# Don't Deceive Me: Mitigating Gaslighting through Attention Reallocation in LMMs

Pengkun Jiao, Bin Zhu, Jingjing Chen, Chong-Wah Ngo, and Yu-Gang Jiang

## Abstract

Large Multimodal Models (LMMs) have demonstrated remarkable capabilities across a wide range of tasks. However, their vulnerability to user gaslighting—the deliberate use of misleading or contradictory inputs—raises critical concerns about their reliability in real-world applications. In this paper, we address the novel and challenging issue of mitigating the negative impact of negation-based gaslighting on LMMs, where deceptive user statements lead to significant drops in model accuracy. Specifically, we introduce GasEraser, a training-free approach that reallocates attention weights from misleading textual tokens to semantically salient visual regions. By suppressing the influence of "attention sink" tokens and enhancing focus on visually grounded cues, GasEraser significantly improves LMM robustness without requiring retraining or additional supervision. Extensive experimental results demonstrate that GasEraser is effective across several leading open-source LMMs on the GaslightingBench. Notably, for LLaVA-v1.5-7B, GasEraser reduces the misguidance rate by 48.2%, demonstrating its potential for more trustworthy LMMs.

## Keywords

## 1 Introduction

Large Multimodal Models (LMMs) [8, 9, 16, 20, 24, 31] combine the language understanding of Large Language Models (LLMs) [10, 25] with powerful visual encoders, such as CLIP [21] and DINO [4], enabling reasoning over visual and textual inputs. The success of LMMs can be primarily attributed to their attention mechanism, which dynamically establishes associations among the input sequence tokens—representing the fundamental units of both visual and textual information [27]. By selectively focusing on the most salient parts of the input sequence, this mechanism enhances the model's ability to generate contextually appropriate and coherent responses.

Despite their impressive capabilities, LMMs implicitly assume that user inputs are "honest"—that is, neutral and factually accurate. However, in real-world applications, users may provide misleading or adversarial inputs—whether intentionally or not—that can distort the model's reasoning process. A particularly subtle and impactful form of such manipulation is negation-based gaslighting: when a deceptive follow-up statement causes the model to contradict its original, correct answer [2, 18, 36]. As illustrated in Figure 1, such gaslighting exploits weaknesses in the model's attention distribution, often leading to responses that align with the false user-provided information rather than the visual evidence. From an attention perspective, gaslighting is facilitated by attention sink [17, 32] that tokens that absorb disproportionately high attention scores despite contributing little or no relevant semantic or visual content.



**Figure 1: Illustration of negation-based gaslighting in Large Multimodal Models (LMMs). A negation-based gaslighting statement refers to a misleading user prompt that contradicts the initial correct answers (e.g., "There are two pineapples in the image," when only one is present). The figure demonstrates how such deceptive inputs can override the model's initially accurate response, leading it to adopt the false premise.**

These tokens, once emphasized, can diminish the model's focus on meaningful image regions, leading to erroneous or incoherent responses (see Figure 3 (b)). Recent studies [6, 13, 15, 28, 32, 35] have identified attention sink patterns in both LLMs and LMMs, showing their detrimental effect on generation quality.

In this paper, we hypothesize that gaslighting is particularly potent when attention sinks are present and unmitigated. To address this, we propose **GasEraser**, a training-free method that reallocates attention from misleading or irrelevant textual inputs to image-centric attention. GasEraser identifies attention regions distorted by gaslighting and strategically redistributes their influence to enhance visual grounding. This redirection is guided by a head selection mechanism that identifies vision-relevant attention heads and suppresses those associated with sink-like behavior. GasEraser is both plug-and-play and training-free, which can be seamlessly integrated into existing LMMs. Through comprehensive experimental evaluations, we demonstrate that GasEraser substantially enhances the reliability and robustness of LMMs on the GaslightingBench [36], as shown in Figure 2.
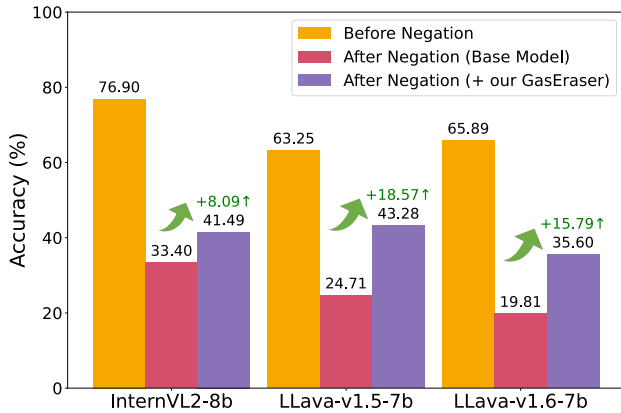
**Figure 2: Performance comparison of three models on GaslightingBench, highlighting the impact of negation-based gaslighting and the effectiveness of the proposed GasEraser. The figure shows the models' accuracy under three conditions: before negation, after negation for base LMMs, and after negation with GasEraser applied to the base LMMs.**

The main contributions of this paper are as follows:

- We propose a pioneering study addressing gaslighting in LMMs from a novel perspective, analyzing the phenomenon through the lens of visual attention sinks.
- We introduce **GasEraser**, a training-free strategy that dynamically reallocates attention weights to mitigate the impact of negation-based gaslighting inputs and enhance model reliability.
- We provide comprehensive experimental results validating the effectiveness of our approach, demonstrating its ability to improve the robustness and accuracy of LMM outputs in the presence of gaslighting or adversarial inputs.

