [35] and InternVL2_5-8B [16]. Furthermore, we observe that InternVL tends to focus more on text, resulting in weaker activation levels; in some cases, this performance is even better than that of LLaVA (e.g., with the broccoli). More extensive examples can be found in Supp. K.

5. Conclusion

This work focuses on the unique nature of MLLMs in visual explainability, which generates multiple tokens progressively, complicating the explanation process. In response, we introduced the Token Activation Map (TAM), a novel method that utilizes estimated causal inference to mitigate interferences from context tokens, with a rank Gaussian filter to reduce activation noise, thereby providing clearer visual explanations. Our results demonstrate that TAM significantly outperforms existing SoTA methods, offering high-quality visualizations applicable to diverse scenarios.

Despite the success, we focus on visual inputs, other modalities, such as audio, remain under-explored. Additionally, interpreting model decisions is an extensible further aspect. The potential applications of our TAM method are extensive, including open-vocabulary segmentation, detection, image grounding, anomaly detection, remote sensing, and medical fields that require pixel-level activations. Our work lays a foundation for advancing their explainability and practical utility across various domains.

Acknowledgment

This work was supported by a research grant from the Joint Research Scheme (JRS) under the National Natural Science Foundation of China (NSFC) and the Research Grants Council (RGC) of Hong Kong (Project No. N_HKUST654/24), as well as a grant from the RGC of the Hong Kong Special Administrative Region, China (Project No. R6005-24).

References

- [1] Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4190–4197, 2020. 2, 5, 6, 7, 14, 15
- [2] Reduan Achitbat, Sayed Mohammad Vakilzadeh Hatefi, Maximilian Dreyer, Aakriti Jain, Thomas Wiegand, Sebastian Lapuschkin, and Wojciech Samek. Attnlrp: Attention-aware layer-wise relevance propagation for transformers. In *International Conference on Machine Learning*, pages 135–168. PMLR, 2024. 2, 5, 6, 7, 14, 15
- [3] Estelle Aflalo, Meng Du, Shao-Yen Tseng, Yongfei Liu, Chenfei Wu, Nan Duan, and Vasudev Lal. Vl-interpret: An interactive visualization tool for interpreting vision-language transformers. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 21406–21415, 2022. 2
- [4] 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. *Advances in neural information processing systems*, 35:23716–23736, 2022. 1
- [5] Ameen Ali, Thomas Schnake, Oliver Eberle, Grégoire Monfrat, Klaus-Robert Müller, and Lior Wolf. Xai for transformers: Better explanations through conservative propagation. In *International Conference on Machine Learning*, pages 435–451. PMLR, 2022. 6, 7, 14, 15
- [6] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kassten, and Tali Dekel. Text2live: Text-driven layered image and video editing. In *European conference on computer vision*, pages 707–723. Springer, 2022. 7
- [7] Gabriela Ben Melech Stan, Estelle Aflalo, Raanan Yehezkel Rohekar, Anahita Bhiwandiwalla, Shao-Yen Tseng, Matthew Lyle Olson, Yaniv Gurwicz, Chenfei Wu, Nan Duan, and Vasudev Lal. Lvlm-intrepert: An interpretability tool for large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8182–8187, 2024. 1, 2, 3
- [8] Chin-Chen Chang, Ju-Yuan Hsiao, and Chih-Ping Hsieh. An adaptive median filter for image denoising. In *2008 Second international symposium on intelligent information technology application*, pages 346–350. IEEE, 2008. 3, 5, 6, 16
- [9] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pages 839–847. IEEE, 2018. 2, 6, 7, 14, 15
- [10] Hila Chefer, Shir Gur, and Lior Wolf. Generic attention-model explainability for interpreting bi-modal and encoder-decoder transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 397–406, 2021. 2
- [11] Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 782–791, 2021. 2, 5, 15
- [12] Ruoyu Chen, Hua Zhang, Siyuan Liang, Jingzhi Li, and Xiaochun Cao. Less is more: Fewer interpretable region via submodular subset selection. In *The Twelfth International Conference on Learning Representations*. 3
- [13] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 2, 5, 6, 7, 13, 16, 17, 20, 21, 22, 23, 24, 26, 29, 30
- [14] Xuhai Chen, Jiangning Zhang, Guanzhong Tian, Haoyang He, Wuhao Zhang, Yabiao Wang, Chengjie Wang, and Yong Liu. Clip-ad: A language-guided staged dual-path model for zero-shot anomaly detection. In *International Joint Conference on Artificial Intelligence*, pages 17–33. Springer, 2024. 7
- [15] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 2, 13
- [16] Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024. 1, 6, 8, 16, 17, 19, 25, 35
- [17] Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, et al. A survey on multimodal large language models for autonomous driving. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 958–979, 2024. 1
- [18] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022. 6
- [19] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. In *The Twelfth International Conference on Learning Representations*. 2, 3, 4
- [20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019. 2
- [21] Xinpeng Ding, Yongqiang Chu, Renjie Pi, Hualiang Wang, and Xiaomeng Li. Hia: Towards chinese multimodal llms for comparative high-resolution joint diagnosis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 575–586. Springer, 2024. 7
- [22] Xinpeng Ding, Jianhua Han, Hang Xu, Xiaodan Liang, Wei Zhang, and Xiaomeng Li. Holistic autonomous driving understanding by bird’s view injected multi-modal large models.

- In 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 13668–13677, 2024. 7
- [23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. 1, 2
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 2
- [25] Aaron Hirst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 2
- [26] Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. Layercam: Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing*, 30:5875–5888, 2021. 2, 6, 14
- [27] Licheng Jiao, Yuhang Wang, Xu Liu, Lingling Li, Fang Liu, Wenping Ma, Yuwei Guo, Puhua Chen, Shuyuan Yang, and Biao Hou. Causal inference meets deep learning: A comprehensive survey. *Research*, 7:0467, 2024. 2, 3
- [28] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really learn. *Nature communications*, 10(1):1096, 2019. 2
- [29] Unggi Lee, Minji Jeon, Yunseo Lee, Gyuri Byun, Yooram Son, Jaeyoon Shin, Hongkyu Ko, and Hyeoncheol Kim. Llava-docent: Instruction tuning with multimodal large language model to support art appreciation education. *Computers and Education: Artificial Intelligence*, 7:100297, 2024. 1
- [30] Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18061–18070, 2024. 1
- [31] Yi Li, Hualiang Wang, Yiqun Duan, Hang Xu, and Xiaomeng Li. Exploring visual interpretability for contrastive language-image pre-training. *arXiv preprint arXiv:2209.07046*, 2022. 2
- [32] Yi Li, Hualiang Wang, Yiqun Duan, Jiheng Zhang, and Xiaomeng Li. A closer look at the explainability of contrastive language-image pre-training. *Pattern Recognition*, page 111409, 2025. 2, 3, 4, 5, 6, 7, 14, 15, 22
- [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 5
- [34] Fenglin Liu, Tingting Zhu, Xian Wu, Bang Yang, Chenyu You, Chenyang Wang, Lei Lu, Zhangdaihong Liu, Yefeng Zheng, Xu Sun, et al. A medical multimodal large language model for future pandemics. *NPJ Digital Medicine*, 6(1):226, 2023. 1
- [35] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 1, 2, 3, 6, 8, 13, 16, 17, 25
- [36] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017. 2, 3
- [37] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 5, 7, 8, 14, 18, 31, 33
- [38] Daniel Omeiza, Skyler Speakman, Celia Cintas, and Komminist Weldermariam. Smooth grad-cam++: An enhanced inference level visualization technique for deep convolutional neural network models. *arXiv preprint arXiv:1908.01224*, 2019. 1, 3
- [39] Vedant Palit, Rohan Pandey, Aryaman Arora, and Paul Pu Liang. Towards vision-language mechanistic interpretability: A causal tracing tool for blip. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2856–2861, 2023. 3
- [40] Judea Pearl. *Causality*. Cambridge university press, 2009. 2, 3
- [41] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1, 2, 22
- [42] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdellrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13009–13018, 2024. 5, 6, 7, 13, 16, 17
- [43] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “ why should i trust you?” explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016. 2, 3
- [44] Raanan Y Rohekar, Yaniv Gurwicz, and Shami Nisimov. Causal interpretation of self-attention in pre-trained transformers. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [45] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 1, 2, 6, 7, 14, 15, 16, 17, 18, 20, 24, 30, 34
- [46] Zenglin Shi, Ying Sun, and Mengmi Zhang. Training-free object counting with prompts. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 323–331, 2024. 7

- [47] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 1
- [48] Joseph Tey. Understanding how vision-language models reason when solving visual math problems. 1, 2
- [49] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 2, 3
- [50] Lehan Wang, Haonan Wang, Honglong Yang, Jiaji Mao, Zehong Yang, Jun Shen, and Xiaomeng Li. Interpretable bilingual multimodal large language model for diverse biomedical tasks. In *The Thirteenth International Conference on Learning Representations*, 2025. 7
- [51] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 1, 2, 5, 6, 7, 8, 13, 14, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38
- [52] Bo Wu, Shoubin Yu, Zhenfang Chen, Joshua B Tenenbaum, and Chuang Gan. Star: A benchmark for situated reasoning in real-world videos. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*. 2, 5, 8, 14, 18, 29, 32, 34
- [53] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. In *The Twelfth International Conference on Learning Representations*. 3
- [54] Liuyi Yao, Zhixuan Chu, Sheng Li, Yaliang Li, Jing Gao, and Aidong Zhang. A survey on causal inference. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(5):1–46, 2021. 3
- [55] Xiaofeng Zhang, Yihao Quan, Chen Shen, Xiaosong Yuan, Shaotian Yan, Liang Xie, Wenxiao Wang, Chaochen Gu, Hao Tang, and Jieping Ye. From redundancy to relevance: Enhancing explainability in multimodal large language models. *Annual Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics*, 2025. 1, 2
- [56] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016. 2, 6, 16, 17, 18, 20, 24, 30, 34
- [57] Zijian Zhou, Zheng Zhu, Holger Caesar, and Miaojing Shi. Openpsg: Open-set panoptic scene graph generation via large multimodal models. In *European Conference on Computer Vision*, pages 199–215. Springer, 2024. 2, 5, 6, 7, 14, 16, 17

Token Activation Map to Visually Explain Multimodal LLMs

Supplementary Material

Section and Textual Content	Graphic and Tabular Content
Supp. A: Tabular Catalog of the Supplementary.	Table 5: The catalog for quick reference.
Supp. B: Examples of Multimodal Activation Map	Fig. 9: How to read these visualization examples?
Supp. C: TAM to Explain All Generated Tokens.	Fig. 10: A complete example explaining all tokens.
Supp. D: Implementation of Baselines.	N/A, MLLM and explainability baselines .
Supp. E: Details of Metrics.	N/A, Invalid faithfulness evaluation for MLLM and other metric details.
Supp. F: Analysis About Motivation.	N/A, Statistical test and causal validation.
Supp. G: Extensive Cases About Method Comparison.	Fig. 12 & Fig. 13: TAM exceeds existing SoTA methods in extensive visualizations.
Supp. H: Visual Comparison About Causal Inference	Fig. 14: How does the estimated causal inference work in visualization?
Supp. I: Visual Comparison Among Denoise Filters	Fig. 15: How does the rank Gaussian filter work in visualization?
Supp. J: Visualization of Ablation Study.	Fig. 16: Two involved modules are mutually beneficial .
Supp. K: Explainability Results on Diverse MLLMs.	Table 6: MLLM quantitative results ; Fig. 11 relation between model size and explainability.
Supp. L: TAM for MLLM Visual Comparison.	Fig. 17 & Fig. 18: TAM supports visual comparison among MLLMs about attributes.
Supp. M: Extensive Cases About Attributes analysis.	Fig. 20 & Fig. 21 & Fig. 22 & Fig. 23: Explaining fine-grained attributes beyond SoTA.
Supp. N: TAM for Biased Scenario.	Fig. 19: TAM supports biased error analysis .
Supp. O: Extensive Failure Cases Study.	Fig. 24 & Fig. 25: TAM supports failure cases analysis for images and videos.
Supp. P: Extensive Success VQA Examples.	Fig. 26: Explanation result on the VQA dataset.
Supp. Q: Examples About Video Visualization.	Fig. 27: Clearer video visualizations with fewer redundant activations and noises.
Supp. R: Corner Case About Reasoning.	Fig. 28: TAM supports failure case analysis for visual reasoning .
Supp. S: TAM for Multi-image Conversation.	Fig. 29: High applicability on multi-image conversation.
Supp. T: TAM for Multi-turn Conversation.	Fig. 30 & Fig. 31: TAM supports multi-turn conversation about attributes and case study.

Table 5. Tabular catalog of the supplementary.

A. Tabular Catalog of the Supplementary

In this supplementary material, we primarily provide extensive qualitative results to demonstrate the effectiveness and wide applicability of TAM. These sections include comparisons with state-of-the-art (SoTA) methods, visualizations about ablation study, attribute explanation, failure case analysis, VQA examples, video visualizations, MLLM comparisons, reasoning analysis, multi-turn conversation, multi-image input, as well as some quantitative results and baseline descriptions. To enhance readability given the extensive content, we provide a tabular catalog in Table 5 for quick reference.

B. Examples of Multimodal Activation Map

In this section, we present a high-resolution example accompanied by detailed captions to facilitate the explanation of the multimodal activation map defined in Eq. 3. The primary element is the activation map at the top, which reflects the degree of vision-text alignment and serves to visually explain the MLLM. All multimodal activation maps in this paper adhere to a consistent format, and we provide high-quality images; please zoom in if any example appears too small to read.

The visual activations and textual relevances are normalized to the same scale as specified in Eq. 3, allowing for a direct comparison between the two modalities to identify where the model focuses—whether on the image or the context. The text is colored by tokens, with some words represented by multiple tokens marked in different colors. The answers following the target are not visible for the current explained token, and are therefore colored in gray. The colors of the candidate responses reflect the prediction confidence of the top three tokens corresponding to the target, which can be useful for analyzing failure cases through potential predictions and confidence levels associated with each token.

