

Understanding Multimodal LLMs: the Mechanistic Interpretability of Llava in Visual Question Answering

Anonymous ACL submission

Abstract

Understanding the mechanisms behind Large Language Models (LLMs) is crucial for designing better models and strategies. While recent studies have yielded valuable insights into the mechanisms of textual LLMs, the mechanisms of Multi-modal Large Language Models (MLLMs) remain underexplored. In this paper, we apply mechanistic interpretability methods to analyze the visual question answering (VQA) mechanisms in an MLLM, Llava. We compare the mechanisms between VQA and textual QA (TQA) in color answering tasks and find that: a) VQA exhibits a mechanism similar to the in-context learning mechanism observed in TQA; b) the visual features exhibit significant interpretability when projecting the visual embeddings into the embedding space; and c) Llava enhances the existing capabilities of the corresponding textual LLM Vicuna during visual instruction tuning. Based on these findings, we develop an interpretability tool to help users and researchers identify important visual locations for final predictions, aiding in the understanding of visual hallucination. Our method demonstrates faster and more effective results compared to existing interpretability approaches. Our code will be available on Github.

1 Introduction

Large Language Models (LLMs) (Brown, 2020; Ouyang et al., 2022; Touvron et al., 2023) have achieved remarkable results in numerous tasks (Xiao et al., 2023; Tan et al., 2023; Deng et al., 2023). However, the underlying mechanisms are not yet well understood. This lack of clarity poses a significant challenge for researchers attempting to address issues such as hallucination (Yao et al., 2023), toxicity (Gehman et al., 2020), and bias (Kotek et al., 2023) in LLMs. Therefore, understanding the mechanisms of LLMs has become an increasingly important area of research. Recently, efforts have been made to explore the mechanisms

behind different LLM capabilities, including factual knowledge (Meng et al., 2022; Geva et al., 2023), in-context learning (Wang et al., 2023; Wei et al., 2023), and arithmetic (Stolfo et al., 2023).

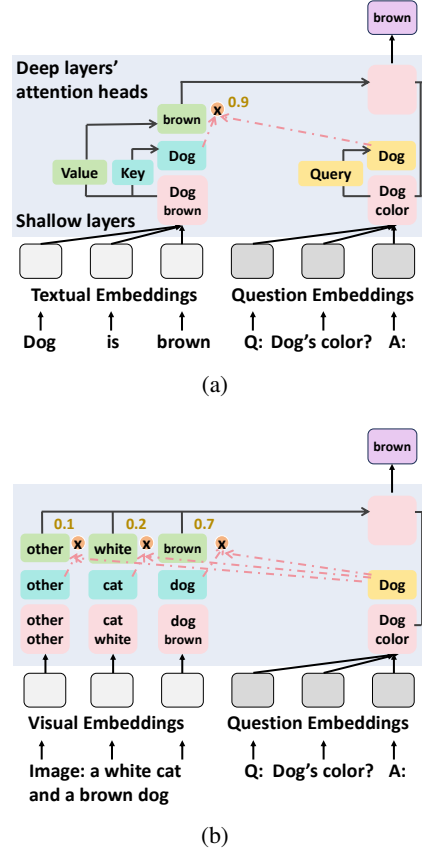


Figure 1: (a) Mechanism of TQA in Vicuna. (b) Mechanism of VQA in Llava.

Although numerous studies have explored the mechanisms of LLMs, they have mainly focused on textual LLMs, often overlooking multi-modal LLMs (MLLMs). It has been demonstrated that features from different modalities, such as images and audio, can significantly enhance the core abilities of LLMs (Zhang et al., 2024). Therefore, investigating the mechanisms of MLLMs is essential. In this paper, we examine the mechanism of VQA in an Multimodal LLM, Llava (Liu et al., 2024b), which is fine-tuned from the existing textual LLM,

Vicuna (Chiang et al., 2023). Our study seeks to address three key questions: a) What is the relationship between the mechanisms of VQA and TQA? b) Are the visual features interpretable under textual LLM’s interpretability analysis method? c) How does Llava acquire its VQA ability during visual instruction tuning?

We investigate the color answering task in VQA, as color is a key feature in images, making it an ideal starting point for VQA analysis. We collect animal photos from the COCO dataset (Lin et al., 2014), each with an animal and its correct color, and pose the question, ‘What is the color of the [animal]?’ For TQA, we generate textual context for each photo, e.g., ‘Dog is brown. Q: What is the color of the dog? A:’, with the correct answer being ‘brown.’ We explore the TQA mechanism in Vicuna using the interpretability method of Yu and Ananiadou (2024a), as shown in Figure 1(a). In shallow layers, the color position (‘brown’) extracts animal features (‘dog’). In deeper layers, attention heads capture color features, and the query-key matrices measure the similarity between the question about the animal and the color position, increasing the probability of ‘brown’ as the answer. When the question aligns with the textual context, the color position receives a high attention score, leading to a large log probability increase for ‘brown’.

