

An Early Evaluation of GPT-4V(ision)

Yang Wu Shilong Wang Hao Yang Tian Zheng
Hongbo Zhang Yanyan Zhao* Bing Qin

Research Center for Social Computing and Information Retrieval

Harbin Institute of Technology

{yw, shilongwang, hyang, tzheng, hbzhang, yyzhao, qinb}@ir.hit.edu.cn

Abstract

In this paper, we evaluate different abilities of GPT-4V including visual understanding, language understanding, visual puzzle solving, and understanding of other modalities such as depth, thermal, video, and audio. To estimate GPT-4V’s performance, we manually construct 656 test instances and carefully evaluate the results of GPT-4V. The highlights of our findings are as follows: (1) GPT-4V exhibits impressive performance on English visual-centric benchmarks but fails to recognize simple Chinese texts in the images; (2) GPT-4V shows inconsistent refusal behavior when answering questions related to sensitive traits such as gender, race, and age; (3) GPT-4V obtains worse results than GPT-4 (API) on language understanding tasks including general language understanding benchmarks and visual commonsense knowledge evaluation benchmarks; (4) Few-shot prompting can improve GPT-4V’s performance on both visual understanding and language understanding; (5) GPT-4V struggles to find the nuances between two similar images and solve the easy math picture puzzles; (6) GPT-4V shows non-trivial performance on the tasks of similar modalities to image, such as video and thermal. Our experimental results reveal the ability and limitations of GPT-4V and we hope our paper can provide some insights into the application and research of GPT-4V¹.

1 Introduction

GPT-4V has shown remarkable capabilities on a wide of tasks (Yang et al., 2023). However, the performance of GPT-4V has not been quantitatively studied. In this paper, we manually construct 656 test examples to quantitatively assess GPT-4V’s abilities and seek answers to the following intriguing questions.

* Corresponding Author

¹Our data are available at <https://github.com/albertwy/GPT-4V-Evaluation>

1. What is the performance of GPT-4V on visual-centric benchmarks such as image captioning and visual question answering? Can GPT-4V surpass the current SOTA multimodal LLMs such as Qwen-VL-Chat (Bai et al., 2023) on these benchmarks? (**Visual Understanding**)
2. After being equipped with visual perception, can GPT-4V maintain its language understanding performance and better capture visual commonsense knowledge and world knowledge (specifically physics knowledge)? (**Language Understanding**)
3. Can GPT-4V benefit from exemplars? (**Visual Understanding, Language Understanding**)
4. How to evaluate multimodal LLMs given the observation that multimodal LLMs have achieved really high performance on the current benchmarks? (**Visual Puzzle Solving**)
5. Can GPT-4V perceive other modalities such as depth, thermal, video, and audio? (**Understanding of Other Modalities**)

We conduct extensive evaluation of GPT-4V and the results not only reveal GPT-4V’s strong abilities, but also point to the following issues that should be addressed in future studies.

1. GPT-4V tends to generate verbose responses, even when provided with exemplars that have short answers, which makes it hard to accurately assess GPT-4V’s performance using current automatic metrics. For example, the CIDEr scores on Nocaps (Agrawal et al., 2019) and Flickr30K (Young et al., 2014)² obtained by GPT-4V are close to 0.
2. GPT-4V shows inconsistent refusal behavior when answering questions related to sensitive

²We utilize the released code by Bai et al. (2023) to estimate the results.

traits such as gender, race, and age. This issue causes an obvious performance drop on GQA. Future research should address this issue carefully when comparing GPT-4V with other multimodal LLMs.

3. GPT-4V performs very well with English text recognition, yet it cannot recognize Chinese texts in images.
4. GPT-4V struggles to solve the easy math picture puzzle (grades five level) although it exhibits strong performance on much harder textual math datasets such as SAT math (OpenAI, 2023).
5. The current version of GPT-4V does not support interleaved images and texts and can only accept a maximum of four images. These constraints limit the design space of prompts.



Figure 1: An example image from GQA.