C. TAM to Explain All Generated Tokens

The proposed TAM demonstrates a significant advantage in explaining multiple generated tokens from MLLMs, in contrast to conventional models that typically focus on a single output. We depict all multimodal activation maps in Fig. 10 and support it. The figure clearly shows that TAM produces considerably fewer redundant activations, particularly for non-object words, thanks to the proposed estimated causal inference.

For instance, the activations for the function word “with” and the punctuation mark “.” in the Class Activation Map (CAM) are extremely high, overshadowing object activations. Additionally, these visual activations often exceed those of text tokens, resulting in excessive redundant activations. In comparison, our TAM is much clearer and fo-

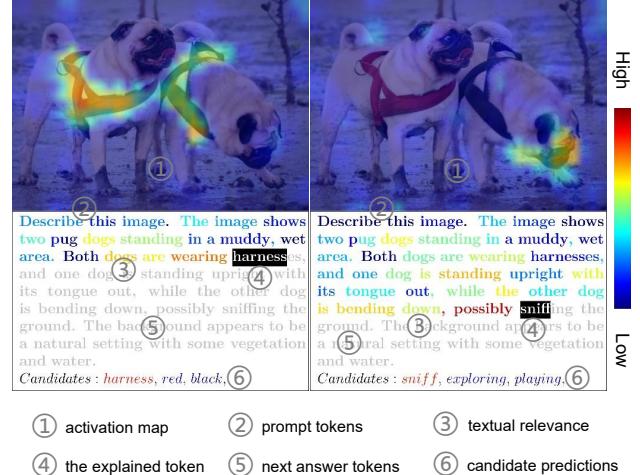


Figure 9. A high-resolution example of the multimodal activation map. This image is processed by the Qwen2-VL-2B model [51]. There are five components to visualize, including the visual activation map, prompt tokens, textual relevance, the explained target token, next answer tokens, and its top predictions (top 3). The colors indicate the corresponding response degree.

cuses more on important objects. These results suggest that TAM produces closer explanations to the understanding of humans than CAM, where words from the image are highlighted while those related to texts show much fewer responses. Besides, there are much fewer activations showing higher visualization quality as well.

D. Implementation of Baselines

We conduct experiments on various MLLMs, including Qwen2-VL-2B [51], Qwen2-VL-7B, LLaVA1.5-7B [35], LLaVA1.5-13B, InternVL2.5-2B [15], InternVL2.5-4B, and InternVL2.5-8B implemented by transformers using weights from huggingface. For Qwen2-VL all the models the weights are the Instruct version (e.g., Qwen2-VL-2B from huggingface “Qwen/Qwen2-VL-2B-Instruct”). The example of LLaVA1.5 weights is from “llava-hf/llava-1.5-7b-hf” and “thisisiron/InternVL2.5-2B” for InternVL2.5. Due to device limitations, very large MLLMs are not used. For image resolution, Qwen2-VL supports raw image size, while LLaVA1.5 and InternVL2.5 fix image sides at 336 and 448, respectively. For the implementation of video caption on Qwen2-VL [51], we extract 10 frames from a short video and repeat frames for the number of temporal_patch_size to ensure each frame has its own activation, instead of activations from other frames. We use the same prompts for the involved MLLMs. These prompts are set according to the average length of captions. For COCO Caption [13] the prompt is “Write a one-sentence caption for this image:”, and the prompts for GranDf [42] and

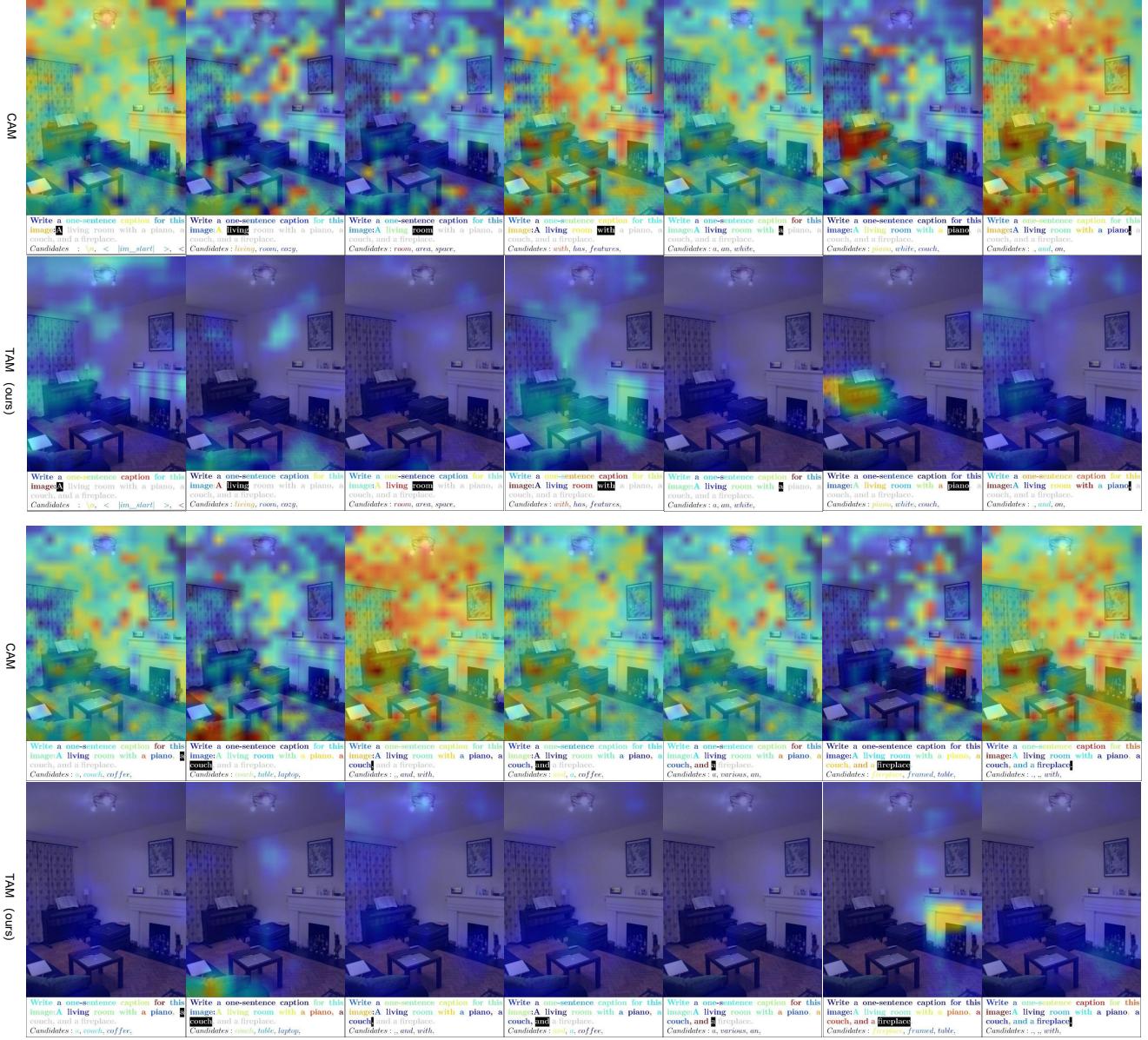


Figure 10. Visualization of one example for all generated tokens on the Qwen2-VL-2B [51] model. *The proposed TAM shows more accurate object localization ability beyond conventional CAM, with much fewer redundant activations in an overall view.*

OpenPSG [57] are “Write a description for this image using around two sentences:”, “Write a description for this image using around three sentences:”, respectively. For attribute analysis, the prompt is set like “What is the [attribute] in this picture?”. In addition, we use the prompts provided from the QK-VQA [37] and STAR [52] datasets, which vary according to the images.

For the explanation baselines [2, 9, 26, 32, 45], we implement them referring to their official codebases. To obtain attention weights for attention-based explainability methods

(such as attention weights and Rollout [1]), we rewrite the SdpaAttention in PyTorch. This is necessary because the original implementations of SdpaAttention and FlashAttention do not provide attention weights. For the CP-LRP [5] and AttnLRP [2], modules of MLLM are replaced by the official implementations of AttnLRP to back-propagate the relevance from output to visual tokens. Besides, these methods need to close the kv_cache to maintain gradients for key and value. Note, the Grad-CAM [45] is equivalent to CAM. Because the weights for the activation map are derived from

the classifier weight in MLLMs. Since there is only a feature vector without the pooling and other structures. The gradient is fully dependent on the classifier weight at the same ratio among channels. Since they are the same, we apply the classifier weight of CAM to achieve TAM to avoid extra back-propagation in Grad-CAM. Other implementation methods of estimated causal inference include the mean of context maps (mean), using the attention weights as relevance in Eq. 4 (AttnWeights).

E. Details of Metrics.

The metrics are based on the part of speech using the pos.tag function from the NLTK Python package. The specific tags are “NN”, “NNS”, “NNP”, and “NNPS” for Obj-IoU. Function words are identified by the tags: “CC”, “DT”, “EX”, “MD”, “POS”, “PRP”, “PRP\$”, “UH”, “WDT”, “WP”, “WP\$”, and “WRB”. Notably, we exclude the tags “IN” and “CD” from function words, as they pertain to location and quantity.

Besides the IoU based metric, in Sec. 3.4, we have discussed the difference between the used plausibility test (how accurately it reflects the true reasoning process), compared with another widely used faithfulness (how accurately it reflects the true reasoning process). While the faithfulness is not suitable for MLLM evaluation. Because the perturbation tests [2, 11] of the faithfulness metric alter the generated texts every time, resulting in inconsistent generated texts that are not stable to evaluate raw generated texts of MLLM. Specifically, masking different input regions in the faithfulness test drastically changes MLLM output tokens, invalidating observations of “decision-making” tied to a fixed class. In conventional models, input changes affect a fixed class’s confidence, but in MLLMs, it causes vanished tokens or shifted context, making confidence comparisons invalid. Besides, its cost is unacceptable, which needs N times repeated inferences ($N = \text{token number} \times \text{regions ratios}$).

Another consideration regarding the metrics is the variation in response levels. Specifically, our ECI involves a subtraction operation between activation maps, which can lead to a lower overall intensity compared to the original responses. We did not overlook this limitation when design the metrics; instead, we implemented a straightforward operation to penalize excessive discrepancies. Specifically, we use the response map of the first prompt token in place of the first generated token in evaluation. Since the first prompt token does not have any earlier text tokens, the map does not incorporate the ECI and reflects the original response level. If the response level of the altered map significantly differs from the processed maps, it can result in inappropriate background thresholds, thereby diminishing the Func-IoU metric. For instance, if the background threshold processed after ECI is 0.1, it may be too low for this map, leading to

false positives and consequently affecting the metrics. Detailed operations can be referenced in our open-source code.

F. Analysis About Motivation

In Fig. 1(c), we randomly pair CAMs and count their L1 distance against text correlation. Higher text correlation corresponds to lower distance, indicating concurrent interferences. In this section, we provide a statistical test to support it. Specifically, the added statistical test is the Pearson correlation at -0.16 with p-value of 1.5E-30. Since most pairs are not related in the random pairing, the correlation is not strong. When pairing the most related tokens, the Pearson correlation comes to -0.359 (p-value 7.9E-32). It confirms that the negative correlation is evident.

We also conducted a causal validation for the causal inference. In this paper, our ECI is based on the potential outcome model (POM). The used causal validation for this model is the Placebo test. Specifically, we validate it by replacing the target CAM to a random earlier CAM as the placebo (not the observed target), and then record the results drop. The Obj-IoU reduced to 6.2% on COCO Caption and 4.4 times lower than the raw result, suggesting the causal effect is significant.

G. Extensive Cases about Method Comparison

In addition to the visual comparison presented in Fig. 2, we offer more complex examples in Fig. 12 and Fig. 13 within this section. The findings are consistent with those discussed in Sec. 4.3: the proposed TAM significantly outperforms existing explainability methods.

Specifically, TAM generates fewer redundant activations and exhibits less noise compared to gradient-based methods [9, 45]. Moreover, it effectively locates objects, contrasting with the scattered activations seen in attention-based methods (e.g., Attention, Attention-Rollout [1], CP-LRP [5], AttnLRP [2]). These results indicate that TAM enhances the localization capabilities of MLLMs, even in complex scenarios. Consequently, TAM can be integrated into existing MLLMs without requiring grounding abilities, thereby facilitating a wide range of potential downstream tasks without additional supervision or alignment.

H. Visual Comparison About Causal Inference

We have validated the effectiveness of the proposed Estimated Causal Inference (ECI) in Table 2. In this section, we present visualization results that illustrate how our ECI outperforms existing methods and alternative implementations, as shown in Fig. 14. The first baseline we consider is feature surgery [32], which is designed to mitigate redundant features along the class dimension. However, the challenge with multi-language models (MLLMs) lies in the correlated activations along the token prediction dimension, which is

fundamentally different. As a result, feature surgery performs significantly worse than our ECI. Given the limited methods addressing correlated activations, we introduce additional baselines derived from other implementations of ECI: ECI-mean and ECI-attnWeights, as details in Supp. D. Although these suboptimal implementations outperform feature surgery, they still yield inferior results compared to the final ECI. Notably, our ECI demonstrates superior performance in handling function words, producing significantly fewer redundant activations while achieving better recall of target objects. These results indicate that our ECI is well-designed and effective for mitigating correlated activations among the generated tokens of MLLMs.

I. Visual Comparison Among Denoise Filters

Image denoising remains a traditional research topic, but it is the first time to be introduced in the visual explanation field. The issue of noise has been addressed in Sec. 2, where various methods aimed at noise reduction in transformers are discussed. However, residual noise persists even after these methods are applied. Consequently, it is essential to introduce denoising filters as a straightforward yet effective solution. Unlike conventional models that produce very small output sizes (e.g., 7×7), the output size of MLLMs is comparatively larger (e.g., 36×36). As a result, scatter-shaped noise is more likely to occur in MLLMs.