Next, we investigate the mechanism of VQA in Llava, starting by using log probability increase scores to identify the most important image regions, which we find to be the image patches related to the animals (as shown in Figure 2). We then apply similar methods used to analyze the value-output matrices and the query-key matrices for these key output vectors. Our analysis reveals that the VQA mechanism is similar to that of TQA: the value-output matrices extract color information, while the query-key matrices compute the similarity between the question content and the animal features. Furthermore, we analyze the visual features by projecting them into the embedding space (Dar et al., 2022), discovering that the visual embeddings exhibit significant interpretability regarding colors and animals, indicating that these embeddings already contain essential information about both. Based on these findings, we conclude the VQA mechanism shown in Figure 1(b). The visual embeddings store information about animals and colors, which is then transferred to deeper layers via the positions’ residual streams. In the deep

layers’ attention heads, the value-output matrices extract color features, while the query-key matrices calculate the similarity between the question and the animal features. Finally, we compare the most important heads across vicuna TQA and Llava VQA, finding that the important attention heads are similar in all scenarios. This result suggests that Llava enhances Vicuna’s existing abilities during visual instruction tuning.



Figure 2: Identifying important image patches for final predictions. mid: log prob increase; right: attn score.

According to these findings, we propose an interpretability tool for users and researchers to understand the important image patches that influence final predictions in Llava’s VQA (Figure 6), which is helpful for understanding visual hallucination. Existing studies typically rely on causal explanations (Rohekar et al., 2024) or average attention scores (Stan et al., 2024) to locate important visual features. However, causal explanation methods require much computational cost, and average attention scores lack strong interpretability. Comparatively, our method computes the log probability increase at each position to identify the important locations in visual features, achieving much lower computational cost than causal explanations and much better interpretability than average attention.

Overall, our contributions are as follows:

1) We investigate the mechanism of TQA in Vicuna and VQA in Llava, finding that the visual embeddings are interpretable when projected into embedding space. We show that the mechanisms of VQA and TQA are similar, and that Llava enhances Vicuna’s existing capabilities during visual tuning.

2) Based on this mechanism analysis, we design an interpretability tool to identify key locations for final predictions, which is valuable for understanding visual hallucinations. Compared to previous methods, our approach provides better interpretability and lower computational cost, making it suitable for real-time interpretations.

2 Mechanism Explorations of TQA and VQA

We investigate the mechanism of TQA and VQA. We introduce the background in Section 2.1, followed by an exploration of the mechanisms of TQA (Section 2.2) and VQA (Section 2.3). Finally, we compare the important attention heads before and after visual instruction tuning in Section 2.4, to explore how Llava obtains its VQA ability.

2.1 Background

Inference pass of decoder-only LLMs. Except the visual encoder and the projection matrix, Llava and Vicuna has the same decoder-only LLM architecture as Llava is a fine-tuned model of Vicuna. So we start from introducing the inference pass of decoder-only LLM with textual inputs. Given $X = [x_1, x_2, \dots, x_T]$ with T tokens, the model predicts an output distribution Y over B tokens in vocabulary V . Every token x_i (at position i) is transformed into a word embedding $h_0^i \in \mathbb{R}^d$ by embedding matrix $E \in \mathbb{R}^{B \times d}$. After that, the word embeddings are sent into $L + 1$ (0th – Lth) transformer layers, where each transformer layer’s output h_i^l (layer l , position i) is the sum of previous layer’s output h_i^{l-1} , this layer’s multi-head self-attention (MHSA) layer output A_i^l , and this layer’s feed-forward network layer (FFN) output F_i^l :

$$h_i^l = h_i^{l-1} + A_i^l + F_i^l \quad (1)$$

To compute the final distribution Y , the final layer’s output at last position h_T^L is multiplied with the unembedding matrix $E_u \in \mathbb{R}^{B \times d}$ and a softmax function over all B tokens:

$$Y = \text{softmax}(E_u h_T^L) \quad (2)$$

As h_T^L is the sum of the last position’s layer outputs and previous studies (Olsson et al., 2022; Wang et al., 2023) find that attention layers play the largest roles for in-context learning, we focus on the last position T ’s attention outputs. Each layer’s MHSA output is computed by the weighted sum of different vectors:

$$A_T^l = \sum_{j=1}^H o_{j,T}^l \quad (3)$$

$$o_{j,T}^l = \sum_{p=1}^T \alpha_{j,T,p}^l \cdot O_j^l V_j^l h_p^{l-1} \quad (4)$$

$$\alpha_{j,T,p}^l = \text{softmax}(Q_j^l h_T^{l-1} \cdot K_j^l h_p^{l-1}) \quad (5)$$

where $o_{j,T}^l$ is the head output in head j , layer l . $\alpha_{j,T,p}^l$ is the attention score at position p , head j , layer l , computed by a softmax function over all positions’ query-key inner products ($Q_j^l h_T^{l-1} \cdot K_j^l h_{pp}^{l-1}$, pp from 1 to T). V_j^l and O_j^l are the value and output matrices in head j , layer l . Generally, A_T^l can be regarded as the weighted sum of $H \times T$ value-output vectors over H heads and T positions, where $O_j^l V_j^l h_p^{l-1}$ is the value-output vector and $\alpha_{j,T,p}^l$ is its weight (attention score).

Identifying important heads and important positions. To explore the mechanism of in-context learning, Yu and Ananiadou (2024a) identify the important heads for the final prediction token b using causal interventions and log probability increase S_j^l of each head output $o_{j,T}^l$:

$$S_j^l = \log(p(b|o_{j,T}^l + h_T^{l-1})) - \log(p(b|h_T^{l-1})) \quad (6)$$

If S_j^l is large, it indicates that the head output $o_{j,T}^l$ contains important information about the final token b . Also, this importance score can be used to identify the important positions in this head by replacing $o_{j,T}^l$ with every position’s weighted value-output vector $\alpha_{j,T,p}^l \cdot O_j^l V_j^l h_p^{l-1}$. They also design logit minus M to evaluate the information storage of $o_{j,T}^l$ for two different tokens $b1$ and $b2$.

$$M = \log(p(b1|o_{j,T})) - \log(p(b2|o_{j,T})) \quad (7)$$

Interpretability analysis: projecting vectors in unembedding space. Geva et al. (2022) and Dar et al. (2022) find that many vectors are interpretable when projecting into the unembedding space E_u by multiplying E_u with the vectors. For instance, $EU_{j,T}^l$ is the projection of $o_{j,T}^l$.

$$EU_{j,T}^l = \text{softmax}(E_u o_{j,T}^l) \quad (8)$$

Yu and Ananiadou (2024a) use this method to analyze the weighted value-output vectors in different positions and find that if S_j^l is large for token b , b usually ranks top in the projection $EU_{j,T}^l$.