2 Visual Understanding

We evaluate GPT-4V on various visual-centric benchmarks such as image captioning and visual question answering to assess its visual understanding ability. Following Qwen-VL-Chat (Bai et al., 2023), we choose Nocaps (Agrawal et al., 2019) and Flickr30K (Young et al., 2014) as the evaluation datasets for image captioning. As for visual question answering, we evaluate GPT-4V on VQAv2 (Goyal et al., 2016), OKVQA (Marino et al., 2019), GQA (Hudson and Manning, 2019), ScienceQA (Lu et al., 2022), and Vizwiz VQA (Gurari et al., 2018).

Metric GPT-4V always tends to generate verbose responses, which makes it hard to accurately evaluate GPT-4V’s performance using current automatic metrics. For example, given the image shown in Figure 1, we ask GPT-4V to find out which kind of watercraft is underneath the airplane and GPT-4V answers correctly with “the watercraft

Table 1: Human evaluation for GPT-4V and Qwen-VL-Chat (Zero-shot).

Task	Dataset	GPT-4V	Qwen-VL-Chat
Image	Nocaps	17/20	15/20
Captioning	Flickr30K	19/20	17/20
	VQAv2	16/20	15/20
Visual	OKVQA	18/20	16/20
	GQA	11/20	15/20
Question Answering	ScienceQA	17/20	14/20
	VizWiz	17/20	14/20
	OCR-VQA	20/20	19/20

Table 2: Automatic evaluation for GPT-4V and Qwen-VL-Chat (Zero-shot). We do not carefully adjust the prompts and we acknowledge that using task-specific prompts to control the output formats could be helpful.

Task	Dataset	Metric	GPT-4V	Qwen-VL-Chat
Image	Nocaps	SPICE	15.9	16.2
Captioning	Flickr30K	SPICE	15.2	17.3
	VQAv2	VQA Score	0.0	85.0
Visual	OKVQA	VQA Score	5.0	38.0
	GQA	EM Accuracy	10.0	40.0
Question Answering	ScienceQA	Accuracy	85.0	70.0
	VizWiz	VQA Score	23.0	36.0
	OCR-VQA	EM Accuracy	0.0	70.0

underneath the airplane are sailboats”. However, if we utilize EM Accuracy as the metric, which is adopted by Qwen-VL-Chat for GQA, the answer of GPT-4V will be considered as incorrect given the ground truth is “sailboat”. To address this problem, we manually evaluate the results of GPT-4V and Qwen-VL-Chat. The automatic evaluation results are also reported to reveal the limitation of current automatic metrics. Besides, we utilize SPICE (Anderson et al., 2016) instead of CIDEr (Vedantam et al., 2015) as the metric for image captioning, because we find that the current implementation of CIDEr adopted by Qwen-VL-Chat gives a large penalty to the difference between candidate and reference sentence lengths. In our experiments, the CIDEr scores obtained by GPT-4V are close to 0 on Nocaps and Flickr30K.

Zero-shot Results. We randomly sample 20 test instances for each dataset and manually evaluate GPT-4V’s performance. We list the human and automatic evaluation results in Table 1 and Table 2. There are mainly two findings. (1) **GPT-4V performs really well on various tasks.** GPT-4V can generate detailed and coherent descriptions for the given images and answer the questions based on the images. GPT-4V is able to accurately recognize the English characters in the images, achieving 100% accuracy on OCR-VQA. We also observe

Image Captioning

Prompt:

Describe this image



GPT-4V:

The image depicts a bride in her wedding attire. She is wearing a white wedding gown with intricate beading and detailing on the bodice. The bride is also adorned with a veil that cascades down her back, and she holds a bouquet of white roses in one hand. Her hair is styled in an updo with a slight bouffant at the crown. Beside her, partially visible, are a few people, including an elderly woman in a pink dress and a man in a suit. In the background, there's a stone wall and blossoming trees, suggesting the photograph might have been taken in a garden or a courtyard setting.

Qwen-VL-Chat:

A bride holding a bouquet of flowers in her hand.

Ground Truth:

A bride in her wedding gown holding a bouquet.

