# Aerial Mirage: Unmasking Hallucinations in Large Vision Language Models - Supplementary Section

## 1. Caption Annotation Process for AeroCaps

A team of 10 annotators aged between 20 and 30 years were selected from different technical backgrounds. These annotators were students from Computer Science, Electrical Engineering, Civil Engineering, Mechanical Engineering, Electronics and Communication, etc., to ensure a wide range of perspectives. The following set of written guidelines was provided to the annotators:

1. Describe the image (preferably in a single sentence) while covering all the important objects and their relationships with each other.
2. The caption length should be at least 5 words.
3. Give captions clearly and descriptively that convey the relevant information about the images.
4. Consider the relationships between different elements in the image and describe their relative positioning.
5. Refrain from using vague concepts like "large", "tall", or "many" wherever you deem fit.
6. Avoid using specific directions such as north, south, east, or west to maintain image neutrality.
7. Annotations should be based on the annotators' own understanding of the image without relying on AI models for caption generation.

As the annotators were from a known student pool, we held a meeting with them to explain the annotation guidelines. With this, each image in the dataset was annotated with four captions, each provided by a different annotator. This approach provided multiple perspectives, sentence formation styles and interpretations for each image, resulting in a diverse set of captions. Giving multiple captions to each image increases the richness of the dataset and facilitates a deeper understanding of the visual content of images. The maximum caption length was 40 words.

## 2. More examples of image-captions from AeroCaps Dataset

Figure 1 contains some images-captions from our proposed AeroCaps dataset.

Figure 1. Sample image-captions from the proposed AeroCaps dataset. Each image contains four caption annotations.

## 3. Details regarding Evaluation of advanced LLMs for judgement (Section-7 of main paper)

The following prompt is given to GPT-3.5-Turbo [6], GPT-4 [1] and Gemini-pro [7] and Gemini-pro-vision [7] for evaluating the captions generated by LLaVA [5] on 200 sampled COCO-2017 (val) images [2, 4]:

We only give the text prompt mentioned above for GPT-3.5-Turbo, GPT-4, and Gemini-pro. For Gemini-pro-vision, we also provide the image for reference and add a sentence in the prompt: *"You are also given an image as reference. Any inconsistencies in the "Candidate Description" that doesn't match the visual content of the image is considered hallucination."* We manually evaluate the answers given by these models as correct or incorrect (0 - correct, 1 - incorrect).

As discussed in Section-7 of the main paper, these models often make mistakes in identifying objects or attributes as hallucinated despite their presence in reference captions and bounding boxes. An example is shown in Figure 2. The contexts given for evaluation and the responses by "judge" models are in Figure 2a and Figure 2b, respectively. We clarify the mistakes below:

- We find that LLaVA's generated answer is mostly correct, except that it says "a" giraffe while there are "two" giraffes. So, this is *count hallucination.* However, none of the models point out the count error. LLaVA also fails to point out the "gazelle" *(Missing Object).*

- GPT-3.5-Turbo and GPT-4 point out that the "giraffe drinking water" is a hallucination, which is incorrect since it is mentioned in the references. GPT-4 also claims that there is "no watering hole" or "pond", which is also incorrect. It is mentioned in the references. Finally, both these models have given a very high score of 8 and 9.5 for hallucinations where the pointed hallucinations are wrong.
  For the count error, a human would typically give a very low hallucination score.

- Gemini-pro also mentions that "watering hole" is a hallucination.

- Gemini-pro-vision also makes the same mistakes. It even adds that LLaVA mentions "impala", but LLaVA actually missed the presence of "gazelle". Despite being given both references and the image, it doesn't identify the presence of "watering hole" and "giraffe is drinking water from the pond".

So, we mark all the models' answers as incorrect for this example.

As discussed in Section-7 in the main manuscript, we observed that LLaVA frequently includes additional information in its captions that is not mentioned in the GT captions or bbox annotations. Upon inspecting the images, we found this additional information is often accurate. However, the "judge" models classify these as hallucinations and penalize the score because the information is not in the reference OR not mentioned explicitly in the references. During manual evaluation, we do not mark these cases as incorrect, as they are consistent with the given image and the information in the prompts. Figure 3a shows an example, where LLaVA's caption about the image is correct about the player's dress colour and another player standing behind the batter. However, this information is not clear from the references, and the "judge" models point them as hallucinations, as shown in Figure 3b. So, we mark the answers given by "judge" models as correct. However, LLaVA hallucinates in "several other people in the scene", "a few more players scattered around the field", and "A baseball glove can be seen on the ground", which is correctly pointed out by GPT-4, Gemini-pro and Gemini-pro-vision, but not by GPT-3.5-Turbo.