2.2 Mechanism Exploration of TQA

In this section, we explore the mechanism of TQA in Vicuna. We analyze 1,000 color-answering sentences of the form ‘[animal] is [color]. Q: What is the color of [animal]? A.’. These sentences

are derived from 1,000 images sampled from the COCO dataset (Lin et al., 2014). For VQA, the input consists of an image and the question ‘Q: What is the color of [animal]? A:’. The only difference between VQA and TQA is that, in the case of TQA, the image is ‘translated’ into a textual context.

Inspired by previous studies (Olsson et al., 2022; Yu and Ananiadou, 2024a), we conclude the mechanism shown in Figure 1(a) for TQA: In shallow layers, the color position extracts the animal information, while the last position encodes the question information. In deep layers’ attention heads, the value-output matrices extract color information from the color position, and the query-key matrices compute the similarity between the last position’s question features and the color position’s animal features. When the question and the textual context refer to the same animal, the similarity score is high, leading to an increased probability of the color token in the final prediction’s distribution.

We identify the most important heads and address four key questions: **a) Does the color position play the largest role in predicting the color token?** **b) Do the value-output matrices extract the color features from the color position?** **c) Does the color position extract the animal features from the textual context?** **d) Does the last position encode the animal features in the question?** To explore these questions, we design two comparison sentences S1: ‘[animal1] is [color]. Q: What is the color of [animal]? A:’ and S2: ‘[animal] is [color]. Q: What is the color of [animal1]? A:’, where [animal1] represents a different animal. We refer to the original sentence ‘[animal] is [color]. Q: What is the color of [animal]? A:’ as S0 and the comparison sentences as S1 and S2. The results are shown in Figure 3.

Evidence a). We calculate the proportion of the log probability increase at the color position relative to the total log probability increase across all positions. The proportion score is 99.82%, indicating that the color position plays the most significant role in predicting the final color token.

Evidence b). We compute the Mean Reciprocal Rank (MRR) of the color token when projecting the color position’s weighted value-output vector (Eq.4) into the unembedding space (Eq.6), yielding an MRR score of 0.463 (equal to ranking 2.16). In comparison, a random color’s MRR score is 0.002, as illustrated in Figure 3(left). The logit difference (Eq.7) between the correct color and a random color

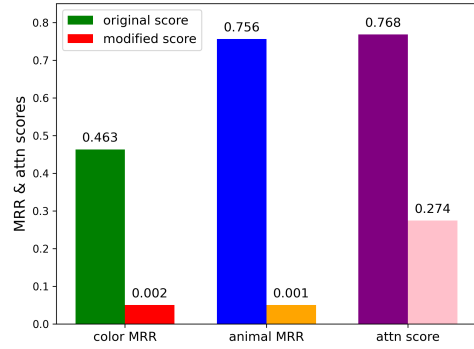


Figure 3: Analysis of color position’s information storage in Vicuna TQA. (left) Color position value-output vector’s information storage for correct color/random color. (mid) Color position layer input vector’s information storage for correct animal/random animal. (right) Color position’s attention score when the question has the same/different animal with the textual context.

at the color position is 2.56. These results confirm that the value-output matrices effectively extract the color features from the color position.

Evidence c). Following Dar et al. (2022), we project the color position’s layer input vector h_p^{l-1} into the unembedding space and calculate the MRR score for the animal tokens [animal] and [animal1]. In S0, the MRR for [animal] is 0.756, while for [animal1], it is 0.001, as shown in Figure 3(mid). The logit difference between [animal] and [animal1] at the color position is 0.32. In S1, the MRR for [animal] is 0.002, the MRR for [animal1] is 0.715, and the logit difference between [animal1] and [animal] is 1.70. These results demonstrate that the layer input vector at the color position, particularly in the most important attention heads, effectively encodes the animal features present in its context.

Evidence d). We calculate the attention scores at the color position for S0, S1, and S2, as queried by the last position. The average attention scores are 0.768, 0.268, and 0.279 for S0, S1, and S2, respectively. When the question involves the same animal as the textual context, the attention score at the color position is high. However, when the animals differ, the attention score drops significantly, as shown in Figure 3(right). This drop in attention scores indicates that the last position encodes the question’s animal features.

Conclusion. Based on the experimental results, we conclude: In shallow layers, the color position extracts the animal features from the textual context (evidence c), while the last position encodes the question features (evidence d). In deep layers’ attention heads, the value-output matrices extract

the color features from the color position (evidence b), and the query-key matrices compute the similarity score between the color position’s animal features and the last position’s question features (evidence d). When the question references the same animal as the textual context, the attention score is significantly high, resulting in the color position’s weighted value-output vector containing substantial color information (evidence a), which is crucial for accurately predicting the color token.