## 2 Related Works

### 2.1 Large Multi-Modal Models (LMMs)

Large Multi-Modal Models (LMMs) enhance the capabilities of Large Language Models (LLMs) (e.g., LLaMA [25], Vicuna [10]) by incorporating visual inputs through the integration of pretrained vision encoders (e.g., CLIP [21], DINO [4]), typically connected via a vision-language projection or cross-modal attention mechanism. This architecture enables unified representation and reasoning across both image and text modalities. Early models such as Florence, BLIP/BLIP-2 [19], PaLI [7], and Flamingo [1] demonstrated effective pretraining strategies and cross-modal alignment. More recent models, including LLaVA [20], Qwen-VL [31], InternVL [9], Gemini [24], and GPT-4o [16], represent the rapid progress toward scalable, aligned, and interactive general-purpose multimodal systems. In this paper, we employ representative LMMs to address the challenges of gaslighting tasks through attention reallocation.

### 2.2 Negation in LLMs and LMMs

Negation, in linguistic terms, refers to the contradiction or denial of a proposition [11]. Recent studies have significantly advanced our understanding of negation, particularly through research such

as [26], which demonstrate that LLMs, including GPT-3 and Instruct-GPT, face considerable challenges in processing negation. These models often struggle to accurately interpret the lexical semantics of negation, fail to maintain logical consistency, and encounter difficulties in reasoning effectively when confronted with negated contexts. Moreover, LLMs frequently show an inability to defend correct beliefs against invalid arguments, raising concerns about their alignment and depth of understanding [29]. In the domain of vision-language models (VLMs), research has explored strategies to improve their understanding of opposing arguments, particularly in CLIP-like models [2, 22, 30, 34]. These studies have revealed significant limitations in VLMs' ability to handle negation, particularly across tasks such as retrieval and multiple-choice questions involving negated statements.

A very recent study, GaslightingBench [36], further extends this inquiry by examining the impact of opposing arguments on language models. This research investigates how LLMs respond to the challenge of maintaining consistent reasoning and logical integrity when confronted with misleading or unfaithful negation arguments. In this paper, we aim to mitigate the effects of gaslighting in LLMs within multi-round conversational settings. Specifically, it explores instances where the models' initial responses are correct, but they are subsequently misled by unfaithful negation arguments during the course of a conversation.

### 2.3 Attention Sink

The Attention Sink phenomenon [32] has been observed in large language models (LLMs), wherein a small subset of tokens—typically the initial few—receive disproportionately high attention scores despite conveying limited semantic information. Xiao *et al.* [32] demonstrated that LLMs allocate substantial attention to these early tokens regardless of their informational value, giving rise to the Attention Sink effect. Several foundational studies have sought to investigate the underlying causes of this phenomenon. Cancedda *et al.* [3] identified that Attention Sink primarily occurs in the first token, attributing this bias to the large norm of its hidden state. In contrast, Sun *et al.* [23] and Yu *et al.* [33] observed that Attention Sink may also manifest in various word tokens with limited semantic relevance, without being confined to a fixed position in the input sequence. This broader manifestation complicates attention distribution and underscores the necessity for more refined attention mechanisms in LLMs. The implications of Attention Sink are far-reaching, with relevance to multiple downstream tasks and model optimizations, including long-context generation [14, 32], key-value (KV) cache optimization [13, 28], efficient inference [5, 35], and model quantization [6].

### 2.4 Visual Attention Sink

The same phenomenon, where tokens with limited information receive disproportionately high attention scores, is also observed in large multimodal models (LMMs). Timothée *et al.* [12] demonstrate that high-norm tokens often appear during inference, primarily in low-informative background areas of images. Seil *et al.* [17] further emphasize that these low-informative background regions can exhibit high norm values, which they refer to as the "visual attention sink." More details can be found in Section 3.3.
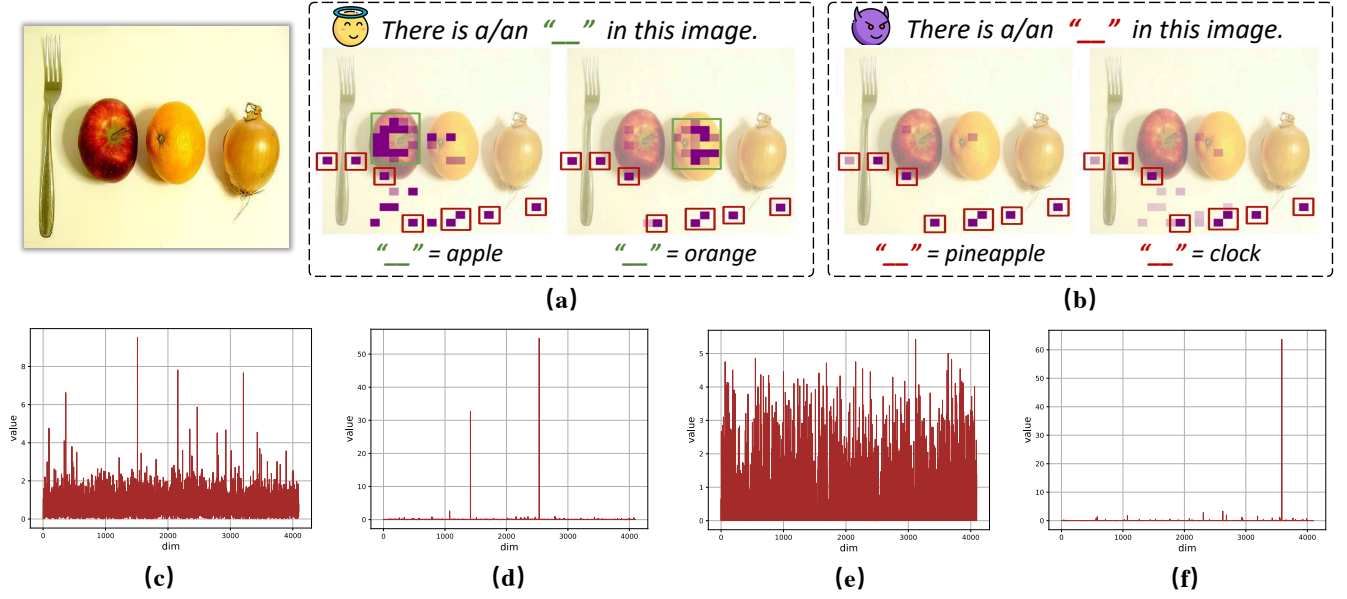
**Figure 3: (a)** The image-relevant token attends to both key and some irrelevant visual features. **(b)** Gaslighting tokens primarily focus on irrelevant visual features. **(c)** and **(e)** show normal token embeddings for LLaVA-v1.5-7B and InternVL2-8B, while **(d)** and **(f)** show the corresponding sink token embeddings, which exhibit significantly higher norms in specific dimensions.