Figure 2: The caption generated by GPT-4V is more detailed and specific than the ground truth. GPT-4V can capture the details in the whole image (marked with red color) and also attempts to infer the location where the photo was taken (marked with blue color).

that GPT-4V exhibits superior performance than Qwen-VL-Chat on all datasets except on GQA. It is because GPT-4V refuses to answer some questions of GQA. We will analyze this later. (2) **Current automatic metrics are not suitable for evaluating the responses generated by GPT-4V.** The automatic evaluation results show that GPT-4V performs worse than Qwen-VL-Chat. However, when we manually evaluate the outputs, we find that GPT-4V's performance is better than Qwen-VL-Chat. We have elaborated on the reasons and show an example of visual question answering in the **Metric** subsection. We further present an example of image captioning in Figure 2. As shown in this figure, the generated descriptions by GPT-4V are more detailed and specific than the ground truths, which makes the generated captions less similar to the ground truths leading to lower metric scores.

Few-shot Results. In-context learning (ICL) as a training-free learning framework enables us to inject human knowledge into models by designing the demonstration and templates (Dong et al., 2022). Given the observation that GPT-4V can discern the order of uploaded images, we devise the

Table 3: Human evaluation for GPT-4V (Zero-shot and Few-shot).

Task	Dataset	Zero-shot	Few-shot
Image Captioning	Flickr30K	19/20	19/20
Visual Question	VQAv2	16/20	17/20
Answering	GQA	11/20	16/20

few-shot prompts to guide GPT-4V to complete the tasks. Considering the maximum number of the uploaded images to GPT-4V is 4, we prompt GPT-4V with 3-shot exemplars. We show an example in Figure 3. The human evaluation results of few-shot prompting are presented in Table 3. GPT-4V with few-shot prompting achieves better performance on the VQA tasks, which indicates GPT-4V has in-context learning ability. Although the exemplars of the image caption task do not help GPT-4V obtain better human evaluation results, they make GPT-4V generate shorter responses leading to an improvement in SPICE from 15.2 to 17.5.

GPT-4V's inconsistent refusal behavior. GPT-4V is instructed to refuse requests for identity, sensitive traits (e.g. age, race), and ungrounded infer-

Visual Question Answering (Few-shot)

Prompt:

Image: The first image

Question: Which type of clothing in this scene is gray?

Answer: Coat

Image: The second image

Question: Which kind of animal is it?

Answer: Dog

Image: The third image

Question: Do you see any people in the living room?

Answer: Yes

Image: The fourth image

Question: Is the tape green?

Answer:

GPT-4V(Few-shot):

Yes, the tape is green.



Figure 3: Few-shot prompting for VQA. We specify the corresponding image to each question in the prompt.

ences. GPT-4V refuses to answer some questions of GQA resulting in low performance on GQA. Specifically, we evaluate 20 test instances sampled from GQA and ask GPT-4V to answer the questions with the zero-shot prompt. 4 out of 20 requests are refused by GPT-4V. These four questions are as follows.

1. Is the player next to the other player female or male?
2. The woman to the right of the camera is watching who?
3. Who wears the shorts?
4. What do you think is the old lady wearing?

The first one is asking for sensitive traits. The second and third questions, to some extent, are asking for information related to identity. The fourth one may be seen as impolite, particularly due to the descriptor "old lady".

To study GPT-4V's refusal behavior, we select some words related to the sensitive traits such as gender, race, and age. We filter out the questions containing these words from the test dataset of

GQA and obtain 20, 11, and 25 test instances related to gender, race, and age respectively. We conduct experiments using these instances. And we observe that 10 out of 20 gender-related requests are refused. 8 out of 10 refused questions are directly asking for gender. Here is an example. "Is the young person female or male?". It is in line with our expectation as GPT-4V is trained to refuse such requests.

But the following examples make us confused, which reveal the inconsistent refusal behavior of GPT-4V.

1. Refused Request: What is the person that is not male standing on?
2. Refused Request: Which side of the photo is the male person on?
3. Approved Request: On which side of the photo is the female pedestrian?
4. Approved Request: Does the male person seem to be sitting?

It seems hard to find out the reasons why does GPT-4V refuse the first and second questions while approve the third and fourth ones.

Chinese Text Recognition (1)

Prompt:

图中是什么汉字?

Translation: What is the Chinese word in the picture?