These noises belong to the salt-and-pepper noises in general, which can be effectively addressed using median and Gaussian filters. While these methods do not represent the optimal solution, as illustrated in Fig. 15. Specifically, the Gaussian filter proves inadequate in mitigating clustered noise, leaving many noises visible in the yellow boxes. The median filter reduces noise effectively, yet it still leaves behind unsolved scatter noise, with additional missing regions indicated by blue boxes in the final row. Similarly, the adaptive median filter [8] exhibits significant scatter noise, particularly near image edges where noise concentration is higher. In contrast, our proposed rank Gaussian filter demonstrates superior performance by amalgamating the strengths of both Gaussian and median filters, along with the novel technical enhancements discussed in Sec. 2.

J. Visualization of Ablation Study

We conducted ablation studies in Table 1. In addition to the quantitative results, we present further visualizations in Fig. 16 to elucidate the effectiveness of these modules. The first column showcases the baseline method, CAM [56] / Grad-CAM [45], which displays numerous redundant activations accompanied by noise, highlighted in white boxes. The proposed estimated causal inference (ECI) method in the second column effectively mitigates most correlated activations, although some persistent noise remains. The rank

Gaussian filter in the third column successfully removes this noise, but redundant activations are still evident. By integrating these two innovative techniques into the proposed TAM, we achieve substantial explanatory results that leverage the strengths of both approaches. These examples illustrate the mutual benefits of the modules, leading to an overall improvement that exceeds the sum of their individual contributions, as shown in Table 1.

K. Explainability Results on Diverse MLLMs

In addition to the explainability improvements highlighted in Table 4, we present specific results in Table 6. This table reveals that the overall F1-IoU of the TAM ranges from 30.68% to 41.45% across three datasets and seven MLLMs. In contrast, the baseline CAM [56], which is considered a SoTA method in terms of performance and practicality (as shown in Table 3), achieves F1-IoU results ranging from 23.63% to 34.39%. These results clearly indicate that TAM demonstrates broader applicability and enhanced explainability across diverse MLLMs.

Moreover, TAM offers a unique perspective on evaluating MLLMs from an explainability standpoint, beyond existing metrics. For instance, LLaVA models [35] and Qwen2-VL models [51] exhibit higher F1-IoU scores than InternVL models [16] on the COCO Caption [13] and OpenPSG [57] datasets, while InternVL models excel on the GranDf dataset [42].

We also observe scalability in explainability across certain model sizes, as shown in Fig. 11a-c. For example, LLaVA models show improvements from 7B to 13B, InternVL models from 2B to 4B, and Qwen2-VL models from 2B to 7B. This trend suggests a positive correlation between the scalability and explainability of MLLMs to a certain extent. When the model size are larger, the model tends to encode objects with fewer tokens, leading to a decrease in recall (see Fig. 11d) and an increase in precision. Subsequently, the Obj-IoU decreases due to a more significant decrease in recall.

L. TAM for MLLM Visual Comparison

Generally, researchers MLLMs using quantitative metrics or textual outputs, while visual comparisons remain underexplored. In contrast, visual evaluations are prevalent in conventional models such as Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs), which effectively illustrate the enhanced representational capabilities of new models. The absence of a dedicated explanation tool for MLLMs may contribute to this gap. The proposed TAM addresses this limitation, enabling researchers to conduct visual comparisons of their MLLMs against existing models, beyond the conventional textual comparison.

We present visual comparisons among MLLMs in Fig.

Method	MLLM	COCO Caption			Obj-IoU	Func-IoU	GranDf	OpenPSG	Obj-IoU	Func-IoU	F1-IoU
		Obj-IoU	Func-IoU	F1-IoU							
CAM	LLaVA1.5-7B [35]	23.17	43.16	30.16	20.07	47.48	28.21	25.11	51.55	33.77	
TAM	LLaVA1.5-7B [35]	27.65	61.43	38.13	20.71	59.15	30.68	28.57	61.06	38.93	
CAM	LLaVA1.5-13B [35]	24.82	51.18	33.43	21.34	43.99	28.74	26.65	48.45	34.39	
TAM	LLaVA1.5-13B [35]	29.12	58.5	38.88	22.1	51.02	30.84	30.88	59.96	40.76	
CAM	InternVL2.5-2B [16]	15.94	45.62	23.63	18.28	37.64	24.61	19.76	46.42	27.72	
TAM	InternVL2.5-2B [16]	21.38	65.1	32.19	20.48	85.93	33.08	23.0	86.86	36.36	
CAM	InternVL2.5-4B [16]	18.23	40.95	25.23	20.91	44.52	28.46	21.28	34.7	26.38	
TAM	InternVL2.5-4B [16]	21.76	63.12	32.36	22.53	89.71	36.02	23.49	89.75	37.23	
CAM	InternVL2.5-8B [16]	14.59	64.41	23.8	18.04	57.42	27.45	18.46	62.21	28.47	
TAM	InternVL2.5-8B [16]	19.98	66.53	30.73	21.56	85.95	34.47	21.73	88.74	34.91	
CAM	Qwen2-VL-2B [51]	21.23	51.93	30.14	17.85	62.15	27.74	22.93	48.5	31.15	
TAM	Qwen2-VL-2B [51]	27.37	68.44	39.1	18.65	88.97	30.83	26.26	92.99	40.95	
CAM	Qwen2-VL-7B [51]	22.51	42.44	29.42	18.6	68.03	29.21	23.41	42.94	30.3	
TAM	Qwen2-VL-7B [51]	28.13	71.85	40.43	19.88	90.57	32.61	26.94	89.88	41.45	

Table 6. TAM shows wide applicability on diverse MLLMs and datasets beyond the CAM [56] for all the experiments on the major F1-IoU (%) metric at large margins. TAM can be used as a visual comparison approach, where Qwen2-VL models [51] show better visual explainability than LLaVA1_6 [35] and InternVL2_5 [16] on the COCO Caption [13] and OpenPSG [57] datasets.

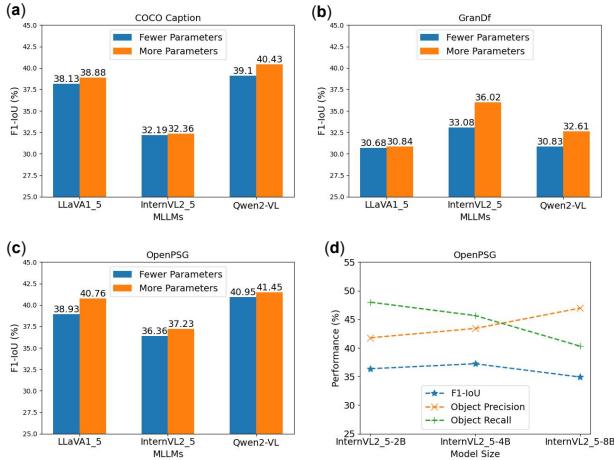


Figure 11. Relation between model size and explainability. (a-c) F1-IoU (%) results on LLaVA1.5 [35] (7B, 13B), InternVL2.5 [16] (2B, 4B), and Qwen2-VL [51] (2B, 7B) across diverse datasets [13, 42, 57] indicates the explainability is increased with more model parameters within a certain range. (d) Increasing the parameters of InternVL2.5 [16] on the OpenPSG dataset [57] improves object precision; however, this comes at the cost of decreased recall, which may negatively impact the F1-IoU score when the recall is too low.

17, focusing on attributes, and Fig. 18, which highlights object recognition. One key finding is that Qwen2-VL-7B [51] surpasses LLaVA1.5-7B [35] in text attributes, as illustrated in Fig. 17. Furthermore, it demonstrates superior performance in the object cases compared with LLaVA. For instance, Qwen2-VL-7B exhibits stronger activations in the top two rows of Fig. 18 and lacks correlation with the sandwich when interpreting the pizza in the third row.

In contrast, InternVL2.5-8B [16] tends to overemphasize textual information, resulting in weaker visual responses.

Although it generates highly detailed textual content that exceeds the length of LLaVA’s outputs, its activation map displays comparatively lower intensity. This is evidenced by the pronounced responses in textual areas (illustrated by the increased red regions). Despite reasonable activations in several successful cases (indicated by the light blue regions), InternVL2.5-8B exhibits lower quantitative performance, as shown in Table 6.

These visual findings underscore that TAM offers a novel perspective for comparing MLLMs, providing deeper insights, particularly when quantitative results are closely matched across certain datasets. Researchers can leverage visual comparisons to highlight the advantages of their models.

M. Extensive Cases about Attributes Analysis

The proposed TAM enables users to analyze the fine-grained attributes of MLLMs. These attributes contribute to a deeper understanding of how the model works. We present various visualizations of attributes, including actions and colors in Fig. 20, text and shapes in Fig. 21, and locations for both images and videos in Fig. 22.

The results indicate that the tested model, Qwen2-VL-7B [51], possesses the capability to comprehend diverse attributes with a high degree of explainability. Furthermore, we compare existing methods [45, 56] with our proposed TAM in Fig. 23, where our method demonstrates significantly superior explanation quality. These activation maps provide visual evidence for the generated content, thereby enhancing the model’s credibility.

N. TAM for Biased Scenario

The Task-Aware Mask (TAM) framework is capable of supporting the analysis of biased scenarios. In Fig. 19, we

investigate whether the background environment unexpectedly influences the classification of target categories. The output text indicates that images of terrestrial birds with synthetic aquatic backgrounds were misclassified as waterbirds, suggesting a significant bias introduced by background features in the model’s predictions. We conducted an in-depth analysis of this phenomenon using the TAM.

TAM effectively separates the contributions of different regions within the image to the classification decision, allowing for precise localization of the source of bias. Our research reveals that the synthetic aquatic background exerts a substantial influence on the model’s internal representations, leading it to favor categorizing images as waterbirds. This finding underscores the importance of considering background information during the model training and evaluation processes. Over-reliance on background features rather than the characteristics of the target itself may result in systematic misjudgments in scenarios that include synthetic or artificially manipulated backgrounds. The TAM-based analysis provides an effective diagnostic tool for identifying issues like background bias.



Figure 19. TAM supports analyzing biased scenarios. The landbirds in these two images were mistakenly classified as waterbirds due to the synthesized water backgrounds. The TAM identified that this biased recognition arises from the influence of the background.

O. Extensive Failure Cases Study

An important function of TAM is to support developers in analyzing failure cases, thereby deepening their understanding of the model’s shortcomings and enabling the development of better MLLMs. Generally, developers analyze errors by comparing the reply and answer, while TAM provides a clear visual view to understand them with more insights. As shown in Fig. 24, we list several failure cases with the error reason and corresponding analysis. We find that sometimes the model can successfully locate the target object, but lacks additional knowledge related to it thus replying falsely or refusing to answer (e.g., the train and cat in the left of Fig. 24). If the model focus on other regions out

of the target, the answer is possibly to be wrong. For example, we the model looks at the wall, it replies “living wall”, instead of the specific plant type the user asked for. Another error type is tolerable, that is synonyms, hypernyms, or hyponyms of answers (e.g., UK vs. England, fabric vs. nylon).

We further conduct case analysis on videos using Qwen2-VL-2B [51] in Fig. 25. Some error types are interesting. In the first row, we find the model already knows the object is a laptop when generating the token “pink”. But it turns to the case sequentially. It indicates the answer is shifted by context (maybe trained with some corpus including “pink case”). Besides, the representation is not strong enough, and the model cannot divide the pattern of the pillow and doll in the third row. In the fifth row, the picture with a border is similar to a book, while it is attached to the wall. From this context, we can know this is a picture instead of a book, indicating the weak capacity to integrate context. For the last row, the attention is located on the hair, suggesting the model predicts the “washing” according to the hair, instead of the window. All these examples prove that TAM can provide more cues and insights to analyze failure cases.

P. Extensive Success VQA Examples

In addition to the failure cases illustrated in Fig. 24, we present extensive success Visual Question Answering (VQA) examples in Fig. 26. These visualization results indicate that the Token Activation Model (TAM) is applicable not only to caption-based datasets but also to VQA datasets, such as QK-VQA [37]. From the figure, we observe that certain images are well-aligned with the generated tokens, which include objects, actions, texts, and patterns (e.g., the Qantas logo), thereby facilitating accurate predictions. However, some cases are not primarily object-determined; they rely heavily on textual cues, as seen with terms like “commercial” and “cross” in the last row. This analysis allows us to discern the sources of predictions based on activation levels: higher responses indicate strong visual relevance, while lower responses suggest a greater reliance on textual information.

Q. Examples about Video Visualization

Video modality is a crucial input type for MLLMs; however, it has seldom been studied in the explainability aspect. We compare our TAM with conventional methods [45, 56], as illustrated in Fig. 27, using Qwen2-VL-2B [51] on the STAR dataset [52] for video understanding. It is evident that TAM produces significantly clearer video visualization results compared to CAM [56] and Grad-CAM [45], both of which are well-established methods, as shown in Table 3. Specifically, TAM effectively reduces redundant activa-

tions and minimizes noise, allowing users to concentrate on target objects and observe the raw video more clearly. Additionally, we provide case studies in Fig. 25 for video error analysis.

R. Corner Case About Reasoning

TAM serves as a valuable tool for analyzing the visual reasoning processes of MLLM. In Fig. 28, we present a corner case of visual reasoning and analyze it using TAM. We find that both Qwen2-VL-7B [51] and InternVL2.5-8B [16] provided incorrect answers in this case. TAM reveals that the primary issue arises from incomplete recognition of all arrows in the problem, with missing arrows leading to subsequent reasoning errors. Furthermore, the reasoning capability of Qwen2-VL-7B is weaker than that of InternVL2.5-8B in this instance, as evidenced by inconsistent context understanding and very low text activations. In contrast, while InternVL2.5-8B exhibited stronger textual responses and correct logical reasoning, it compromised visual perception and interpretability, ultimately resulting in an incorrect answer as well. These analyses offer valuable insights for future model optimization and highlight the broad applicability of TAM.