In this paper, we leverage the characteristics of attention sinks to reallocate gaslighting attention to image-centric attention, thereby enhancing the representation of salient visual features.

## 3 Preliminary

### 3.1 Gaslighting in LMMs Task

In the Gaslighting Task, each instance consists of a reference image, a neutral multiple-choice question, and a gaslighting statement. Formally, the input is represented as a triplet $(I, T_q, T_g)$, where $I$ denotes the image, $T_q$ is a neutral question (e.g., "How many people are in this image?" with options "A. One, B. Two, C. Three, D. Zero"), and $T_g$ is a misleading statement that contradicts the visual content (e.g., "There are two people in this image"), with the ground truth answer being "A. One." Let $F$ denote the multimodal model, comprising a vision encoder $V$ and a language model $G$. The image $I$ is encoded into visual tokens $t_v = V(I)$, while the textual inputs $(T_q, T_g)$ are tokenized into textual tokens $t_t$ via a tokenizer. Let $d_t$ denote the embedding dimension. The concatenated sequence $[t_v, t_t]$ is then processed by $G$ for joint reasoning.

While large multimodal models (LMMs) are typically capable of answering the neutral question $T_q$ correctly when considered in isolation, they often produce incorrect responses when conditioned on the misleading statement $T_g$. Our goal is to improve the robustness of these models by reducing their susceptibility to adversarial or manipulative inputs, thereby enabling them to consistently produce correct responses.

### 3.2 Self-Attention

Self-attention is a fundamental mechanism for modeling contextual dependencies within sequences. It operates over the combined input sequence $x = [t_v, t_t]$, where $t_v$ and $t_t$ represent the visual and textual tokens, respectively. The scaled dot-product attention is defined as:

$$\text{Attn}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right) V, \quad (1)$$

where $Q$, $K$, and $V$ are the query, key, and value matrices projected from the input token embeddings, and $d$ is the embedding dimension.

**Multi-Head Self-Attention** extends the basic self-attention mechanism by projecting the input into multiple subspaces, allowing the model to capture diverse patterns and dependencies. It is formally defined as:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \ldots, \text{head}_h)W^O, \quad (2)$$

$$\text{where} \quad \text{head}_i = \text{Attn}(QW_i^Q, KW_i^K, VW_i^V). \quad (3)$$

Here, $W_i^Q$, $W_i^K$, and $W_i^V$ are learned projection matrices specific to the $i$-th attention head, and $W^O$ is the output projection matrix. In the context of LMMs, this architecture enables the model to simultaneously attend to and integrate information from multiple representational subspaces, facilitating both visual perception and language understanding.

### 3.3 Visual Attention Sink

In practice, certain tokens may receive disproportionately high attention despite lacking visual or semantic relevance, a phenomenon referred to as sink tokens [17, 32]. These sink tokens contribute minimally to inference and often distort the attention distribution, as illustrated in Figure 3(b). These tokens exhibit abnormally large values in specific embedding dimensions, as shown in Figures 3(d)
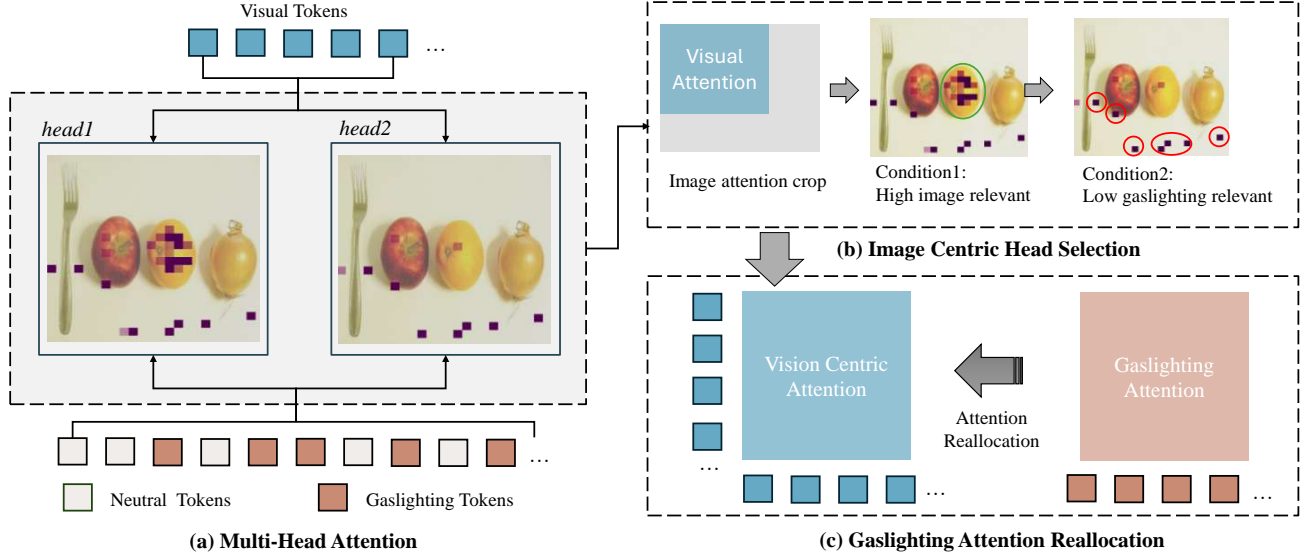
**Figure 4: Illustration of our GasEraser. (a) Multi-head attention applies multiple attention mechanisms in parallel, allowing the model to capture different perspectives of the information. (b) We evaluate the relevance between image and text tokens to identify which visual-textual associations are important. (c) We then relocate attention from less important associations that have high attention scores to those that are more relevant.**