GPT-4V:

图中的汉字是“愛”，这是繁体中文的“爱”，意思是“love”。

Translation: The Chinese character in the picture is “愛”， which is the traditional Chinese character for “爱”， meaning “love”

Figure 4: GPT-4V fails to recognize the Chinese word in the given image.

Chinese Text Recognition (2)

Prompt:

请复述题干与选项

Translation: Please extract the question stem and options.

3.哺乳动物因长时间未饮水导致机体脱水时，会发生的生理现象是
A.血浆渗透压降低
B.抗利尿激素分泌增加
C.下丘脑渗透压感受器受到的刺激减弱
D.肾小管和集合管对水的重吸收作用减弱

GPT-4V:

Totally Wrong

当然可以。题干与选项如下： 题干： 3.当涉及到物体的长度或者大小时，会发生的单位转换是选项： A. 仍旧涉及长度单位B. 把利润换算为货币单位C. 下午涉及长度单位时需要考虑货币单位D. 再小量单位换算为更大的重量单位

Figure 5: GPT-4V fails to extract the question and choices from the given image.

As for the questions related to race, 4 out of 11 requests are refused.

1. Refused Request: Is the white person next to the windows wearing shorts?
2. Approved Request: Does the white person near the plants seem to be standing?

As for the questions related to age, 4 out of 25 requests are refused.

1. Refused Request: Does the old man appear to be waiting?

2. Approved Request: Are there any old women or men?

The inconsistent refusal behavior is also observed in the early version of GPT-4. They find that GPT-4 tends to become overly cautious in certain ways such as refusing innocuous requests. We consider that this issue is vital for future research and should be systematically studied.

GPT-4V fails to recognize the Chinese text in images. Impressed by the strong English OCR performance of GPT-4V, we wonder whether GPT-4V can recognize the Chinese text in images. We

Table 4: Results on MMLU, HellaSwag, and WinoGrande (Zero-shot).

Dataset	GPT-4V	GPT-4 API
MMLU	16/20	17/20
HellaSwag	14/20	18/20
WinoGrande	15/20	19/20

Table 5: Results on MMLU, HellaSwag, and WinoGrande (Few-shot).

Dataset	GPT-4V	GPT-4 API
MMLU (5-shot)	17/20	18/20
HellaSwag (5-shot)	16/20	16/20
WinoGrande (5-shot)	15/20	17/20

devise the following two tasks: (1) Given an image with only one Chinese word, identify this word; (2) Given an image, extract the question and choices from it. The first task is much easier than the second one. However, GPT-4V fails to complete either the first task or the second task. Specifically, we create 10 instances for each task and show the examples in Figure 4 and Figure 5. We evaluate GPT-4V on these instances, and it achieves 0% accuracy on both tasks, revealing that GPT-4V could not recognize the Chinese text in images.

3 Language Understanding

We evaluate GPT-4V on a wide range of benchmarks to answer two intriguing questions. After being equipped with visual perception, can GPT-4V (1) maintain its language understanding performance and (2) better capture visual commonsense knowledge, world knowledge (specifically physics knowledge)?

As for the first question, we conduct the experiments on MMLU (challenging subjects: `abstract_algebra`, `anatomy`, `astronomy`, `business_ethics`), HellaSwag, and WinoGrande to evaluate the language understanding ability of GPT-4V. Specifically, 20 test instances are sampled for each dataset. Considering that OpenAI may utilize different models to process text-only inputs and text-image inputs, we upload a white image along with the text input. We acknowledge that it is possible that GPT-4V could be affected by the input white image if GPT-4V is not robust enough. We manually obtain and evaluate the results. The results of GPT-4V and GPT-4 (API) are shown in Table 4 and Table 5. We observe that GPT-4V obtains worse results than GPT-4 (API).

Table 6: Results on ViComTe (Zero-shot).

Type	GPT-4V	GPT-4 API
Color	10/10	10/10
Shape	9/10	10/10
Material	10/10	10/10
Size	10/10	10/10
Visual co-occurrence	10/10	10/10

Table 7: Results on UTOPIA (Zero-shot).

Scene	GPT-4V	GPT-4 API
Collision	6/10	9/10
Free fall	3/10	2/10
Friction	10/10	10/10
Incline	10/10	10/10
Motion	10/10	10/10
Projection	7/10	7/10

But the few-shot results indicate that GPT-4V’s performance could be further boosted by more advanced prompting methods.