2.3 Mechanism Exploration of VQA

In this section, we aim to explore the mechanism of VQA in Llava. For VQA, we identify the most important heads and address the following questions: **a) What are the most important positions for predicting the correct color?** **b) Do the value-output matrices play a similar role as in TQA?** **c) Do the query-key matrices play a similar role as in TQA?** The results are shown in Figure 4.

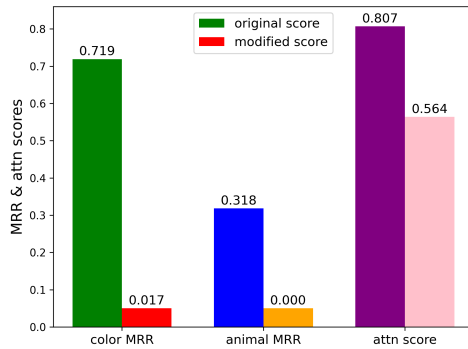


Figure 4: Analysis of top20 important positions’ information in Llava VQA. (left) Top20 position value-output vectors’ information storage for correct color/random color. (mid) Top20 position layer input vectors’ information storage for correct animal/random animal. (right) Top20 positions’ sum attention score when the question has the same/different animal with the image.

Evidence a). We calculate the log probability increase for all positions and visualize these increases as heat maps overlaid on the corresponding images, similar to Figure 2. After randomly sampling 200 cases and analyzing the heat maps on a case-by-case basis, we observe that the positions with the largest log probability increases are those corresponding to image patches related to the animals. Take Figure 2 as an example. When the question is ‘What is the color of the dog?’, the image patches related to the dog’s head exhibit the largest log probability increase. This indicates that the identified image patches contain crucial information for

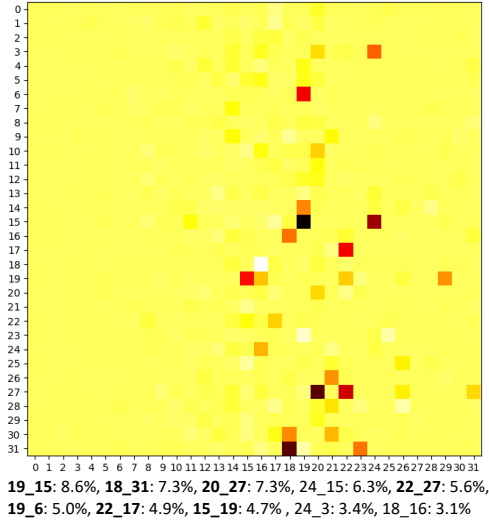
predicting the correct color, demonstrating strong interpretability. This observation inspired the design of an interpretability tool in Section 3, which helps explain why the model arrives at its final predictions. In contrast, the average attention score across all heads typically does not offer the same level of interpretability. Additional examples are provided in Appendix A.

Evidence b). After identifying the most important positions, we analyze whether the value-output matrices extract the color features from the top 20 important positions using a method similar to that used in TQA. When projecting the weighted value-output vector from the color position into the unembedding space, the MRR score for the correct color is 0.719 (equivalent to a ranking of 1.4) and the random color’s MRR is 0.017, as shown in Figure 4(left). The logit difference between the correct color and a random color is 0.09. These results indicate that the value-output matrices effectively extract the color features from the top 20 important positions for the predicted color.

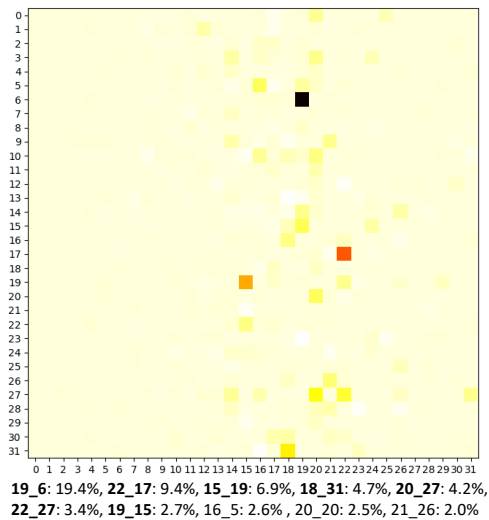
Evidence c) and d). We project the layer inputs of the top 20 important positions into unembedding space and compute the MRR and logit differences between the correct animal and a different animal. The correct animal’s MRR score is 0.318, while the other animal’s MRR score is 0.0004, as shown in Figure 4(mid). The logit difference is 1.53, confirming that the important heads’ layer inputs contain crucial information about the animals. Furthermore, when the animal in the question is replaced with another animal, the attention score at the top 20 positions drops significantly from 0.807 to 0.564 (see Figure 4 right), indicating that the last position encodes information about the question.

Similarity between VQA and TQA. Our findings indicate that the mechanisms underlying VQA and TQA in deep layers are strikingly similar. In both cases, the layer inputs at key positions (the color position in TQA and the animal patch positions in VQA) contain essential information about the animal and color. The value-output matrices are responsible for extracting color information, while the query-key matrices compute the similarity of the animal information between these important positions and the last position. When the attention score is high, more of the color information from these positions is transferred to the last position, which, in turn, increases the likelihood of accurately predicting the color token.

Evidence e). A key difference between VQA and TQA lies in the input embeddings at the 0th layer. Vicuna uses only word embeddings, while Llava combines image and word embeddings. In Vicuna, the color and animal positions encode respective information, with the color token ranking first when projected into the embedding matrix E. To analyze visual embeddings, we projected the top 20 positions into E and computed MRR scores. For the correct color versus a random color, the MRR was 0.455 versus 0.013, and for the correct animal versus a random animal, 0.076 versus 0.0003. At random positions, the MRR scores for the correct color and animal were 0.003 and 0.004, respectively. These results suggest the top 20 positions encode significant information about the correct animal and color, while random positions do not.