S. TAM for Multi-image Conversation

Conventional models generally have a single input and a single output, whereas the characteristic of MLLM is that it supports multiple inputs and multiple token outputs. We provide examples of multi-image conversation in this chapter, as shown in Fig. 29. In the first example, Qwen2-VL-2B [51] can identify the main elements from four different images, accompanied by high-quality explanation results provided by TAM. In the second example, TAM explains from a visual perspective why Qwen2-VL-2B considers the fourth image the most interesting. TAM presents detailed activation maps, explaining specific focus points on attributes such as object (pug), artwork (The Scream), image style (cartoon), and painting style (anthropomorphic). These results demonstrate the broad applicability of TAM, which can support various new capabilities of MLLMs.

T. TAM for Multi-turn Conversation

TAM supports multi-turn conversation for MLLM, which is a new capability compared to conventional models. We first present a qualitative example in Fig. 30. Qwen2-VL-2B [51] can effectively generate the image description, and TAM provides accurate response maps for various attributes, such as objects, actions, and text. Subsequently, the user engaged in multi-turn conversation, inquiring about a fatter dog and the color of a chair. TAM effectively interpreted these fine-grained tokens, including positional information, adjectives, and colors. This example demon-

strates TAM’s broad applicability and offers strong interpretability analysis for new features like multi-turn dialogue in MLLM.

Additionally, we provide an analysis of a faulty example in Fig. 31, showing that TAM helps locate model errors and provides visual insight for developers. Although Qwen2-VL-2B can recognize why this image is distinctive and demonstrates strong interpretability for object tokens, it made errors in understanding speed and motion blur. Specifically, the taxi exhibited motion blur indicating higher speed, but it incorrectly identified it as an SUV. In the second round of dialogue, we speculated that it might not have recognized the blur, or it could have recognized the blur but failed to understand the relation between blur and speed. Thus, in the third round of dialogue, we asked which vehicle exhibited blur, and the clues provided by TAM indicated that the failure to recognize motion blur was the main reason for this faulty example. TAM supports multi-turn conversation, allowing for more detailed analysis of errors and offering developers precise reasons and a deeper understanding for issues.

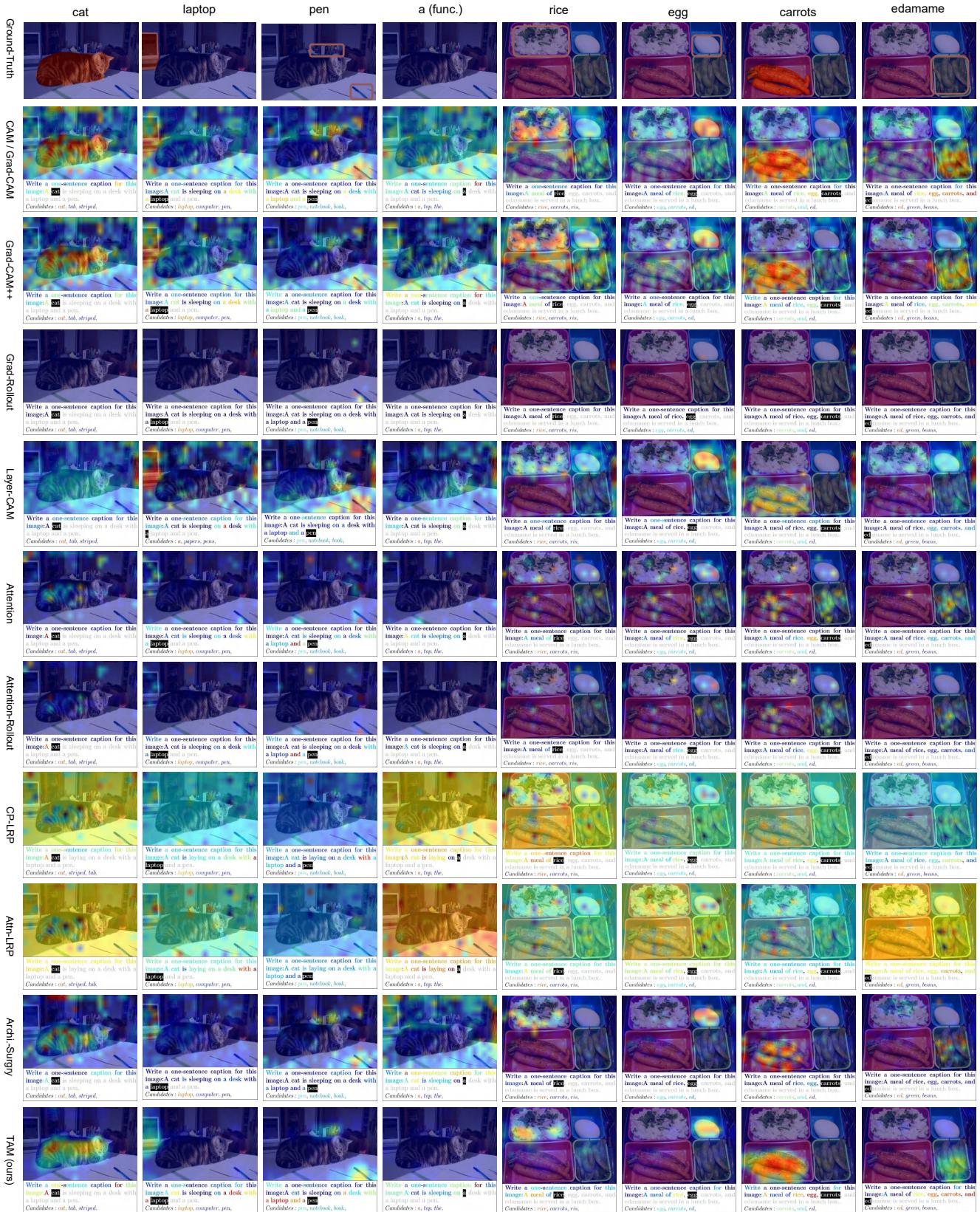


Figure 12. Visual comparison between our TAM and SoTA methods on the COCO Caption dataset [13] using the Qwen2-VL-2B [51] model. Objects without ground-truth are marked by red boxes. **TAM performs best beyond previous SoTA methods.** "CAM / Grad-CAM" indicates CAM [56] and Grad-CAM [45] are equivalent for MLLM as discussed in Supp. D.

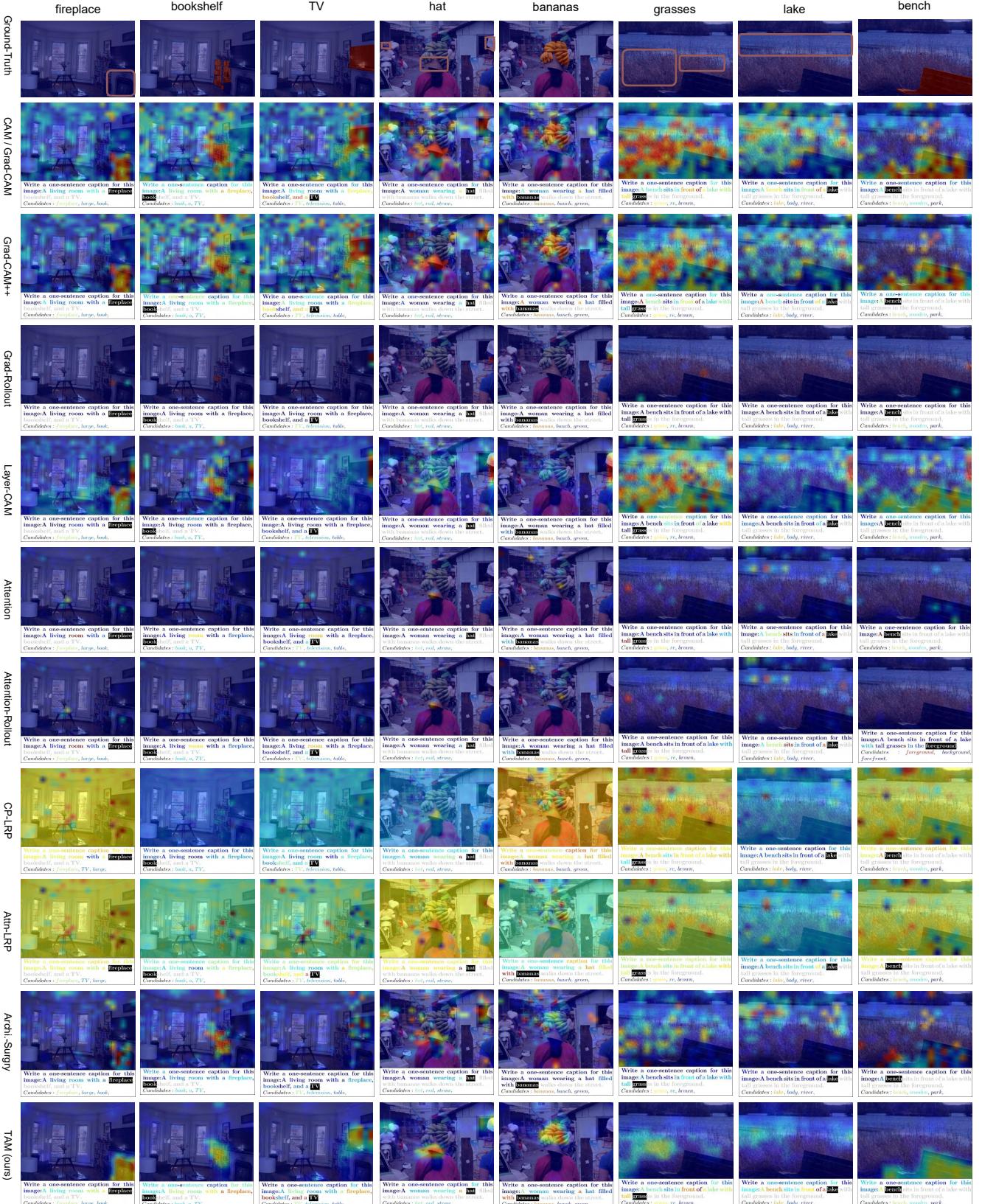


Figure 13. Visual comparison between our TAM and SoTA methods on the COCO Caption dataset [13] using the Qwen2-VL-2B [51] model. Objects without ground-truth are marked by red boxes. “func.” indicates function words assigned as the background class. **TAM** performs best beyond previous SoTA methods.

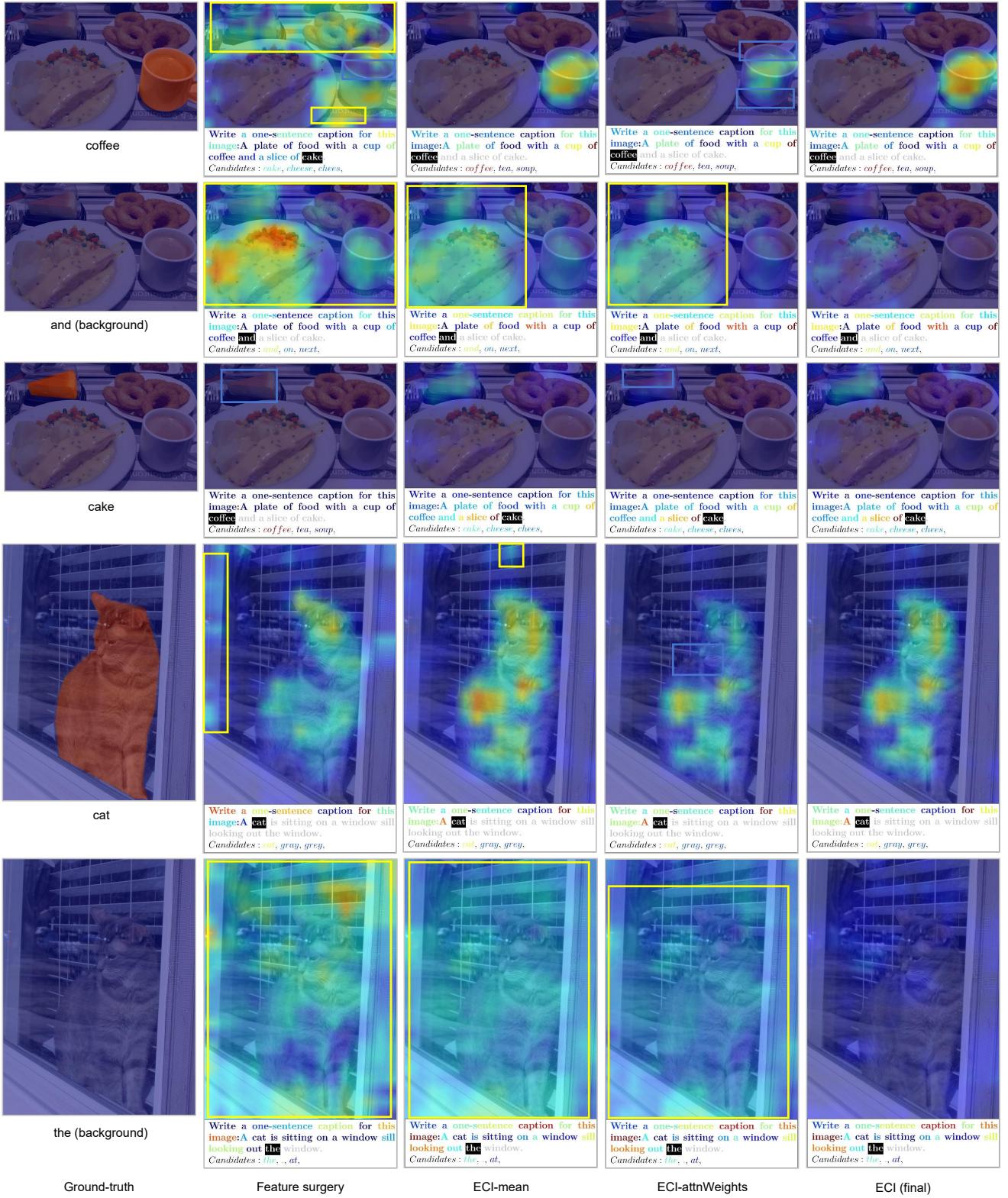


Figure 14. **The proposed estimated causal inference is well-designed beyond other methods and implementations.** The yellow boxes indicate correlated activations, and the blue boxes mean missed activations. Feature surgery [32] is designed for CLIP [41] to mitigate redundant features along the class dimension, while ECI-mean and ECI-attnWeights are other implementations of our estimated causal inference. The used model is Qwen2-VL-2B [51] on the COCO Caption dataset [13].