Let us turn to the second question. We choose ViComTe (Zhang et al., 2022) as our benchmark to find out whether GPT-4V can capture a broad range of visually salient attributes. ViComTe covers 5 property types (color, shape, material, size, and visual co-occurrence) and we sample 10 test instances for each property type to construct our evaluation dataset. We also upload a white image along with the test question. The results are listed in Table 6. The results show that both GPT-4V and GPT-4 (API) perform well on this task. To evaluate GPT-4V’s ability to capture physics knowledge, we utilize UTOPIA (Liu et al., 2023a) as our benchmark. This task requires the models to understand and reason over some basic laws of physics. UTOPIA covers six common scenes including collision, free fall, friction, incline, motion, and projection. We sample 10 test instances for each type and evaluate GPT-4V on them. We also upload a white image along with the input question. The results are presented in Table 7. GPT-4V does not show superiority over GPT-4 (API).

4 Visual Puzzle Solving

GPT-4V has achieved very high performance on the standard visual-centric benchmarks, which indicates these benchmarks could be insufficient to measure its capabilities and limitations. To address this problem, in contrast to constructing more com-

Spot the Differences (1)

Prompt:

Spot ten differences between them



GPT-4V:

Here are the ten differences between the two images:

1. The first image has two birds flying in the sky, while the second image has only one bird.
2. In the first image, the bird on the left is dark-colored, whereas in the second image, the bird is lighter in color.
3. The shadow of the tree on the right side of the first image is longer than in the second image.
4. The elephant's tail in the first image is pointing downwards, while in the second image, it's pointing upwards.
5. There's a small plant near the elephant's foot in the second image, which is not present in the first image.
6. The shape and curve of the tree on the right in the second image are slightly different from the first image.
7. In the first image, the horizon has a more defined mountain range compared to the second image.
8. The cloud formations in the two images are different.
9. The tree trunk on the left side of the first image is straight, whereas in the second image, it has a curve.
10. The tree branches on the left in the first image are slightly more spread out compared to the second image.

Figure 6: GPT-4V finds three differences (marked with red color). GPT-4V hallucinates that the elephant's tail in the second image is pointing upwards.

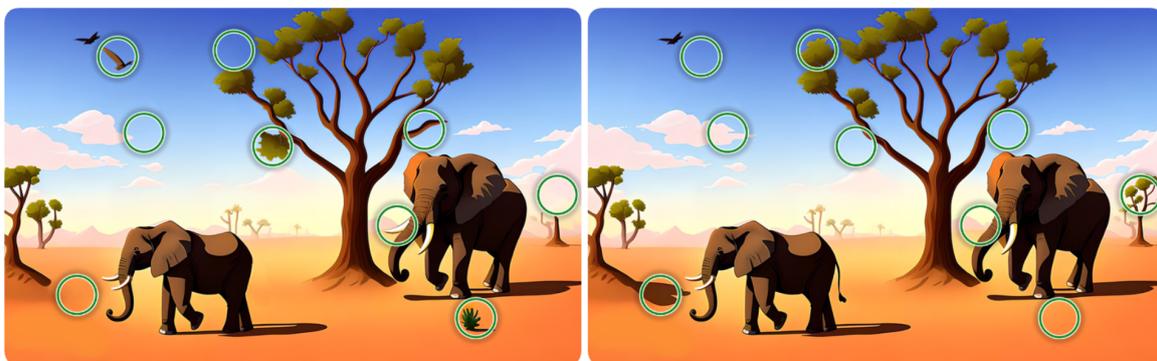


Figure 7: Solution to the level-1 puzzle.

prehensive evaluation benchmarks like MME (Fu et al., 2023) and MMBench (Liu et al., 2023b), we ask GPT-4V to solve interesting but challenging puzzles including Spot the Difference, Draw & Guess, and Math Picture Puzzle.