(a)



(b)

Figure 5: Important heads in Vicuna (a) and Llava (b).

Results on other questions. To examine whether the mechanism applies to other questions, we replace the original question, ‘What is the color of the [animal]?’, with ‘What is the animal in this picture?’ for comparison. Our findings show that the mechanisms are similar. Analysis of the identified important image patches reveals that these patches are closely associated with the animals. The sum of attention scores on the animal patches is 0.74, indicating that a substantial amount of information is extracted from these patches. Several examples are provided in Appendix A, and additional examples can be explored using the interpretability tool introduced in Section 3.

Conclusion. Based on the experimental results, we conclude the mechanism of VQA illustrated in Figure 1(a). The visual embeddings generated by the projection matrix and the CLIP visual encoder already contain information about the animal and the color (evidence e). This information is propagated through the positions’ residual streams into the deep layers. In the deep layers’ attention heads, the value-output matrices extract color information (evidence b), while the query-key matrices compute the similarity between the animal information and the question information at the last position (evidence c and d). When the similarity is high, the color information related to the animal in the question is more effectively transferred to the last position, thereby increasing the probability of correctly predicting the color token.

2.4 Llava’s Visual Instruction Tuning Enhances Existing Abilities of Vicuna

In this section, we investigate how Llava acquires its VQA capabilities for color prediction. Building on our previous analysis, which highlighted the significant role of deep-layer attention heads in storing VQA abilities, we examine how the important heads evolve after visual instruction tuning. We compute the normalized importance scores for all heads and sort these scores for Vicuna TQA and Llava VQA. Figure 5 displays the importance of all 1,024 heads. The horizontal axis represents the layer number, while the vertical axis denotes the head number. The color intensity indicates the importance of each head, with darker colors signifying greater importance. Additionally, we list the top 10 heads for comparison, where a label like 19_15 refers to the 15th head in the 19th layer. 19.4% is this head’s logit importance.

When comparing Llava VQA with Vicuna TQA, we observe that 7 out of the top 10 attention heads are identical. All the top 7 heads in Llava also appear in the top 10 heads in Vicuna, while the remaining 3 heads rank within the top 20 in the other model. A significant difference is the sharp increase in the importance of head 19_6 (layer 19, head 6), which rises from 5.0% to 19.4%. This suggests that the importance of heads in Llava VQA is more concentrated compared to Vicuna TQA. Based on these results, we conclude that: a) The important heads remain largely consistent between Llava VQA and Vicuna TQA. b) While the most crucial heads are generally similar between Llava VQA and Vicuna TQA, some heads, such as 19_6, become significantly more critical for VQA. c) Visual instruction tuning enhances the existing color-predicting ability of Vicuna’s heads.

Overall, we explore the mechanism of TQA in Vicuna in Section 2.2 and that of VQA in Llava in Section 2.3. We find the mechanism of VQA and TQA is similar in the deep layers’ attention heads. Furthermore, we analyze the projections of visual embeddings in the embedding matrix and find the visual embeddings already contain the information about the animals and the colors. Finally, we compared the most important heads in Vicuna TQA and Llava VQA, and find that Llava enhances the existing heads’ color predicting ability in Vicuna during visual instruction tuning.

3 Interpretability Tool for VQA

In this section, we present our interpretability tool for identifying the key image patches that influence the final predictions.

Interface of the interpretability tool. The interface is illustrated in Figure 6, developed using Gradio (Gradio, 2024). On the left side of the screen, users can upload an image and input a question. On the right side, the first box displays the prediction token, while the second box highlights the top 10 important heads related to the prediction. The third box shows the cropped image (the actual input to Llava) along with the important image patches identified by log probability increase and average attention scores. Each image is divided into 24×24 image patches, with lighter areas indicating a larger score in log probability or attention. Although the visualization appears small within the interface, a button allows users to enlarge the images, resembling the images in Figure 2.

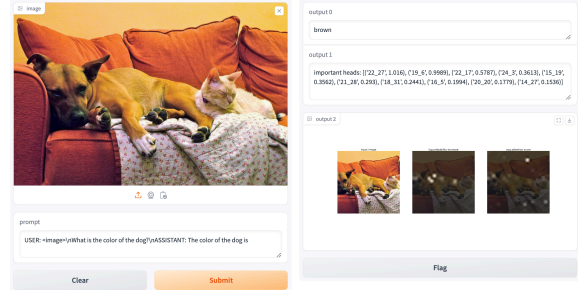


Figure 6: Interface of the interpretability tool. Left: image/question. Right: answer/visualization.

Advantage 1: low computational cost. The first advantage of our method is its low computational cost compared to causal explanations (Rohkar et al., 2024). Causal explanations typically require intervening on each image patch and calculating the impact on the final prediction, necessitating $24 \times 24 + 1$ inference computations. In contrast, our method only requires a single inference computation, with the internal vectors generated during the model’s inference, resulting in minimal additional computation. With our approach, all computations can be completed within 2 seconds with one A100 GPU, offering a promising pathway for real-time explanations.