Figure 15. The proposed rank Gaussian filter is more effective than existing methods. The yellow boxes indicate insufficient denoising, and the blue boxes mean over-denoising. The used model is Qwen2-VL-2B [51] on the COCO Caption dataset [13].

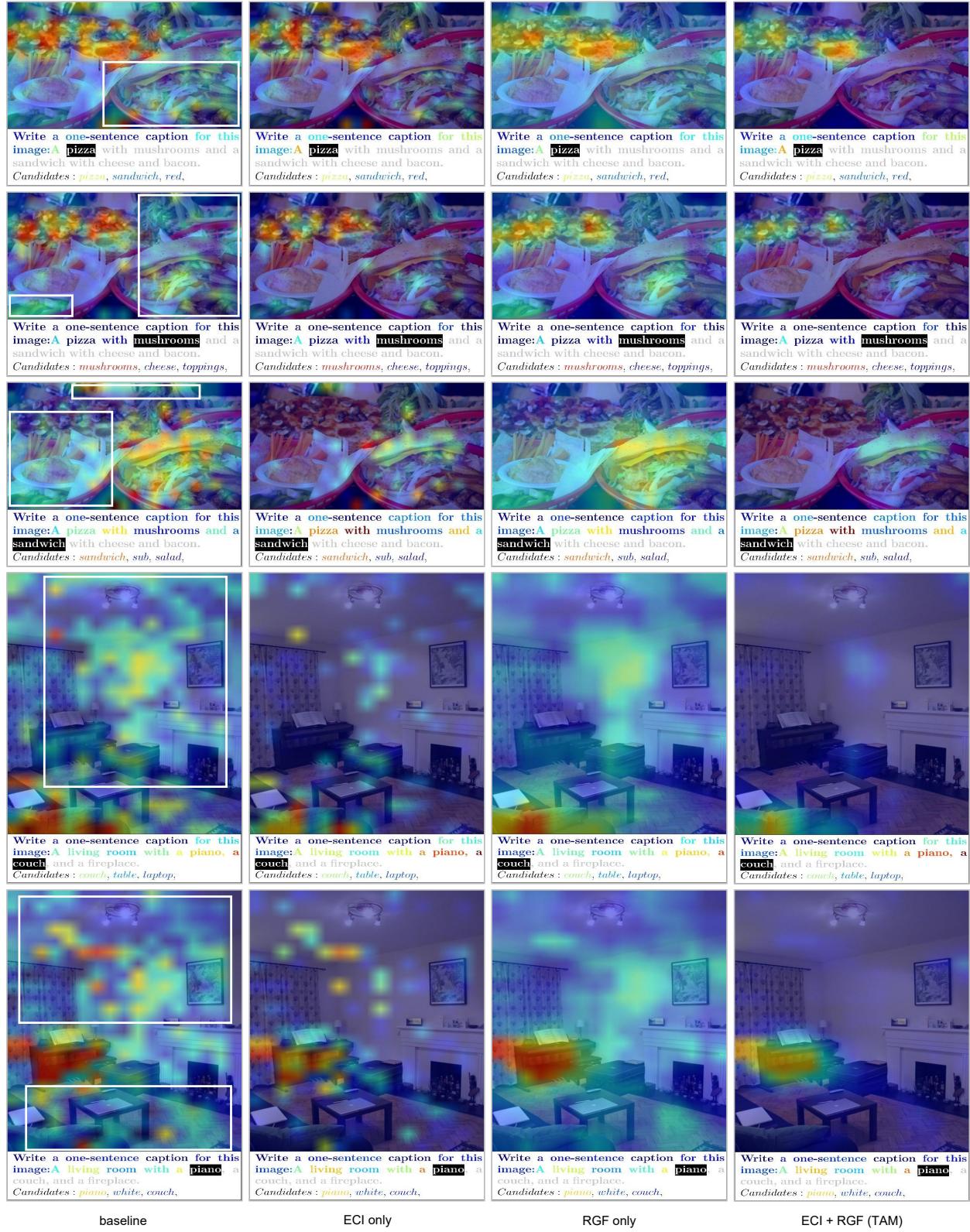


Figure 16. TAM involves two mutually beneficial modules. (ECI only) The estimated causal inference can mitigate most of the interference activations, while some stubborn noise remains. (RGF only) The rank Gaussian filter can remove these noises, while redundant activations are obvious. (ECI + RGF) TAM integrates the advantages of both and achieves good explanation results far beyond the baseline [45, 56]. The white boxes indicate correlated activations solved by our method. The used model is Qwen2-VL-2B [51] on the COCO Caption dataset [13].

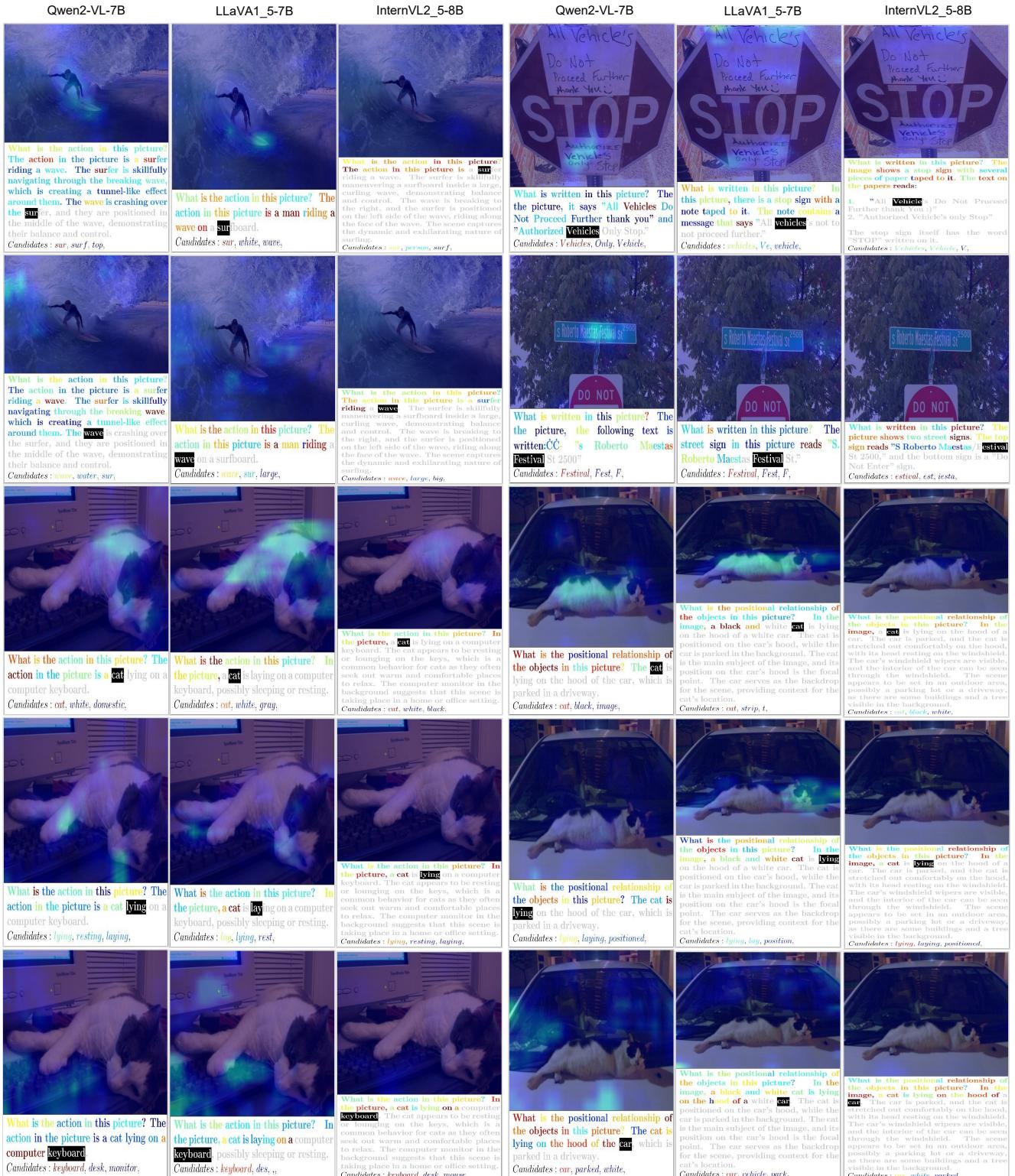


Figure 17. TAM supports visual comparison among MLLMs about attributes. Qwen2-VL-7B [51] presents good visual explainability beyond LLaVA1_5-7B [35] on texts. InternVL2_5-8B [16] focuses on textual content with more red texts and weaker visual activations.

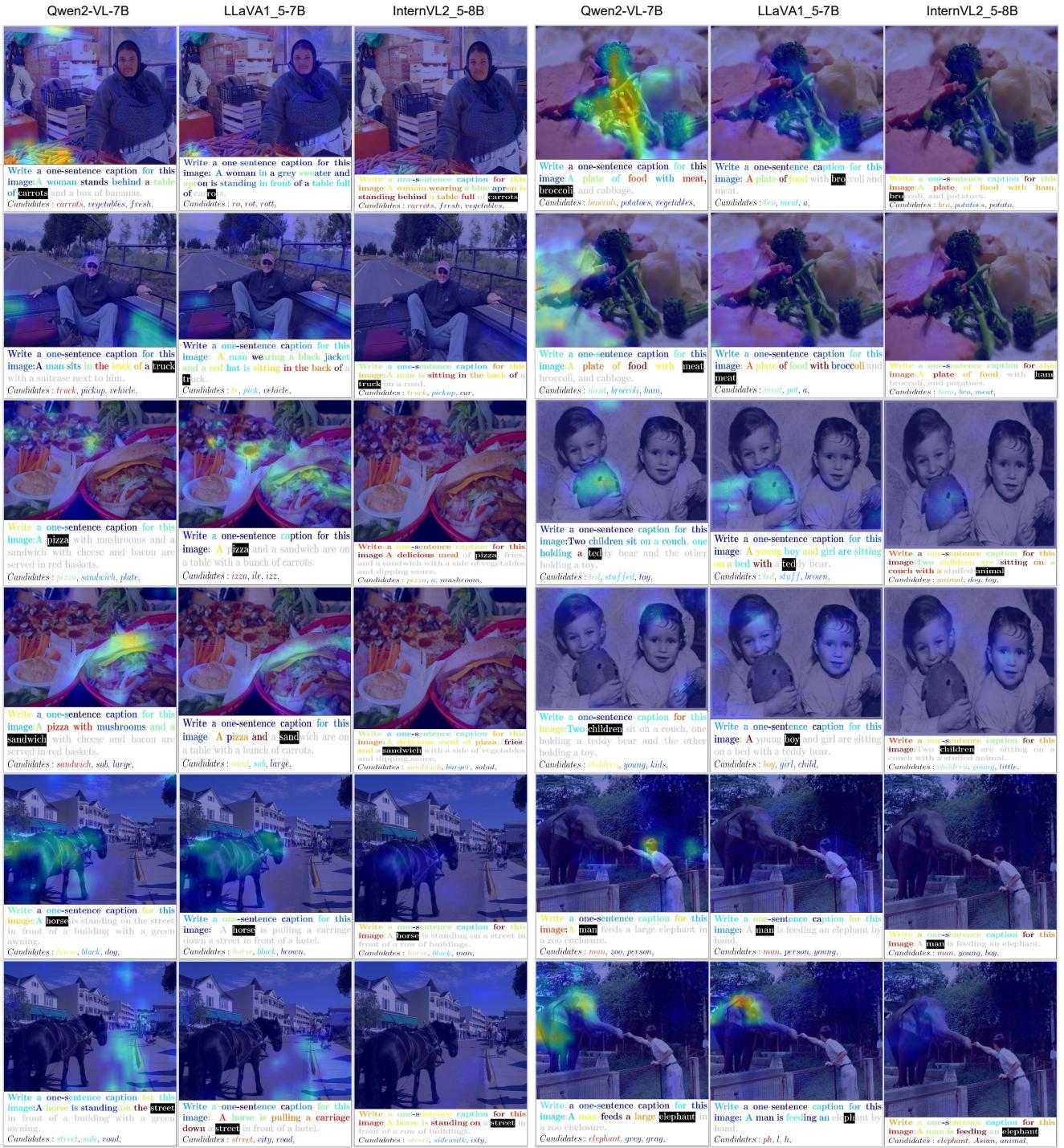


Figure 18. TAM supports visual comparison among MLLMs about objects on the COCO Caption [13] dataset. Qwen2-VL-7B [51] presents the best visual results with less correlation (e.g., pizza vs. sandwich in the third row) and shows a stronger activation degree.