and (f). For example, the sink token shows high norms in dimensions 1,415 and 2,533 for LLaVA-v1.5-7B, and in dimensions 2,624 and 3,584 for InternVL2-8B, consistent with the findings in [17].

To identify sink tokens, we adopt the method proposed in [17, 32], which selects sink tokens based on their high norm in specific dimensions $\mathcal{D}$. Specifically, sink tokens can be defined as those for which the following condition holds:

$$\mathcal{I}_{\text{sink}} = \{i \mid \max\left(\|\mathbf{e}_i\|_{d\in\mathcal{D}}\right) > \tau\}, \quad \mathbf{e}_i = \frac{1}{\sqrt{d}}\sqrt{\sum_{j=1}^{d}|x_{ij}|^2} \quad (4)$$

where $\mathbf{e}_i$ represents the normalized embedding vector for token $i$, and tokens whose normalized embeddings exceed the threshold $\tau$ are classified as sink tokens.

Specifically, let $\mathcal{V}_{\text{sink}}$ and $\mathcal{T}_{\text{sink}}$ denote the sets of indices corresponding to sink visual and sink textual tokens, respectively, defined as follows:

$$\mathcal{V}_{\text{sink}} = \mathcal{I}_{\text{sink}}[\mathcal{I}_{\text{start}} : \mathcal{I}_{\text{end}}], \quad (5)$$

$$\mathcal{T}_{\text{sink}} = \mathcal{I}_{\text{sink}} \setminus \mathcal{V}_{\text{sink}}, \quad (6)$$

where $\mathcal{I}_{\text{start}}$ and $\mathcal{I}_{\text{end}}$ denote the start and end positions of the image (visual) tokens within $\mathcal{I}_{\text{sink}}$, respectively.

## 4 GasEraser: Gaslighting Attention Reallocation

Consider the attention mechanism within an input sequence: when LLMs are exposed to gaslighting, certain tokens may receive disproportionately high attention scores. This distortion impairs the model's ability to establish accurate visual-textual attention relationships. By identifying these misleading high attention scores and reallocating them to their relevant counterparts, the model can generate more accurate predictions that better align with the

content of the image. To address this, we introduce **GasEraser**, a Visual-Text Attention Reallocation method designed to enhance the robustness of Large Multi-Modal Models (LMMs) against negation inputs. The core idea is to mitigate the influence of attention sinks on non-essential visual-textual tokens, while amplifying attention on relevant tokens that reinforce the meaningful image-text correspondence. We will detail GasEraser in the following sections.

### 4.1 Gaslighting Attention in the Visual Attention Sink View

Gaslighting tokens manipulate attention distributions by assigning disproportionately high attention to irrelevant visual features, as shown in Figure 3 (b), where image-irrelevant tokens occupy unusually high attention scores in the background.

To formalize this, let the multi-head attention maps in a single transformer layer be denoted as $\mathbf{A} \in \mathbb{R}^{H \times S \times S}$, where $H$ represents the number of attention heads and $S$ is the sequence length. The attention distribution $A_{h,v,t}$ in head $h$ between the $t$-th text token and the $v$-th visual token can be altered by potential gaslighting text tokens, leading to erroneous attention scores and undermining the model's focus on visual information.

Consequently, our goal is to identify the gaslit attention patterns between gaslighting text tokens and visual tokens. Afterward, we reallocate this attention to vision-centric heads that are more closely aligned with the relevant visual features, thereby enhancing the model's ability to focus on meaningful visual information.

### 4.2 Vision-Centric Head Selection

In transformer architectures, attention heads within each layer capture various aspects of the input. Our objective is to identify

*image-centric* attention heads that are crucial for visual perception. To achieve this, we propose a **Vision-Centric Head Selection** strategy, inspired by [17], to isolate the heads most relevant to visual grounding. The goal is to select the heads that exhibit the strongest alignment with visual features.

We begin by extracting the attention weights corresponding to visual tokens and computing the image relevance score $\delta_{h,s} \in \mathbb{R}$ for each head $h \in \{1, \ldots, H\}$ and source position $s \in \{1, \ldots, S\}$:

$$\delta_{h,s} = \sum_{i=\mathcal{I}_{\text{start}}}^{\mathcal{I}_{\text{end}}} \mathbf{A}_{h,s,i}. \tag{7}$$

Next, we compute the *sink-likelihood score* $\xi_{h,s} \in \mathbb{R}$, which quantifies the normalized attention paid to the sink token:

$$\xi_{h,s} = \frac{\sum_{j \in \mathcal{V}_{\text{sink}}} \mathbf{A}_{h,s,j}}{\delta_{h,s} + \varepsilon}, \tag{8}$$

where $\varepsilon \in \mathbb{R}^+$ is a small constant added to prevent division by zero.

Using these computed scores, we define the set of *visual-centric heads* as:

$$\mathcal{H}_{\text{visual}} = \left\{ (h,s) \mid \delta_{h,s} \leq \rho \ \wedge \ \xi_{h,s} \geq \alpha \right\}, \tag{9}$$

where $\rho$ and $\alpha$ are predefined thresholds. The condition $\delta_{h,s} \leq \rho$ ensures that the head does not focus exclusively on non-image information, while $\xi_{h,s} \geq \alpha$ ensures that the head is not likely to exhibit gaslighting attention.