Advantage 2: good interpretability. Average attention score (Stan et al., 2024) across all attention heads is a widely used method for visual explanation. However, we have observed that this approach does not always provide reasonable explanations. For example, in Figure 2, when asked the question, ‘What is the color of the dog?’, the average attention score is higher on the pillow rather than on the dog itself. This suggests that the average attention score may fail to pinpoint the true reason behind the final prediction. In contrast, our method can accurately identify the important image patches related to the dog. The interpretability of these patches, identified by the log probability increase score, is grounded in the analysis from Section 2, offering a more reliable and robust understanding. More examples demonstrating this trend are provided in Appendix A.

Advantage 3: understanding visual hallucination. Hallucination in vision-language models is a significant issue that has been extensively studied (Li et al., 2023; Zhou et al., 2023; Bai et al., 2024; Liu et al., 2024a). Understanding the precise cause of visual hallucination is crucial. For example, Figure 7 illustrates a hallucination case from

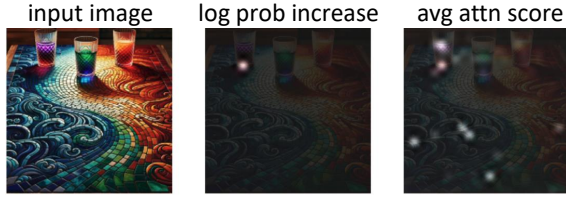


Figure 7: Understanding visual hallucination. Q: What is the color of the left bottle? A: **Red**

Huang et al. (2024). When asked, ‘What is the color of the left bottle?’, Llava incorrectly answers ‘Red’. The exact cause of the hallucination is unclear—whether the model misunderstood the word ‘left’ and provided the color of the right bottle, or if it simply returned the wrong color for the left bottle. Our method’s interpretation clarifies that the model focuses on the bottom of the left bottle, revealing that the hallucination stems from the model failing to consider enough relevant image patches for the color, rather than from a misunderstanding of ‘left’ and ‘right’. Furthermore, our interpretability method can be applied to questions beyond color identification, as provided in Appendix A.

4 Related Works

4.1 Understanding Textual LLMs

Causal intervention (Vig et al., 2020) is a common method for identifying important modules in LLMs (Zhang and Nanda, 2023; Makelov et al., 2023), by computing the change of the final prediction when intervening the module. Meng et al. (2022) find the medium FFN layers in GPT2 store important parameters for knowledge. Stolfo et al. (2023) find similar stages in arithmetic tasks. Wang et al. (2022) proposes activation patching method, using another sentence’s hidden states to replace the original sentence.

A series of studies (Merullo et al., 2023; Lieberum et al., 2023) focus on constructing the internal circuit in transformers from input to output, taking the attention heads and FFN layers as basic units. Elhage et al. (2021) and Olsson et al. (2022) find that the induction heads are helpful for predictions like [A][B] ... [A] => [B]. Hanna et al. (2024) explore how GPT2 computes greater-than algorithm. Gould et al. (2023) find the successor heads help predict the next number like Monday => Tuesday. Wang et al. (2022) study how GPT2 performs the indirect object identification task. Prakash et al. (2024) investigate the circuit

after fine-tuning and find fine-tuning enhances existing mechanisms. Conmy et al. (2023) propose a method to construct the circuits automatically.

Another type of works aim to explore the neurons’ interpretability (Dai et al., 2021; Sajjad et al., 2022; Nanda et al., 2023; Gurnee et al., 2023). Geva et al. (2022) find that FFN neurons are interpretable when projecting into unembedding space. Dar et al. (2022) observe that other vectors are also interpretable in the unembedding space. Yu and Ananiadou (2024b) calculate log probability increase and inner products to identify the important neurons related to the final predictions.

4.2 Understanding Multimodal LLMs

Compared with textual LLMs, only a few studies have investigated the mechanisms of MLLMs. Stan et al. (2024) design an interpretability tool for vision-language models using average attention, relevancy map and causal interpretation. Basu et al. (2024) apply causal intervention methods to understand the information storage and transfer in MLLMs. Tong et al. (2024) study the shortcomings of the visual encoder CLIP. Gandelsman et al. (2023) explore the interpretability of CLIP.

5 Conclusion

In this paper, we utilize mechanistic interpretability methods to investigate the mechanism of VQA in Llava. We find that the mechanism of VQA in Llava is similar to that of TQA in Vicuna. The visual embeddings encode the information of the animals and the colors, and the last position encodes the information of the question in shallow layers. In deep layers’ attention heads, the value-output matrices extract the color information from the visual embeddings, and the query-key matrices compute the similarity between the last position’s question features and the visual positions’ animal features, controlling the probability of the final prediction. Moreover, we find that Llava enhances existing abilities of Vicuna during visual instruction tuning. Based on this analysis, we design an interpretability tool for locating the important image patches related to the final prediction, which has low computational cost, better interpretability and can be utilized for understanding visual hallucination. Overall, our method and analysis is helpful for understanding the mechanism of VQA, paving the way for future studies.

6 Limitations

One limitation of our work is that our experiments are conducted exclusively between Vicuna and Llava. Another limitation is that our work relies on existing interpretability methods in textual LLMs. Nevertheless, we consider the applicability of our methods to MLLMs to be an important and noteworthy finding. Additionally, as the mechanism of vision instruction tuning has not yet been thoroughly studied, our work offers insights that may inspire further exploration in this area.

References

Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. 2024. Hallucination of multimodal large language models: A survey. *arXiv preprint arXiv:2404.18930*.

Samyadeep Basu, Martin Grayson, Cecily Morrison, Besmira Nushi, Soheil Feizi, and Daniela Massiceti. 2024. Understanding information storage and transfer in multi-modal large language models. *arXiv preprint arXiv:2406.04236*.