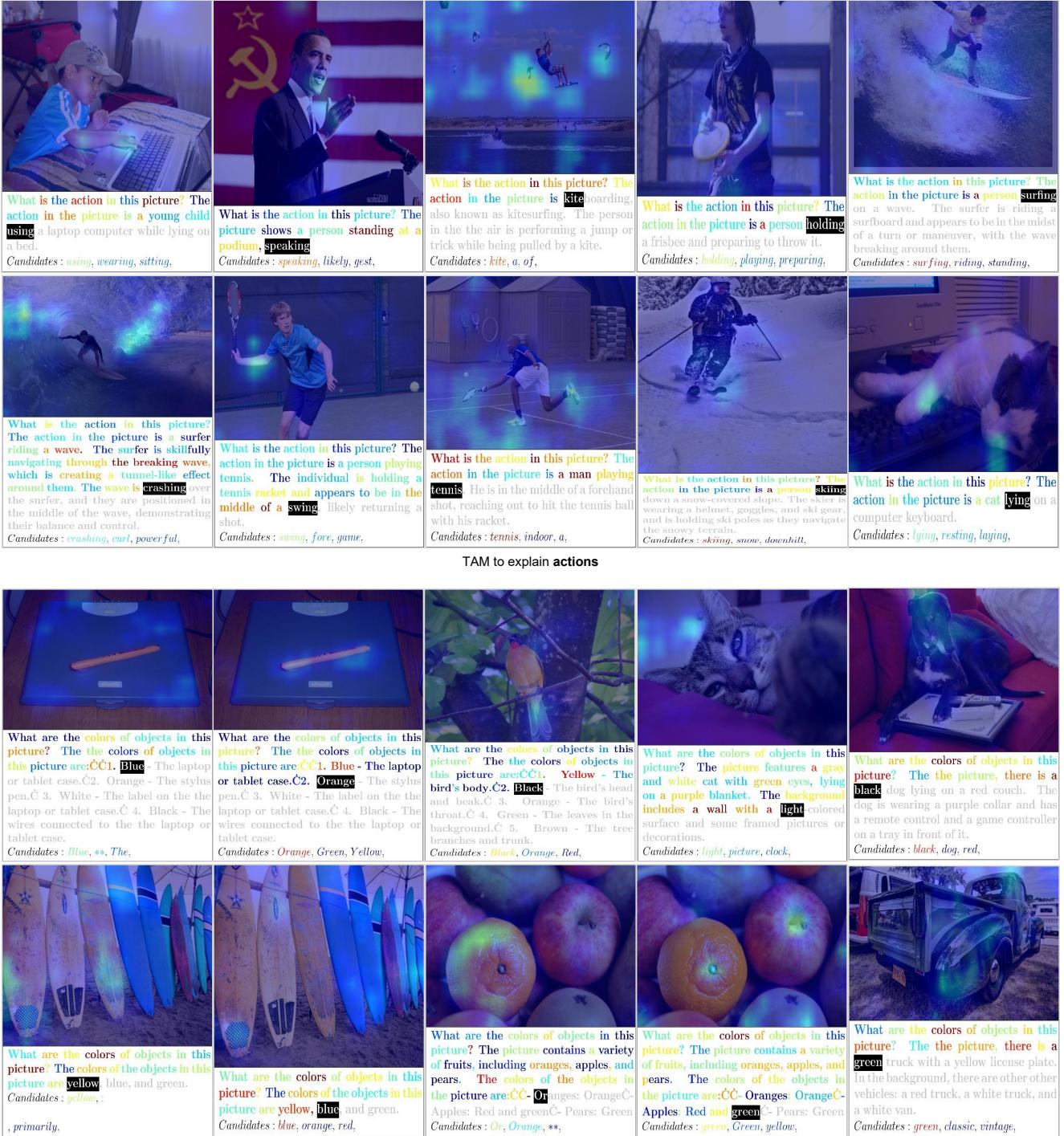
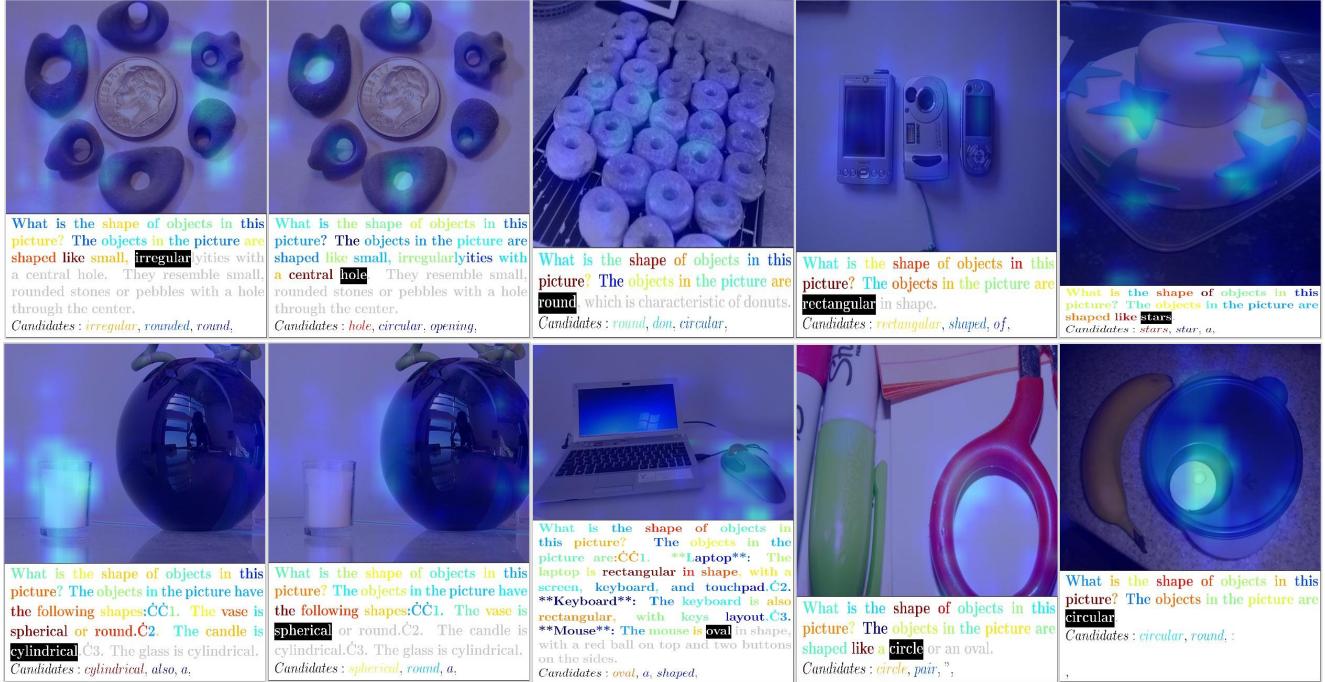


Figure 20. TAM supports explaining attributes of MLLMs at high-quality for the Qwen2-VL-7B [51] about action and colors.



TAM to explain text



TAM to explain shape

Figure 21. TAM supports explaining attributes of MLLMs at high-quality for the Qwen2-VL-7B [51] about text and shape.



TAM to explain location in image



TAM to explain location in video

Figure 22. TAM supports explaining attributes of MLLMs at high-quality for both images and videos. The images are processed by the Qwen2-VL-7B [51] from the COCO Caption dataset [13] and we use the Qwen2-VL-2B for videos from the STAR dataset [52].

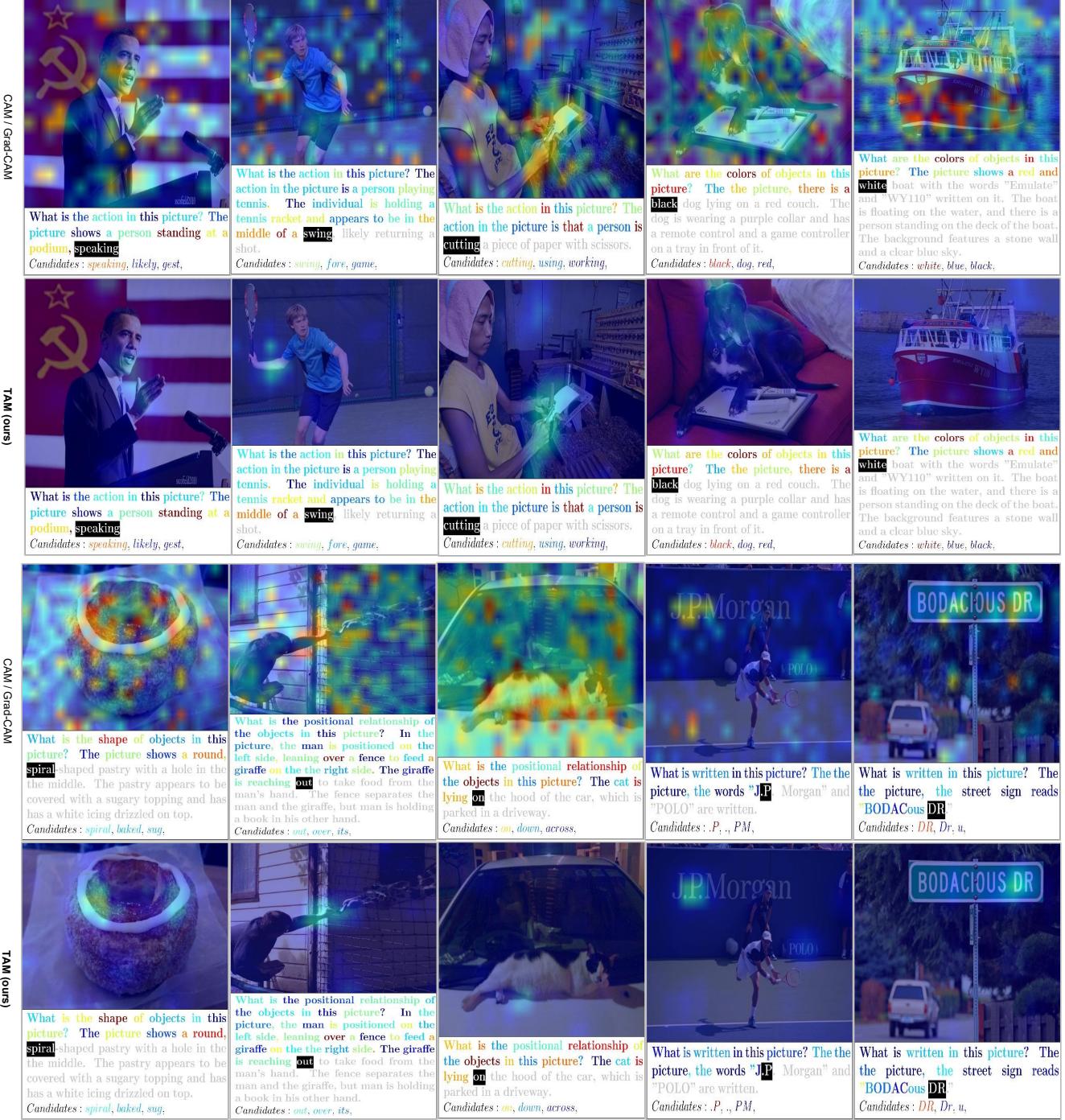


Figure 23. **TAM exceeds existing methods on attribute explanation.** The images are processed by the Qwen2-VL-7B [51] from the COCO Caption dataset [13] compared with the baseline. "CAM / Grad-CAM" indicates CAM [56] and Grad-CAM [45] are equivalent for MLLM, as discussed in Supp. D.

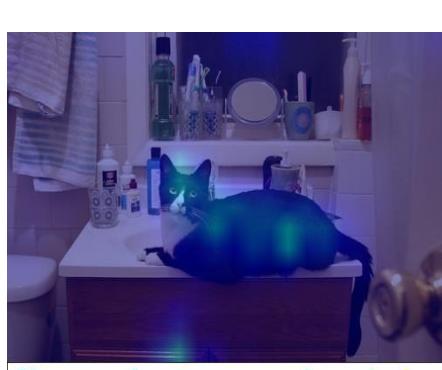
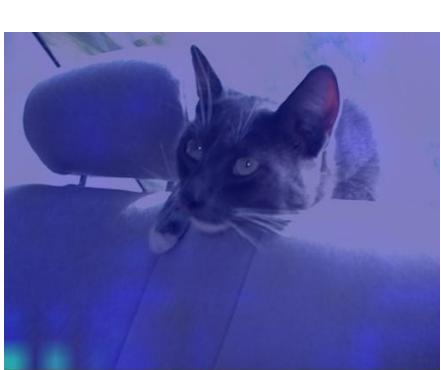
		
<p>How fast is the vehicle pictured capable of reaching? I don't know how fast the vehicle pictured is capable of reaching. It could be slow, medium, or fast.</p> <p>Candidates : <i>vehicle</i>, <i>train</i>, <i>image</i>,</p> <p>Answers: 80 mph, 10mph, 100 mph</p>	<p>What other big vehicle is often painted about the same shade as this vehicle? The other big vehicle often painted about the same shade as this vehicle is a train</p> <p>Candidates : <i>train</i>, <i>locom</i>, <i>bus</i>,</p> <p>Answers: fire truck, firetruck</p>	<p>Name the type of plant this is? The plant in the image appears to be a type of indoor plant commonly known as a "living wall" or "vertical garden." These plants are typically grown on a frame or support structure and can be used to add greenery and aesthetic appeal to a space. They can be used in various indoor settings, including bathrooms, to improve air quality and provide a natural touch.</p> <p>Candidates : <i>living</i>, <i>vertical</i>, <i>green</i>,</p> <p>Answers: climb, vine, look like some kind of ivy, ficus</p>
<p>Reson: lack of additional knowledge</p> <p>Details: the train is located, but lack of additional knowledge about its speed.</p>	<p>Reson: additional knowledge / question understanding</p> <p>Details: the train is recognized, but the question is other vechicle. Missed question understanding and lack of additional knowledge are both possible.</p>	<p>Reson: weak vision-text alignment</p> <p>Details: the plant is not well aligned with the specific plant type. Besides, it focuses on the wall according to the activation, thus reply the "living wall".</p>
		
<p>Name one famous person whom also has a black and white one of these? There is no existence of a famous person with a black and white cat in the image description.</p> <p>Candidates : <i>cat</i>, <i>one</i>, <i>pet</i>,</p> <p>Answers: ernest hemingway, hillary clinton, taylor swift, russell brand</p>	<p>What country does this appear to be? The image appears to be from the United Kingdom. The style of the car, the type of sheep, and the overall landscape are characteristic of the British countryside.</p> <p>Candidates : <i>United</i>, <i>UK</i>, <i>British</i>,</p> <p>Answers: scotland, usa, england</p>	<p>Name the material used to make this car seat shown in this picture? The car seat shown in the picture appears to be made of a fabric material. The texture and color suggest it could be a type of upholstery fabric commonly used in cars.</p> <p>Candidates : <i>fabric</i>, <i>material</i>, ,,</p> <p>Answers: leather, nylon, cloth, polyester</p>
<p>Reson: lack of additional knowledge</p> <p>Details: the cat is located, but lack of additional knowledge about famous person with it.</p>	<p>Reson:synonyms of answers</p> <p>Details: the answer is right, but the answer england is not exactly the same to the United Kingdom.</p>	<p>Reson:synonyms of answers</p> <p>Details: even though the fabric is well located for the right one, but the answer is not specific enough.</p>

Figure 24. TAM supports failure case analysis for deeper understanding with error reason and analysis details using Qwen2-VL-2B [51] on the QK-VQA dataset [37].