The attention heads identified by $\mathcal{H}_{\text{visual}}$ are then selected for enhanced attention, thereby enhancing the model's ability to ground vision-language representations.

## 4.3 Gaslighting Attention Reallocation

Once the high-relevance vision-centric attention and its gaslighting counterparts are identified, we proceed to reallocate the attention scores accordingly. For each attention map in a transformer layer, we extract the relevant slice of the visual-centric attention map:

$$\hat{A} = A[\mathcal{H}_{\text{visual}}, :], \tag{10}$$

The attention values at the text sink token positions are scaled by a factor $p$, where $0 < p < 1$:

$$\hat{A}[:, \mathcal{T}_{\text{sink}}] \leftarrow \hat{A}[:, \mathcal{T}_{\text{sink}}] \cdot p, \tag{11}$$

Next, we compute the gaslighting attention budget $\Omega$ that is removed from the text tokens:

$$\Omega = \sum_{i \in \mathcal{T}_{\text{sink}}} \hat{A}[:, i] \cdot (1 - p), \tag{12}$$

We then set all attention weights for the visual tokens to zero $(\hat{A}_h[:, \mathcal{V}_{\text{sink}}] = 0)$ to eliminate the influence of sink image tokens. Next, we calculate the vision-centric attention ratio for each head:

$$R_{\mathcal{V}} = \frac{\hat{A}_h[:, \mathcal{I}_{\text{start}} : \mathcal{I}_{\text{end}}]}{\sum_{i=\mathcal{I}_{\text{start}}}^{\mathcal{I}_{\text{end}}} \hat{A}[:, i]}, \tag{13}$$

Finally, we reallocate the attention weights, and the updated attention map is written back to the original map:

$$A[\mathcal{H}_{\text{visual}}, \mathcal{I}_{\text{start}} : \mathcal{I}_{\text{end}}] \leftarrow \hat{A}[:, \mathcal{I}_{\text{start}} : \mathcal{I}_{\text{end}}] + \Omega \cdot R_{\mathcal{V}}. \tag{14}$$

By following these steps, the model ensures that relevant visual-textual attention is focused on the correct parts of the image, while suppressing misleading tokens that have received exaggerated attention due to gaslighting or other factors. This approach improves the robustness and interpretability of multimodal models.

Note that Gaslighting attention is primarily induced by misleading textual statements. Therefore, we focus on reallocating attention from sink text tokens. Nonetheless, experiments and further discussion on reallocating attention using both image and text tokens are provided in Section 5.3.1.

## 4.4 Integration with Inference

Our method can be seamlessly integrated as a plug-in within the attention layer, without the need for retraining the LMM. Building on recent findings suggesting that visual perception is more prominent in the earlier layers of transformer blocks in LMMs [5], we evaluate the performance of our method across different injected layers, as discussed in Section 5.3.2. The results show that the performance gain is most significant in the front layers. Therefore, in this paper, we integrate the top 16 layers of LLMs.

## 5 Experiments

## 5.1 Experimental Setup

*5.1.1 Benchmark.* We utilize **GaslightingBench** [36], the only existing multimodal gaslighting benchmark, for our evaluation. This benchmark consists of 20 categories and 1,287 samples. Each sample is presented in a multiple-choice format and includes an image, a paired question, several answer options, and a deliberately misleading statement. In the first round of interaction, the model is prompted to respond based solely on the original question. In the second round, a misleading statement is introduced to evaluate the model's robustness to gaslighting attempts.

Notably, our approach differs from the prompt strategy used in GaslightingBench, which elicits more open-ended responses in the second round. In contrast, we design our prompts to require the model to strictly choose from predefined options that align with the instruction. Interestingly, we observe that LMMs demonstrate greater robustness to gaslighting when allowed to respond with open-form answers—such as explicitly stating "The statement is wrong." We discuss this observation in greater detail in Section 5.4.

*5.1.2 Prompt Design.* The prompt for the first-round conversation follows the format: `[system prompt]` `user:` `[image]` `[question]` `[options]`. In the second round, a misleading statement is added to test the model's robustness. The format becomes: `[system prompt]` `user:` `[image]` `[question]` `[options]` `assistant:` `[answer1]` `user:` `[gaslighting statement]`. A complete example of the conversation structure is shown in Figure 5. For a more detailed discussion of the prompt design, please refer to the Supplemental materials A.

*5.1.3 Base Models and Configuration.* We evaluate our proposed approach using three representative open-source large multimodal models (LMMs):

(1) **LLaVA-1.5-7B** [20], which incorporates the vision encoder CLIP-L-336px and the LMM LLaMA-2-7B-Chat;

(2) **LLaVA-1.6-Vicuna-7B** [20], which integrates the vision encoder CLIP-L-336px and the LMM Vicuna-7B;

| Method | vision encoder | LLM | before negation | after negation | gain |
|---|---|---|---|---|---|
| LLaVA-v1.5-7B [20] | CLIP-L-patch14-336px | LLaMA-2-7B-Chat | 63.25 | 24.71 | - |
| + GasEraser (ours) | | | 63.25 | **43.28** | **+18.57** |
| LLaVA-v1.6-7B [20] | CLIP-L-patch14-336px | vicuna-7b-v1.5 | 65.89 | 19.81 | - |
| + GasEraser (ours) | | | 65.89 | **35.60** | **+15.79** |
| InternVL2-8B [9] | InternViT-300M-448px | InternLM2-5-7B-Chat | 76.90 | 33.40 | - |
| + GasEraser (ours) | | | 76.90 | **41.49** | **+8.09** |