Tom B Brown. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

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):6.

Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. 2023. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural Information Processing Systems*, 36:16318–16352.

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2021. Knowledge neurons in pretrained transformers. *arXiv preprint arXiv:2104.08696*.

Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. 2022. Analyzing transformers in embedding space. *arXiv preprint arXiv:2209.02535*.

Xiang Deng, Vasilisa Bashlovkina, Feng Han, Simon Baumgartner, and Michael Bendersky. 2023. Llm to the moon? reddit market sentiment analysis with large language models. In *Companion Proceedings of the ACM Web Conference 2023*, pages 1014–1019.

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1(1):12.

Yossi Gandelsman, Alexei A Efros, and Jacob Steinhardt. 2023. Interpreting clip’s image representation via text-based decomposition. *arXiv preprint arXiv:2310.05916*.

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realltoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*.

Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting recall of factual associations in auto-regressive language models. *arXiv preprint arXiv:2304.14767*.

Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. 2022. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. *arXiv preprint arXiv:2203.14680*.

Rhys Gould, Euan Ong, George Ogden, and Arthur Conmy. 2023. Successor heads: Recurring, interpretable attention heads in the wild. *arXiv preprint arXiv:2312.09230*.

Gradio. 2024. Gradio. <https://www.gradio.app>.

Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*.

Michael Hanna, Ollie Liu, and Alexandre Variengien. 2024. How does gpt-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. *Advances in Neural Information Processing Systems*, 36.

Wen Huang, Hongbin Liu, Minxin Guo, and Neil Zhenqiang Gong. 2024. Visual hallucinations of multimodal large language models. *arXiv preprint arXiv:2402.14683*.

Hadas Kotek, Rikker Dockum, and David Sun. 2023. Gender bias and stereotypes in large language models. In *Proceedings of the ACM collective intelligence conference*, pages 12–24.

Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.

Tom Lieberum, Matthew Rahtz, János Kramár, Neel Nanda, Geoffrey Irving, Rohin Shah, and Vladimir Mikulik. 2023. Does circuit analysis interpretability scale? evidence from multiple choice capabilities in chinchilla. *arXiv preprint arXiv:2307.09458*.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. 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.

752	Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen,	Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya	807
753	Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li,	Sachan. 2023. A mechanistic interpretation of arith-	808
754	and Wei Peng. 2024a. A survey on hallucination	metic reasoning in language models using causal me-	809
755	in large vision-language models. <i>arXiv preprint</i>	diation analysis. <i>arXiv preprint arXiv:2305.15054</i> .	810
756	<i>arXiv:2402.00253</i> .		
757	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	Yiming Tan, Dehai Min, Yu Li, Wenbo Li, Nan Hu,	811
758	Lee. 2024b. Visual instruction tuning. <i>Advances in</i>	Yongrui Chen, and Guilin Qi. 2023. Can chatgpt	812
759	<i>neural information processing systems</i> , 36.	replace traditional kbqa models? an in-depth analysis	813
760	Aleksandar Makelov, Georg Lange, Atticus Geiger, and	of the question answering performance of the gpt llm	814
761	Neel Nanda. 2023. Is this the subspace you are look-	family. In <i>International Semantic Web Conference</i> ,	815
762	ing for? an interpretability illusion for subspace ac-	pages 348–367. Springer.	816
763	tivation patching. In <i>The Twelfth International Confer-</i>		
764	<i>ence on Learning Representations</i> .	Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma,	817
765	Kevin Meng, David Bau, Alex Andonian, and Yonatan	Yann LeCun, and Saining Xie. 2024. Eyes wide	818
766	Belinkov. 2022. Locating and editing factual associ-	shut? exploring the visual shortcomings of multi-	819
767	ations in gpt. <i>Advances in Neural Information Pro-</i>	modal llms. In <i>Proceedings of the IEEE/CVF Con-</i>	820
768	<i>cessing Systems</i> , 35:17359–17372.	<i>ference on Computer Vision and Pattern Recognition</i> ,	821
769	Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. 2023.	pages 9568–9578.	822
770	Circuit component reuse across tasks in transformer	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	823
771	language models. <i>arXiv preprint arXiv:2310.08744</i> .	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	824
772	Neel Nanda, Senthoooran Rajamanoharan, János Kramár,	Baptiste Rozière, Naman Goyal, Eric Hambro,	825
773	and Rohin Shah. 2023. Fact finding: At-	Faisal Azhar, et al. 2023. Llama: Open and effi-	826
774	tempting to reverse-engineer factual recall on the	cient foundation language models. <i>arXiv preprint</i>	827
775	neuron level. URL https://www. alignmentfo-	<i>arXiv:2302.13971</i> .	828
776	<i>rum. org/posts/iGuwZTHWb6DFY3sKB/fact-finding-</i>	Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov,	829
777	<i>attempting-to-reverse-engineer-factual-recall</i> .	Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart	830
778	Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas	Shieber. 2020. Investigating gender bias in language	831
779	Joseph, Nova DasSarma, Tom Henighan, Ben Mann,	models using causal mediation analysis. <i>Advances</i>	832
780	Amanda Askell, Yuntao Bai, Anna Chen, et al. 2022.	<i>in neural information processing systems</i> , 33:12388–	833
781	In-context learning and induction heads. <i>arXiv</i>	12401.	834
782	<i>preprint arXiv:2209.11895</i> .	Kevin Wang, Alexandre Variengien, Arthur Conmy,	835
783	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Buck Shlegeris, and Jacob Steinhardt. 2022. In-	836
784	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	terpretability in the wild: a circuit for indirect ob-	837
785	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	ject identification in gpt-2 small. <i>arXiv preprint</i>	838
786	2022. Training language models to follow instruc-	<i>arXiv:2211.00593</i> .	839
787	tions with human feedback. <i>Advances in neural in-</i>	Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou,	840
788	<i>formation processing systems</i> , 35:27730–27744.	Fandong Meng, Jie Zhou, and Xu Sun. 2023. Label	841
789	Nikhil Prakash, Tamar Rott Shaham, Tal Haklay,	words are anchors: An information flow perspective	842
790	Yonatan Belinkov, and David Bau. 2024. Fine-tuning	for understanding in-context learning. <i>arXiv preprint</i>	843
791	enhances existing mechanisms: A case study on en-	<i>arXiv:2305.14160</i> .	844
792	tity tracking. <i>arXiv preprint arXiv:2402.14811</i> .	Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert	845
793	Raanan Y Rohekar, Yaniv Gurwicz, and Shami Nisi-	Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu,	846
794	mov. 2024. Causal interpretation of self-attention in	Da Huang, Denny Zhou, et al. 2023. Larger language	847
795	pre-trained transformers. <i>Advances in Neural Infor-</i>	models do in-context learning differently. <i>arXiv</i>	848
796	<i>mation Processing Systems</i> , 36.	<i>preprint arXiv:2303.03846</i> .	849
797	Hassan Sajjad, Nadir Durrani, and Fahim Dalvi. 2022.	Changrong Xiao, Sean Xin Xu, Kunpeng Zhang, Yufang	850
798	Neuron-level interpretation of deep nlp models: A	Wang, and Lei Xia. 2023. Evaluating reading com-	851
799	survey. <i>Transactions of the Association for Computa-</i>	prehension exercises generated by llms: A showcase	852
800	<i>tional Linguistics</i> , 10:1285–1303.	of chatgpt in education applications. In <i>Proceed-</i>	853
801	Gabriela Ben Melech Stan, Raanan Yehezkel Rohekar,	<i>ings of the 18th Workshop on Innovative Use of NLP</i>	854
802	Yaniv Gurwicz, Matthew Lyle Olson, Anahita Bhi-	<i>for Building Educational Applications (BEA 2023)</i> ,	855
803	wandiwalla, Estelle Aflalo, Chenfei Wu, Nan Duan,	pages 610–625.	856
804	Shao-Yen Tseng, and Vasudev Lal. 2024. Lvlm-	Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan	857
805	intrepret: An interpretability tool for large vision-	Ning, and Li Yuan. 2023. Llm lies: Hallucinations	858
806	language models. <i>arXiv preprint arXiv:2404.03118</i> .	are not bugs, but features as adversarial examples.	859
		<i>arXiv preprint arXiv:2310.01469</i> .	860