**Context given in Prompt**

**Reference captions:** ['A family of giraffes by a few zebras and a gazelle.', 'A group of giraffes <u>drink water</u> in the wilderness.', 'Giraffe and zebra grazing in a field next to plants.', 'Several giraffes and zebras walking in a dirt field.', 'A group of zebras and giraffes by a <u>water hole</u>.']

**Bounding boxes (bbox):** ['zebra': [170.42, 236.77, 161.45, 99.19], 'zebra': [4.06, 232.31, 181.99, 101.91], 'zebra': [103.39, 191.55, 54.15, 41.1], 'giraffe': [261.02, 5.71, 350.27, 369.46], 'giraffe': [338.15, 139.91, 96.61, 231.81]]

**Candidate Description (generated by LLaVA):** ["The image features a group of zebras and a giraffe gathered around a watering hole. The zebras are spread out in the scene, with some standing closer to the water and others further away. The giraffe is standing next to the zebras, drinking water from the pond. The animals appear to be enjoying the water and each other's company."]

(a) The contexts used to prompt the "judge" models. The captions and bounding boxes are the groundtruths from COCO. The "Candidate Description" is the caption generated by LLaVA for the given image. Note that the image is used to prompt Gemini-pro-vision only along with the texts, and NOT the other three models. We show it here as a reference.

**Judge LLM's Responses**

**GPT-3.5-Turbo Response:**
I would rate the "Candidate description" a score of 8 for hallucination.
Mistakes:
1. The candidate description mentions that the giraffe is drinking water from the pond, which is not explicitly stated in the reference captions or visible in the image. This is a hallucination.
2. The candidate description mentions that the zebras are spread out in the scene, with some standing closer to the water and others further away. However, this level of detail is not specified in the reference captions or evident in the image. This is another hallucination.

**GPT-4 Response:**
Score: 9.5
Mistakes:
The candidate description includes a lot of hallucinations. Firstly, there is no mention of a watering hole in the image according to bounding box information. The narrative of zebras being spread out and standing closer or farther from the water is also not supported by the bounding box data. It also hallucinates the giraffe drinking water from a non-existing pond. Lastly, candidate description says about the animals enjoying water and each other's company, however, this action is not specified or supported by the bounding box information or reference captions. Only the presence giraffes and zebras can be confirmed from the ground truth annotations.

**Gemini-pro Response:**
Score: 4
The candidate description contains some hallucination.
1. It mentions that the animals are gathered around a watering hole, but there is no ground truth bounding box for a water hole.
2. It describes that the animals appear to be enjoying the water and each other's company, but such kind of interaction is not indicated in the ground truth.

**Gemini-pro-vision Response:**
Score: 3
There are several factual mistakes in the "Candidate description".
1. The candidate description hallucinates the presence of a "watering hole".
2. It hallucinates that the "giraffe is drinking water from the pond".
3. It hallucinates the presence of an "impala".
4. It hallucinates that "the animals appear to be enjoying the water and each other's company".

(b) Answers given by "judge" models for the given prompt. The mistakes made by them are marked in red.

Figure 2. Example showing that the "judge" models make mistakes in identifying objects or attributes as hallucinated despite their presence in reference captions or bounding boxes.

**Judge LLM's Responses**

**GPT-3.5-Turbo Response:**
I would rate the "Candidate description" with a score of 2 for hallucination.
Mistakes:
1. The candidate description mentioned that the batter is wearing a black shirt and white pants, however, the ground truth annotations did not specify any clothing colors for the player. This is a hallucination as it adds details not present in the image.
2. The candidate description mentioned a player standing behind the batter, but there is no bounding box for this specific player in the image. This could be considered as an imagined detail not supported by the ground truth annotations.

**GPT-4 Response:**
Score: 3
Mistakes:
1. The candidate description mentions that the baseball player is wearing a black shirt and white pants, but there is no such information given in the bounding boxes or reference captions.
2. It mentions about a player standing behind the batter, this detail is not confirmed by the bounding boxes or reference captions.
3. The description states that a baseball glove can be seen on the ground. There's no ground-level indication from the bounding box information.
4. There is a mention of several people present in the scene while only three person bounding boxes are given. More so, no information about these additional people is provided in the reference captions.
5. Candidate description says the scene is filled with action and lively, which is a subject interpretation and can be interpreted differently.