Table 1: Performance comparison on GaslightingBench after incorporating our proposed GasEraser into three representative MLLMs. "Before negation" refers to the accuracy of the model's initial answers, while "after negation" denotes the accuracy following the introduction of the gaslighting statement.

| image tokens | text tokens | LLaVA-v1.5-7B [20] | | LLaVA-v1.6-7B [20] | | InternVL2-8b [9] | |
|---|---|---|---|---|---|---|---|
| | | before negation | after negation | before negation | after negation | before negation | after negation |
| × | × | 63.25 | 24.71 | 65.89 | 19.81 | 76.90 | 33.40 |
| ✓ | × | 63.25 | 25.87 | 65.89 | 14.60 | 76.90 | 34.81 |
| × | ✓ | 63.25 | 43.28 | 65.89 | **35.60** | 76.90 | 40.33 |
| ✓ | ✓ | 63.25 | **43.50** | 65.89 | 32.78 | 76.90 | **41.19** |

Table 2: Performance comparison of attention relocation token sources (image and text) in our method, before and after negation, on GaslightingBench. Here, "image tokens" and "text tokens" indicate whether the image sink token and text sink token are used, respectively.

(3) **InternVL2-8B** [9], which uses the vision encoder InternViT-300M-448px and the LLM InternLM2-5-7B-Chat.

Our approach is training-free, with all model parameters frozen. Experiments were conducted on A6000 GPUs.

*5.1.4 Hyperparameter Selection.* For **LLaVA-v1.5-7B**, we set the hyperparameters as follows: $\tau = 20$, $\rho = 0.6$, $\alpha = 0.005$, and $p = 0.6$. For **LLaVA-v1.6-Vicuna-7B**, we set the hyperparameters as follows: $\tau = 20$, $\rho = 0.6$, $\alpha = 0.01$, and $p = 0.6$. For **InternVL2-8B**, we set the hyperparameters as follows: $\tau = 20$, $\rho = 0.6$, $\alpha = 0.1$, and $p = 0.6$. Additional results with varying hyperparameters can be found in the Supplementary Material.

## 5.2 Result Analysis

*5.2.1 Performance Comparison.* Table 1 presents a performance comparison between baseline large multimodal models (LMMs) and their counterparts enhanced with the proposed method (+ GasEraser ) on the GaslightingBench dataset. Relative to their original performance prior to negation, all baseline models exhibit substantial accuracy degradation following exposure to negated or misleading statements. For instance, the accuracy of LLaVA-v1.5-7B drops from 63.25% to 24.71%, LLaVA-v1.6-7B declines from 65.89% to 19.81%, and InternVL2-8B decreases from 76.90% to 33.40%. These pronounced declines highlight a consistent vulnerability among existing LMMs in handling negated semantics, underscoring their susceptibility to linguistic gaslighting.

By integrating the proposed GasEraser, all evaluated models demonstrate improved resilience to negation. Specifically, LLaVA-v1.5-7B achieves a post-negation accuracy of 43.28%, representing an absolute improvement of 18.57%. Similarly, LLaVA-v1.6-7B improves by 15.79%, and InternVL2-8B gains 8.09%. These results affirm the effectiveness of GasEraser in mitigating the semantic

disruption caused by negation, enhancing both interpretability and robustness. Overall, the proposed method serves as a valuable enhancement for strengthening the reliability of LMMs in adversarial language scenarios.

*5.2.2 Qualitative Results.* Figure 5 presents examples illustrating how LMMs respond to negation arguments across different tasks, *i.e.,* image emotion recognition, image topic classification, and object localization. In each case, the models initially produce correct responses. However, when negation arguments are introduced, many models revise their answers incorrectly. With our proposed approach, the model consistently maintains the correct response, demonstrating improved robustness against misleading inputs.

## 5.3 Analysis of GasEraser Design Choices

We further analyze the proposed approach by leveraging image sink tokens and text sink-token sources to reallocate attention scores, aiming to identify the primary contributors to gaslighting-induced attention shifts. Additionally, we examine the effects of integrating GASERASER at different layers of the model to determine which stages benefit most from enhanced visual grounding. An ablation study is also conducted to assess the contribution of key components within the GASERASER framework. These analyses are presented in detail in the following sections.

*5.3.1 Image or Text Tokens Matter More in Mitigating Gaslighting?* We investigate the impact of relocated attention sources, specifically image and text tokens. The results from using various sources are presented in Table 2. Our findings clearly demonstrate the superiority of using text tokens as a source, compared to image tokens. When only image tokens are used, the accuracy improves by 1.16 points, whereas using text tokens alone results in a significant

**Figure 5: Qualitative examples using LLaVA-1.5-7B as the base model. The base model generates incorrect answers when misled by gaslighting negation statements, whereas our method effectively mitigates the impact of such misleading content. The ground truth option is highlighted in green.**

gain of 18.57 points. This disparity can be attributed to the fact that performance degradation caused by negation primarily arises from gaslighting tokens, which predominantly originate from text tokens. In contrast, visual tokens are more neutral and contribute less to the model's susceptibility to gaslighting. Consequently, when the visual budget is removed, its impact on overall performance is minimal.

*5.3.2  Which Layers Should GasEraser Be Integrated at?* To evaluate the impact of layer selection on model performance, we apply our method to different layers of the LLaVA-v1.5-7B model and present the results in Figure 6 (detailed results are provided in the supplementary material C). As illustrated, the majority of performance improvements are concentrated in the earlier layers, with the top 16 layers achieving the highest accuracy. This trend suggests that the lower layers of the transformer architecture are particularly sensitive to visual inputs. At these stages, the model primarily focuses on processing visual features, which are essential for establishing precise visual-textual alignments.