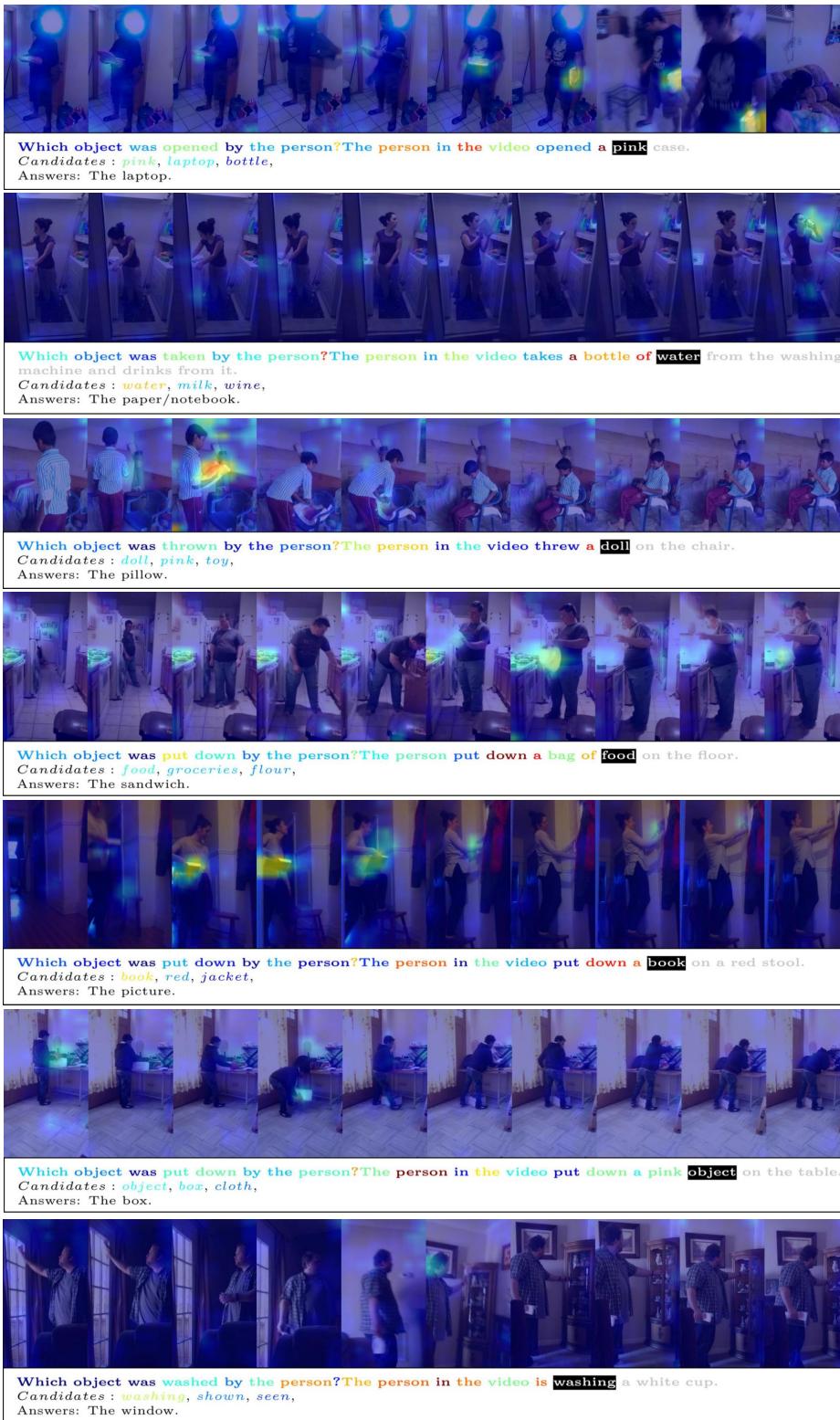
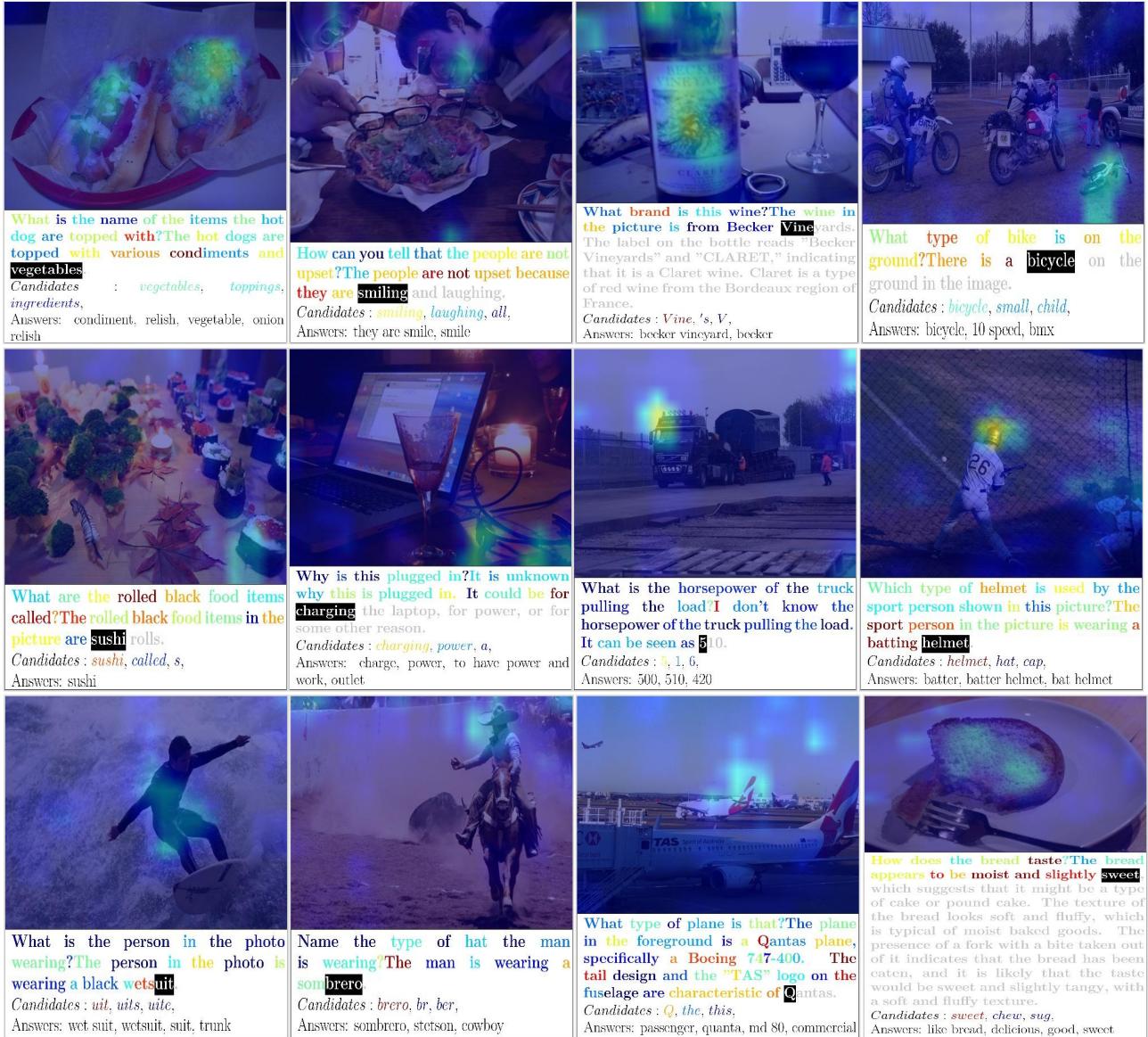


Figure 25. TAM supports failure case analysis for videos with error reason and analysis details using Qwen2-VL-2B [51] on the STAR dataset [52].



Object-determined success cases with strong activations



Textual-determined success cases with weaker activations

Figure 26. TAM presents good visual explanation result for the VQA dataset with extensive successful examples on the QK-VQA dataset [37] using Qwen2-VL-2B [51]. These cases are dependent on different information, divided into “Object-determined” type and “Textual-determined” type, with higher and lower activation degrees, respectively.

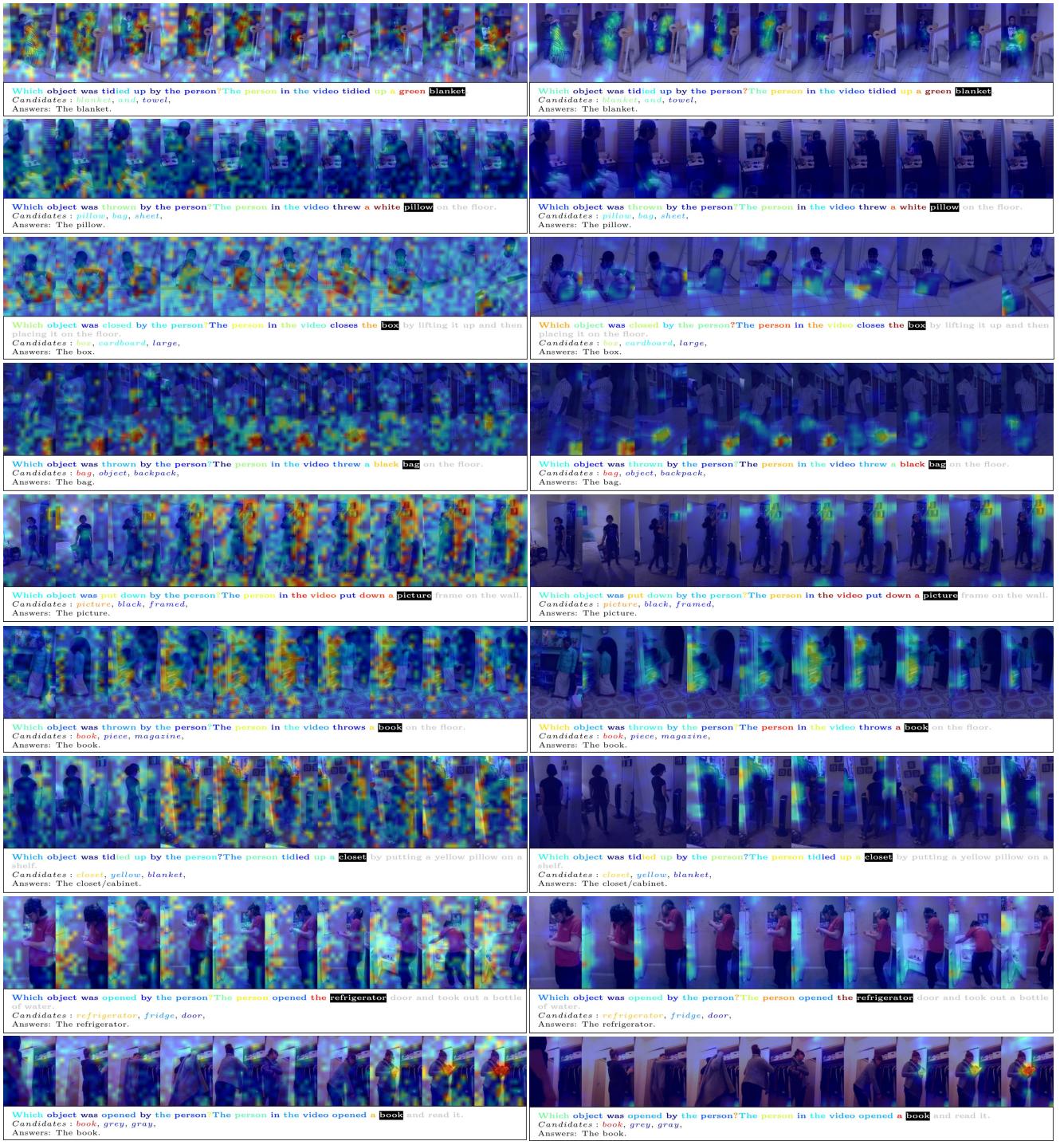
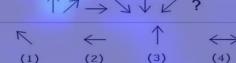
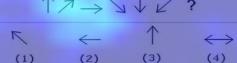
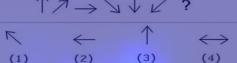
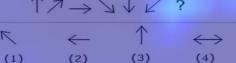


Figure 27. Comparison about video visualization between our TAM and CAM [56] / Grad-CAM [45] (equivalent to MLLM) on the STAR dataset [52] using Qwen2-VL-2B [51]. **TAM presents much clearer visualization results with fewer redundant activations and noises.**