**Gemini-pro Response:**
Hallucination score: 1
Mistakes:
1. The candidate description mentions "a player standing behind the batter", but there is no such person in the image.
2. The candidate description mentions "a few more players scattered around the field", but there are no other players visible in the image.
3. The candidate description mentions "A baseball glove can be seen on the ground", but the baseball glove is actually in the batter's hand.

**Gemini pro-vision Response:**
Score: 2
The description is mostly accurate, but there are a few mistakes:
1. The batter is not wearing a black shirt and white pants, he is wearing a black and yellow shirt and white pants.
2. There is only one other player visible in the image, not several.
3. The baseball glove is not on the ground, it is in the hand of the catcher.

(b) Answers given by "judge" models about Figure 3a. Some additional correct information generated by LLaVA, that are pointed out as hallucinations are marked in blue.

Figure 3. Example showing additional information generated by LLaVA that is not mentioned explicitly in the GT captions or bbox annotations, but is actually correct. They are pointed out as mistakes by the "judge" models.

The image depicts a cozy bedroom with a large bed situated in the center of the room. The bed is covered with a brown blanket, and there is a clock on the wall above it. The room also features a fireplace, which adds warmth and ambiance to the space. In addition to the bed, there are two chairs in the room, one located near the left side of the bed and the other on the right side. A television is mounted on the wall above the right chair, providing entertainment for the occupants. The room is well-lit, with natural light coming in through the windows, creating a comfortable and inviting atmosphere.

The image features two people sitting on a bench on the side of a street. They appear to be taking a break or waiting for something. Behind them, there are two motorcycles parked on the sidewalk, one closer to the bench and the other slightly further away. In addition to the people and motorcycles, there are a few other objects in the scene. A handbag is placed on the ground near the bench, and a backpack is located further away from the bench. A chair is also visible in the scene, placed near the bench.

The image shows two young boys wearing baseball uniforms, standing next to a fence. One of the boys is holding a baseball bat, while the other is wearing a backpack. They appear to be preparing for a baseball game or practice. There are several baseball bats visible in the scene, with one placed near the boys and others scattered around the area. A baseball glove can also be seen, indicating that they are well-equipped for the game. The boys seem to be focused on their game or practice, ready to engage in the sport.

The image features a white bus with green and red accents parked at a bus stop. The bus is stopped at the curb, and its doors are open, waiting for passengers to board. There are several people in the vicinity of the bus, some standing close to the bus stop and others further away. In addition to the bus, there are two bicycles in the scene. One bicycle is parked near the right side of the image, while the other is located on the far right side. A car can also be seen in the background, parked further away from the bus stop.

Figure 4. Examples of LLaVA's hallucination on sample images from COCO-2017 (val). Inappropriate details are marked in red.

## 4. Some examples of LLaVA's hallucination on COCO

We evaluated LLaVA's generated captions for the prompt: *"Describe the image briefly"*, on the 200 sampled COCO-2017 (val) images (as mentioned in Section-7 of main paper). We mark 1 for hallucination and 0 for no hallucination. We found that it hallucinated in **77.11**% of the responses. Thus, LLaVA hallucinates not only on aerial imagery, as discussed in our study, but also on the COCO dataset it was trained on. Figure 4 and Figure 5 demonstrates LLaVA's hallucination on images from the COCO-2017 validation set.

## 5. Example of Labelled Illusion Dataset (LID)

Figure 6 demonstrates an example image from VisDrone (val) [8] where InstructBLIP [3] generates the caption. The annotations for hallucinations are given in our proposed Labelled Illusion Dataset (LID). For 'Any hallucina-

The image features a cozy living room with a television placed on a stand in the corner. The TV is turned on, displaying a football game. A comfortable chair is positioned in front of the television, providing a perfect spot for watching the game. In addition to the television and chair, the room is filled with numerous books scattered throughout the space. Some books are placed on the floor, while others are stacked on various surfaces. The room also has a window with curtains, which adds to the overall ambiance of the living room.