This observation is consistent with recent findings suggesting that visual tokens exert the greatest influence in the early layers of multimodal models [5]. These insights have important implications for designing more computationally efficient inference strategies. By prioritizing optimization in the early layers, it is possible to reduce processing overhead while preserving the efficacy of visual grounding mechanisms.



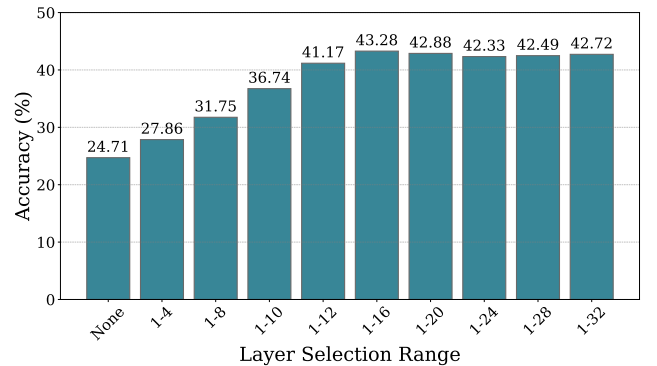**Figure 6: Performance across different layer selections on LLaVA-v1.5-7B.**

| Image-centric head selection | Gaslighting token selection | LLaVA-v1.5-7B | InternVL2-8B |
|:---:|:---:|:---:|:---:|
| × | × | 24.71 | 33.40 |
| × | ✓ | 24.86 | 35.83 |
| ✓ | × | 37.45 | 36.68 |
| ✓ | ✓ | **43.28** | **41.49** |

**Table 3: Ablation study on image-centric head selection and gaslighting token selection.**
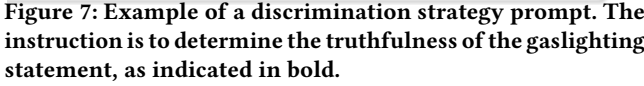
**Figure 7: Example of a discrimination strategy prompt. The instruction is to determine the truthfulness of the gaslighting statement, as indicated in bold.**

| Method | before negation | after negation |
|---|---|---|
| LLava-v1.5-7b [20] | 63.25 | 42.91 |
| **+ our GasEraser** | 63.25 | **45.38** |
| LLava-v1.6-7b [20] | 65.89 | 50.81 |
| **+ our GasEraser** | 65.89 | **51.20** |
| InternVL2-8b [9] | 76.90 | 70.57 |
| **+ our GasEraser** | 76.90 | **71.25** |

**Table 4: Accuracy comparison with and without our method, before and after negation, under the discrimination strategy.**

*5.3.3 Ablation Study.* We perform an ablation study to evaluate the effectiveness of two components: image-centric head selection and gaslighting token selection, using LLaVA-v1.5-7B and InternVL2-8B. The results are summarized in Table 3. Excluding the image-centric head and distributing attention across all image-related heads leads to a drop in performance. This may be due to the equal amplification of both relevant and irrelevant visual information, resulting in insufficient focus on critical visual cues. Similarly, when gaslighting token selection is removed and attention is distributed across all textual tokens, performance also declines. This is likely due to the dilution of attention, which, while enhancing alignment with visual cues, reduces the model's sensitivity to the gaslighting tokens crucial for detecting and mitigating gaslighting. The best performance is achieved when both components are included, highlighting their complementary roles in enhancing model effectiveness.

## 5.4 Discrimination Strategy

In addition to the standard option selection setup, we explore a discrimination strategy, where the model is explicitly prompted to be aware of potentially misleading statements. The prompt design is illustrated in Figure 7, and the results, both with and without our proposed method, are presented in Table 4. In this setting, we observe that models tend to perform more reliably when tasked with binary judgment (i.e., classifying statements as either "Correct" or "Wrong"), rather than choosing from multiple options. This improvement may stem from the fact that validating a statement (i.e., determining whether it is true or false) is a simpler task than answering a question in a more complex QA scenario. The answer space in this mode is limited to two options—*Correct* or *Wrong*—which likely reduces ambiguity and cognitive load. However, it is important to recognize that such a discrimination format may not fully reflect practical or natural user interactions, where end users typically seek direct answers rather than explicit validation of intermediate statements. Nonetheless, our method demonstrates consistent improvements under the discrimination setting, further supporting its robustness across varying prompt designs and interaction paradigms. While promising, future work should investigate how insights from this simplified evaluation setting can be incorporated into more realistic dialog systems.

## 6 Conclusion

Gaslighting in Large Multimodal Models (LMMs) remains a significant yet underexplored challenge. In this paper, we propose GasEraser, a novel, training-free approach designed to mitigate the adverse effects of negation in MLLMs. Extensive empirical evaluations demonstrate the effectiveness of our method. Our analysis reveals that gaslighting-related attention predominantly originates from text tokens, aligning with the intuition that negation is primarily introduced through textual input. Moreover, we find that reallocating attention in the early layers—responsible for low-level visual processing—significantly enhances model robustness. This finding underscores the critical role of early-stage visual perception in resisting deceptive or contradictory inputs. Overall, our findings highlight the potential of GasEraser to improve the reliability of MLLMs in gaslighting scenarios, where misleading or adversarial prompts are prevalent. By providing an attention-based perspective, this work contributes to the advancement of more robust and trustworthy multimodal reasoning systems.

**Future Work.** We observe that LMMs demonstrate greater robustness in distinguishing gaslighting statements under the Discrimination Strategy. It would be interesting to explore how dialectical reasoning can be integrated into general dialogue tasks. In future work, we plan to further investigate methods for incorporating dialectical thinking into the inference of LMMs.

## References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems* 35 (2022), 23716–23736.

[2] Kumail Alhamoud, Shaden Alshammari, Yonglong Tian, Guohao Li, Philip Torr, Yoon Kim, and Marzyeh Ghassemi. 2025. Vision-language models do not understand negation. *arXiv preprint arXiv:2501.09425* (2025).