Zeping Yu and Sophia Ananiadou. 2024a. How do large language models learn in-context? query and key matrices of in-context heads are two towers for metric learning. *arXiv preprint arXiv:2402.02872*.

Zeping Yu and Sophia Ananiadou. 2024b. Neuron-level knowledge attribution in large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3267–3280.

Duzhen Zhang, Yahan Yu, Chenxing Li, Jiahua Dong, Dan Su, Chenhui Chu, and Dong Yu. 2024. Mm-llms: Recent advances in multimodal large language models. *arXiv preprint arXiv:2401.13601*.

Fred Zhang and Neel Nanda. 2023. Towards best practices of activation patching in language models: Metrics and methods. *arXiv preprint arXiv:2309.16042*.

Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. 2023. Analyzing and mitigating object hallucination in large vision-language models. *arXiv preprint arXiv:2310.00754*.

A Appendix A: Example Images' Interpretability

We provide more examples (Figure 8-15) to verify the usage of our interpretability tool. Our method is not only suitable for identifying the important image patches about color questions, but also for other questions. The questions are listed in the titles of the following images, where the answers are marked as bold. In each figure, the left picture is the input image, the mid picture is the visualization of our method, and the right picture is the visualization of average attention score.



Figure 8: Q: What is the color of the cat? A: The color of the cat is **white**



Figure 9: Q: What is the color of the pillow? A: The color of the pillow is **orange**

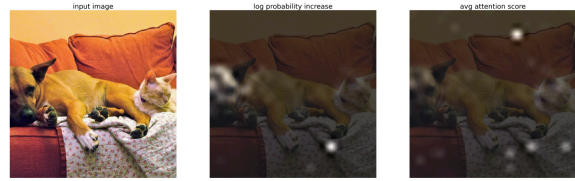


Figure 10: Q: What is the left animal? A: The left animal is a **dog**



Figure 11: Q: What is the right animal? A: The right animal is a **cat**



Figure 12: Q: What is in the painting? A: The painting features a **woman**



Figure 13: Q: What is in the painting? A: The painting features a **dog**

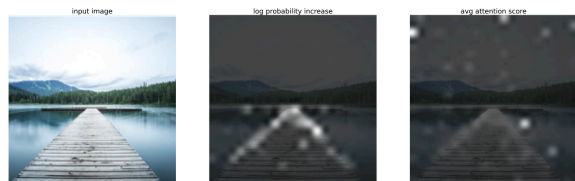


Figure 14: Q: What is in the picture? A: The picture features a **pier**



Figure 15: Q: What is the table made of? A: The table is made of **glass**