The image features a unique statue of a person holding a baby, with the statue being covered in teddy bears. The teddy bears are scattered all over the statue, creating a heartwarming scene. The statue is placed on a patch of grass, and there are several teddy bears of various sizes surrounding it. In addition to the statue, there is a cross nearby, which adds a religious element to the scene. The combination of the statue, teddy bears, and the cross creates an interesting and thought-provoking display.
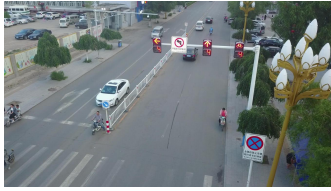
The image features a red stop sign with the word "SOS" written on it, placed at an intersection. The stop sign is positioned on a pole, and it is surrounded by a few trees and bushes.
In the background, there are several cars and trucks visible, indicating that the area is a busy street. Some of the vehicles are parked, while others are in motion, creating a lively atmosphere.

Figure 5. More examples of LLaVA's hallucination on sample images from COCO-2017 (val). Inappropriate details are marked in red.

tion', 'Object Hallucination', 'Position hallucination', 'Count hallucination', and 'Missing objects', the *presence of hallucination* is denoted by 1 and it's absence as 0.

## 6. Few-shot prompting

Figure 7 illustrates a two-shot prompt example. We manually provided two similar images with their corresponding captions as prompts. LLaVA has to generate caption for the third image. LLaVA's generated answers were directly copied from these examples, showing no improvement with this method.

**InstructBLIP's Output:**
The image depicts a busy city street with multiple vehicles, including cars, motorcycles, and bicycles. There are at least 10 cars in the scene, some of which are parked on the side of the road. Additionally, there are several motorcycles and bicycles traveling through the street. Some of the motorcycles are parked on the side of the road, while others are moving along the street. The bicycles are scattered throughout the scene, with some ridden by individuals and others left on the side of the road. There are also several pedestrians walking along the street, adding to the lively atmosphere.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Annotations of hallucination (from LID)**:
Any hallucination:1, Object Hallucination:1, Position hallucination:1, Count hallucination:0, Missing objects:1, Hallucination level:3, Categories hallucinated: 'motorcycles', Num of categories hallucinated:1, Num of missing categories:2, Names of categories missed: ['trees', 'street lights'], Hallucinated word-position index: motorcycles-55, Total number of tokens:114.

Figure 6. An Example of the hallucination annotations from LID (lower part). The description (upper part) is generated by InstructBLIP on an image from VisDrone (val). (InstructBLIP, VisDrone, long)



**Prompt:**
USER: ⟨*image*⟩ Describe the image briefly.

ASSISTANT: The image depicts a parking lot filled with various cars parked in rows. There are at least 13 cars visible in the scene, with some parked closer to the foreground and others further back. The parking lot is situated next to a street, and a tree can be seen in the background, providing a touch of greenery to the urban setting. In addition to the cars, there are two people in the scene. The overall scene showcases a busy urban environment with a mix of vehicles and a hint of nature.

USER: ⟨*image*⟩ Describe the image briefly.

ASSISTANT: The image depicts a parking lot filled with various cars parked in rows. Some cars are parked closer to the foreground while others are further back. The parking lot is located near a building, possibly a shopping center or a business establishment. In addition to the cars, there are a few people in the scene who might be walking around or attending to their vehicles. The overall atmosphere of the scene suggests a busy day at the parking lot, with many cars parked and people going about their activities.

USER: ⟨*image*⟩ Describe the image briefly.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**LLaVA's Output:**
ASSISTANT: The image depicts a parking lot filled with various cars parked in rows. There are at least 13 cars visible in the scene, with some parked closer to the foreground and others further back. The parking lot is situated near a building, possibly a shopping center or a business establishment. In addition to the cars, there are a few people in the scene, likely walking around or attending to their vehicles. The overall atmosphere of the scene suggests a busy day at the parking lot, with many cars parked and people going about their activities.

Figure 7. Example of few-shot with LLaVA. The first two images and the corresponding prompt was given as input. LLaVA generated a caption for the third image. It seems they were directly copied from the prompt.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2

[2] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 2

[3] Wenliang Dai, Junnan Li, DONGXU LI, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung,

and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *Advances in Neural Information Processing Systems*, volume 36, 2023. 6

[4] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 2

[5] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26296–26306, June 2024. 2

[6] OpenAI. Gpt-3.5 turbo, 2023. 2

[7] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 2

[8] Pengfei Zhu, Longyin Wen, Dawei Du, Xiao Bian, Heng Fan, Qinghua Hu, and Haibin Ling. Detection and tracking meet drones challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7380–7399, 2021. 6