[3] Nicola Cancedda. 2024. Spectral Filters, Dark Signals, and Attention Sinks. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 4792–4808. doi:10.18653/v1/2024.acl-long.263

[4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*. 9650–9660.

[5] Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. 2024. An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. In *European Conference on Computer Vision*. Springer, 19–35.

[6] Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang, Peng Gao, Kaipeng Zhang, and Ping Luo. 2024. Efficientqat: Efficient quantization-aware training for large language models. *arXiv preprint arXiv:2407.11062* (2024).

[7] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish V Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme Ruiz, Andreas Peter Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. 2023. PaLI: A Jointly-Scaled Multilingual Language-Image Model. In *The Eleventh International Conference on Learning Representations*. https:

//openreview.net/forum?id=mWVoBz4W0u

[8] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xinghai Yu, and Chong Ruan. 2025. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811* (2025).

[9] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024. 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*. 24185–24198.

[10] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)* 2, 3 (2023), 6.

[11] William Croft. 1991. The evolution of negation. *Journal of linguistics* 27, 1 (1991), 1–27.

[12] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. 2024. Vision Transformers Need Registers. In *The Twelfth International Conference on Learning Representations*. https://openreview.net/forum?id=2dnO3LLiJ1

[13] Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. 2024. Model Tells You What to Discard: Adaptive KV Cache Compression for LLMs. In *The Twelfth International Conference on Learning Representations*. https://openreview.net/forum?id=uNrFpDPMyo

[14] Chi Han, Qifan Wang, Hao Peng, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. 2024. LM-Infinite: Zero-Shot Extreme Length Generalization for Large Language Models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. 3991–4008.

[15] Wei Huang, Haotong Qin, Yangdong Liu, Yawei Li, Xianglong Liu, Luca Benini, Michele Magno, and Xiaojuan Qi. 2024. SliM-LLM: Salience-driven mixed-precision quantization for large language models. *arXiv preprint arXiv:2405.14917* (2024).

[16] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276* (2024).

[17] Seil Kang, Jinyeong Kim, Junhyeok Kim, and Seong Jae Hwang. 2025. See What You Are Told: Visual Attention Sink in Large Multimodal Models. In *The Thirteenth International Conference on Learning Representations*. https://openreview.net/forum?id=7uDI7w5RQA

[18] Nora Kassner and Hinrich Schütze. 2020. Negated and Misprimed Probes for Pretrained Language Models: Birds Can Talk, But Cannot Fly. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 7811–7818.

[19] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*. PMLR, 19730–19742.

[20] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 26296–26306.

[21] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PmLR, 8748–8763.

[22] Jaisidh Singh, Ishaan Shrivastava, Mayank Vatsa, Richa Singh, and Aparna Bharati. 2024. Learn" no" to say" yes" better: Improving vision-language models via negations. *arXiv preprint arXiv:2403.20312* (2024).

[23] Mingjie Sun, Xinlei Chen, J Zico Kolter, and Zhuang Liu. 2024. Massive Activations in Large Language Models. In *ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models*. https://openreview.net/forum?id=1ayU4fMqme

[24] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805* (2023).

[25] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).

[26] Thinh Hung Truong, Timothy Baldwin, Karin Verspoor, and Trevor Cohn. 2023. Language models are not naysayers: an analysis of language models on negation benchmarks. In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)*, Alexis Palmer and Jose Camacho-collados (Eds.). Association for Computational Linguistics, Toronto, Canada, 101–114. doi:10.18653/v1/2023.starsem-1.10

[27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

[28] Zhongwei Wan, Ziang Wu, Che Liu, Jinfa Huang, Zhihong Zhu, Peng Jin, Longyue Wang, and Li Yuan. 2024. LOOK-M: Look-Once Optimization in KV Cache for

Efficient Multimodal Long-Context Inference. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 4065–4078. doi:10.18653/v1/2024.findings-emnlp.235

[29] Boshi Wang, Xiang Yue, and Huan Sun. 2023. Can ChatGPT Defend its Belief in Truth? Evaluating LLM Reasoning via Debate. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 11865–11881. doi:10.18653/v1/2023.findings-emnlp.795

[30] Hualiang Wang, Yi Li, Huifeng Yao, and Xiaomeng Li. 2023. Clipn for zero-shot ood detection: Teaching clip to say no. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1802–1812.

[31] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191* (2024).

[32] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient Streaming Language Models with Attention Sinks. In *The Twelfth International Conference on Learning Representations*. https://openreview.net/forum?id=NG7sS51zVF

[33] Zhongzhi Yu, Zheng Wang, Yonggan Fu, Huihong Shi, Khalid Shaikh, and Yingyan Celine Lin. 2024. Unveiling and Harnessing Hidden Attention Sinks: Enhancing Large Language Models without Training through Attention Calibration. In *Forty-first International Conference on Machine Learning*. https://openreview.net/forum?id=DLTjFFiuUJ

[34] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2023. When and Why Vision-Language Models Behave like Bags-Of-Words, and What to Do About It?. In *The Eleventh International Conference on Learning Representations*. https://openreview.net/forum?id=KRLUvxh8uaX

[35] Zhenyu Zhang, Shiwei Liu, Runjin Chen, Bhavya Kailkhura, Beidi Chen, and Atlas Wang. 2024. Q-hitter: A better token oracle for efficient llm inference via sparse-quantized kv cache. *Proceedings of Machine Learning and Systems* 6 (2024), 381–394.

[36] Bin Zhu, Huiyan Qi, Yinxuan Gui, Jingjing Chen, Chong-Wah Ngo, and Ee-Peng Lim. 2025. Calling a Spade a Heart: Gaslighting Multimodal Large Language Models via Negation. *arXiv preprint arXiv:2501.19017* (2025).