Qwen2-VL-7B

| $\uparrow \nearrow \rightarrow \downarrow \downarrow \swarrow ?$ |
|---|---|---|---|---|
| 
(1) (2) (3) (4) | 
(1) (2) (3) (4) | 
(1) (2) (3) (4) | 
(1) (2) (3) (4) | 
(1) (2) (3) (4) |
| Please think step by step and reply which answer is right. To solve this problem, we need to identify the pattern in the sequence of arrows. Let's analyze the given sequence step by step:
1. The first arrow is an upward arrow.
2. The second arrow is a right-headed arrow pointing to the right.
3. The third arrow is a downward arrow.
4. The fourth arrow is a question mark. | Please think step by step and reply which answer is right. To solve this problem, we need to identify the pattern in the sequence of arrows. Let's analyze the given sequence step by step:
1. The first arrow is an upward arrow.
2. The second arrow is a right-headed arrow pointing to the right.
3. The third arrow is a downward arrow.
4. The fourth arrow is a question mark. | Please think step by step and reply which answer is right. To solve this problem, we need to identify the pattern in the sequence of arrows. Let's analyze the given sequence step by step:
1. The first arrow is an upward arrow.
2. The second arrow is a right-headed arrow pointing to the right.
3. The third arrow is a downward arrow.
4. The fourth arrow is a question mark. | Please think step by step and reply which answer is right. To solve this problem, we need to identify the pattern in the sequence of arrows. Let's analyze the given sequence step by step:
1. The first arrow is an upward arrow.
2. The second arrow is a right-headed arrow pointing to the right.
3. The third arrow is a downward arrow.
4. The fourth arrow is a question mark. | Please think step by step and reply which answer is right. To solve this problem, we need to identify the pattern in the sequence of arrows. Let's analyze the given sequence step by step:
1. The first arrow is an upward arrow.
2. The second arrow is a right-headed arrow pointing to the right.
3. The third arrow is a downward arrow.
4. The fourth arrow is a question mark. |
| Now, let's look at the options provided:
1. An upward arrow.
2. The second arrow is a two-headed arrow pointing to the left.
3. An upward arrow.
4. A two-headed arrow pointing to the right. | Now, let's look at the options provided:
1. An upward arrow.
2. The second arrow is a two-headed arrow pointing to the left.
3. An upward arrow.
4. A two-headed arrow pointing to the right. | Now, let's look at the options provided:
1. An upward arrow.
2. The second arrow is a two-headed arrow pointing to the left.
3. An upward arrow.
4. A two-headed arrow pointing to the right. | Now, let's look at the options provided:
1. An upward arrow.
2. The second arrow is a two-headed arrow pointing to the left.
3. An upward arrow.
4. A two-headed arrow pointing to the right. | Now, let's look at the options provided:
1. An upward arrow.
2. The second arrow is a two-headed arrow pointing to the left.
3. An upward arrow.
4. A two-headed arrow pointing to the right. |
| The pattern seems to alternate between upward and downward arrows, and then a question mark. The next arrow should follow the same pattern. | The pattern seems to alternate between upward and downward arrows, and then a question mark. The next arrow should follow the same pattern. | The pattern seems to alternate between upward and downward arrows, and then a question mark. The next arrow should follow the same pattern. | The pattern seems to alternate between upward and downward arrows, and then a question mark. The next arrow should follow the same pattern. | The pattern seems to alternate between upward and downward arrows, and then a question mark. The next arrow should follow the same pattern. |
| Given this, the correct answer is the upward arrow. | Given this, the correct answer is the upward arrow. | Given this, the correct answer is the upward arrow. | Given this, the correct answer is the upward arrow. | Given this, the correct answer is the upward arrow. |
| So, the final answer is (1). Candidates: <i>upward, up, ↑</i> . | So, the final answer is (1). Candidates: <i>-headed, -sided, -way,</i> | So, the final answer is (1). Candidates: <i>third, fourth, second,</i> | So, the final answer is (1). Candidates: <i>question, two, left,</i> | So, the final answer is (1). Candidates: <i>↑, < in_end>,</i> |
| correct activations | missing the second arrow; read the third arrow instead | down arrow is the fifth arrary; not the third | question mark is the seventh symbol; three arrow are missing | weak context response, the answer "1" is not consistent to the context "upward arrow" |

InternVL2_5-8B

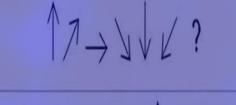
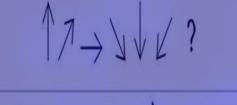
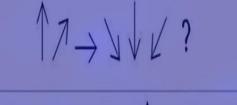
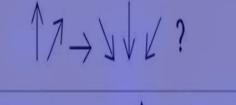
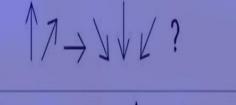
| $\uparrow \nearrow \rightarrow \downarrow \downarrow \swarrow ?$ |
|---|---|---|---|---|
| 
(1) (2) (3) (4) | 
(1) (2) (3) (4) | 
(1) (2) (3) (4) | 
(1) (2) (3) (4) | 
(1) (2) (3) (4) |
| Please think step by step and reply which answer is right. To solve this puzzle, we need to identify the pattern in the sequence of arrows and apply it to find the missing arrow.
1. The first arrow is pointing up.
2. The second arrow is pointing right.
3. The third arrow is pointing down.
4. The fourth arrow is pointing left. | Please think step by step and reply which answer is right. To solve this puzzle, we need to identify the pattern in the sequence of arrows and apply it to find the missing arrow.
1. The first arrow is pointing up.
2. The second arrow is pointing right.
3. The third arrow is pointing down.
4. The fourth arrow is pointing left. | Please think step by step and reply which answer is right. To solve this puzzle, we need to identify the pattern in the sequence of arrows and apply it to find the missing arrow.
1. The first arrow is pointing up.
2. The second arrow is pointing right.
3. The third arrow is pointing down.
4. The fourth arrow is pointing left. | Please think step by step and reply which answer is right. To solve this puzzle, we need to identify the pattern in the sequence of arrows and apply it to find the missing arrow.
1. The first arrow is pointing up.
2. The second arrow is pointing right.
3. The third arrow is pointing down.
4. The fourth arrow is pointing left. | Please think step by step and reply which answer is right. To solve this puzzle, we need to identify the pattern in the sequence of arrows and apply it to find the missing arrow.
1. The first arrow is pointing up.
2. The second arrow is pointing right.
3. The third arrow is pointing down.
4. The fourth arrow is pointing left. |
| The pattern seems to be a clockwise rotation of the arrow direction:
- Up → Right → Down → Left | The pattern seems to be a clockwise rotation of the arrow direction:
- Up → Right → Down → Left | The pattern seems to be a clockwise rotation of the arrow direction:
- Up → Right → Down → Left | The pattern seems to be a clockwise rotation of the arrow direction:
- Up → Right → Down → Left | The pattern seems to be a clockwise rotation of the arrow direction:
- Up → Right → Down → Left |
| Following this pattern, the next arrow should be pointing up again. | Following this pattern, the next arrow should be pointing up again. | Following this pattern, the next arrow should be pointing up again. | Following this pattern, the next arrow should be pointing up again. | Following this pattern, the next arrow should be pointing up again. |
| So, the correct answer is (3). Candidates: <i>up, upwards, upward,</i> | So, the correct answer is (3). Candidates: <i>right, left, to,</i> | So, the correct answer is (3). Candidates: <i>down, up, left,</i> | So, the correct answer is (3). Candidates: <i>left, right, down,</i> | So, the correct answer is (3). Candidates: <i>3, 1, 4,</i> |
| weak visual responses | missing the second arrow; read the third arrow instead | down arrow is the fifth arrary; not the third | wrong perception with very weak activation map | strong context reasoning (Up -> Right -> Down -> Left -> Up), but the visual perception is weak |

Figure 28. Visual reasoning corner case analyzed using TAM. The analysis reveals that both Qwen2-VL-7B [51] and InternVL2_5-8B [16] incorrectly answered the question due to incomplete recognition of all the arrows. Missing arrows led to erroneous reasoning. Qwen2-VL-7B demonstrated weaker reasoning capabilities compared to InternVL2_5-8B, exhibiting inconsistent context understanding and low text activations. In contrast, while InternVL2_5-8B provided a stronger text response with correct reasoning logic, its visual perception and activation degree are weaker.



Figure 29. TAM supports multi-image conversation of MLLM showing wide applicability. It generates high-quality activation maps for four input images in the first case using the Qwen2-VL-2B [51]. In the second case, TAM explains why the model regards the last image as the most interesting one, providing visual cues for key tokens. These results showcase the model’s effectiveness in multi-image scenarios, highlighting its versatility in handling complex visual data.

Multi-turn Conversation

The image shows two pugs sitting in colorful, round chairs. The chairs are designed with a playful and whimsical pattern, featuring the word "Peppa" on the fabric. The pugs are sitting on a wooden floor, and the background includes various household items such as a microwave, a food storage container, and a shelf with various containers and boxes. The overall setting appears to be a kitchen or a similar living space.

Which pug is fatter?

The pug on the left is fatter.

What color is the chair that this fatter dog is sitting on?

The chair that the fatter dog is sitting on is blue.

Token Activation Maps

Figure 30. TAM supports multi-turn conversation of MLLM for diverse attributes. TAM presents high-quality visual explanation results on Qwen2-VL-2B [51] regarding attributes such as objects, actions, and text in the first round. Then, the user inquires about the fatter dog and the color of a chair in the second and third rounds, respectively. Activation maps suggest TAM is capable of explaining fine-grained tokens like positional information, adjectives, and colors in multi-turn conversation.

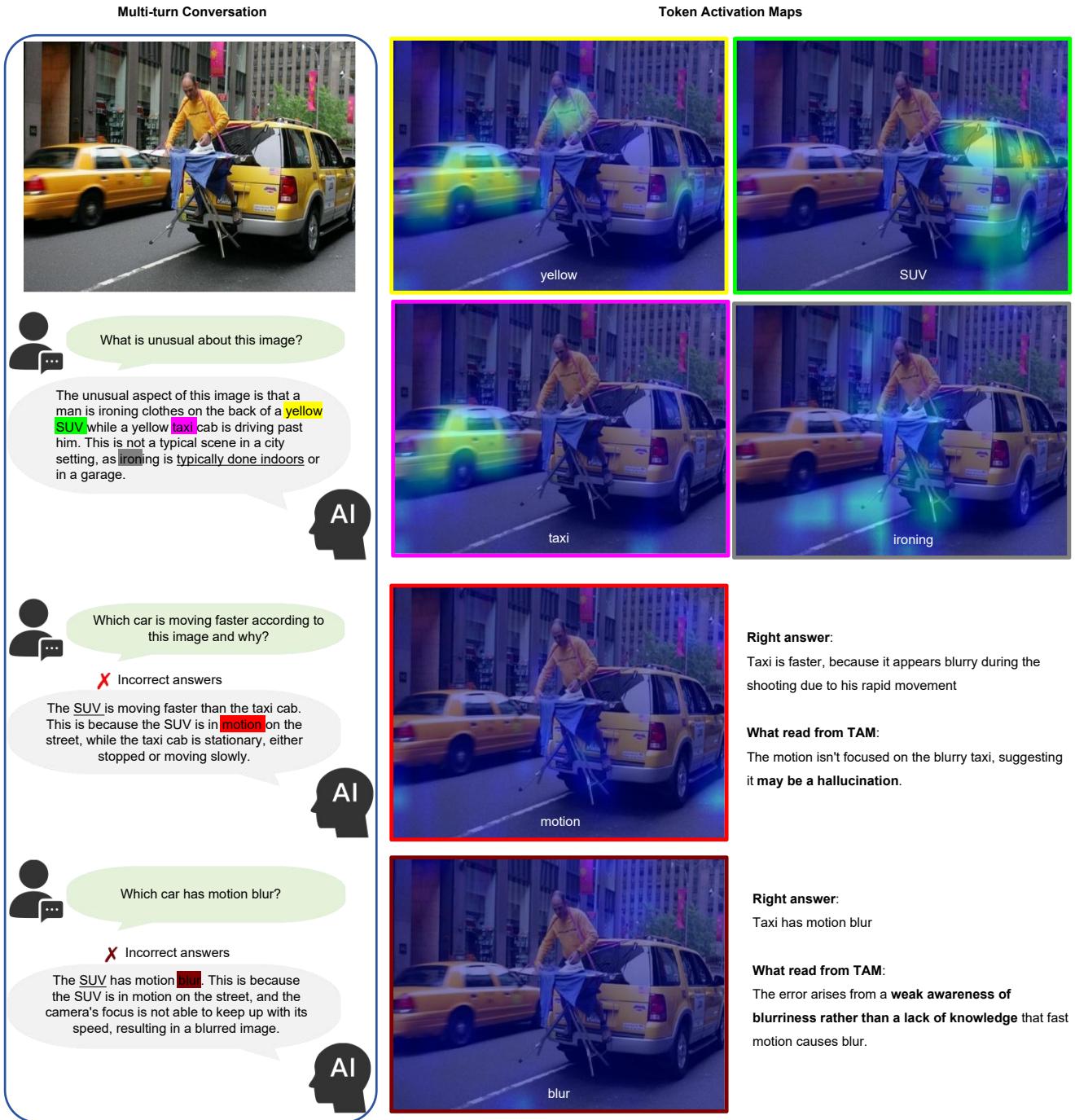


Figure 31. **TAM enables failure case analysis in multi-turn conversation.** Although Qwen2-VL-2B [51] well recognizes objects with good explanation results in the first round chat, it fails to identify motion blur related to speed and mistakenly regards the SUV as the faster car. The clues provided by TAM reveal that the failure to recognize motion blur is the primary reason for this error, highlighting TAM’s effectiveness in supporting detailed error analysis from multi-turn conversation.