Spot the Difference. Given two similar images, the goal is to find 10 differences between them. This challenging task requires the model to capture the details of the input images and find out the differences by comparing two similar images. We

manually obtain 10 different level puzzles (from 1 to 10) from CrazyGames³. To evaluate GPT-4V's performance on these puzzles, we devise two types of prompts. As for the first one, we simply position two images within a single image: one on the left and the other on the right.⁴ Then we upload this obtained image to GPT-4V and ask it to spot ten differences. We show the level-1

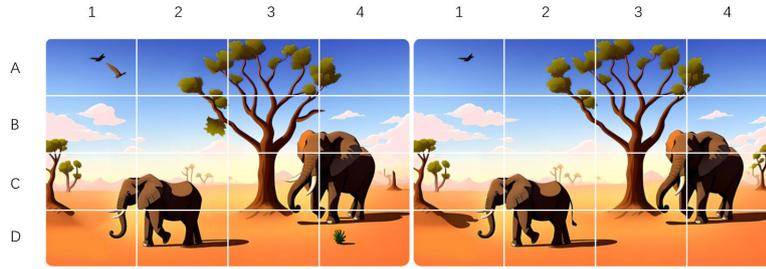
³<https://www.crazygames.com/game/find-the-difference>

⁴We also tried to upload two images separately but did not find any significant difference.

Spot the Differences (2)

Prompt:

Spot ten differences between them and provide the corresponding block numbers



GPT-4V:

Sure, here are ten differences between the two images, along with their corresponding block numbers:

1. A1: In the left image, there are two birds. In the right image, there's only one bird.
2. A3: In the left image, the top branches of the tree have two clusters of leaves. In the right image, there are three clusters of leaves.
3. A4: In the left image, there's one bird. In the right image, there's no bird.
4. B1: In the left image, there are three clouds. In the right image, there are two clouds.
5. C1: In the left image, there's a small tree in the background. In the right image, there's no small tree.
6. C4: In the left image, the elephant's tail is pointing downward. In the right image, the tail is pointing up.
7. D2: In the left image, there's a small green plant near the bottom edge. In the right image, the plant is missing.
8. D3: In the left image, the shadow of the elephant's front leg is longer. In the right image, the shadow is shorter.
9. D4: In the left image, there's no small plant near the elephant's back leg. In the right image, there's a small green plant.
10. B2 and B3: In the left image, the top of the tree where it splits has a more pronounced gap between the branches. In the right image, the gap is narrower.

Figure 8: GPT-4V should not only spot the differences, but also provide the corresponding block numbers. Only one difference is found (marked with red color).

Draw & Guess (1)

Prompt:

Let's play Draw & Guess. You are given a drawing. Please guess a word (Hints: Profession) based on the drawing.



GPT-4V:
Magician!

Figure 9: GPT-4V guesses the word "magician" correctly.

puzzle to GPT-4V and the result is presented in Figure 6. We also present the solution to this puzzle in Figure 7 for reference. As shown in Figure 6, GPT-4V finds three differences but the answers

are rough. Therefore, we design another prompt. We label different areas of the picture with block numbers and ask GPT-4V to spot the differences and provide the corresponding block numbers. In

Draw & Guess (2)

Prompt:

Let's play Draw & Guess. You are given a drawing. Please guess a word (Hints: Activity) based on the drawing.



GPT-4V:

Based on the drawing, my guess is "golf."

Figure 10: GPT-4V guesses the word "golf" correctly.

Draw & Guess (3)

Prompt:

Let's play Draw & Guess. You are given a drawing. Please guess a word (Hints: Object, Machine) based on the drawing.



GPT-4V:

Based on the drawing, my guess is "Washing Machine".

Figure 11: GPT-4V wrongly recognizes the claw machine as the washing machine.

this way, we can evaluate GPT-4V's results more accurately. We show an example in Figure 8. GPT-4V only correctly finds one difference. We test 10 puzzles and manually check the correctness of the answers. GPT-4V finds 14 differences in total with the first prompting method and 8 differences with the second one. The evaluation results show that GPT-4V struggles to capture the small differences between two similar images, which could limit its application. For example, it may be unsuitable that utilizing GPT-4V as a strict discriminator to evaluate the predictions of shadow removal models.

Draw & Guess. Draw & Guess is a casual drawing game. One player selects one word and draws it. The other players should guess the word based on the drawing and some hints such as describing an object. We collect 10 drawings and the corresponding hints. We want to know whether GPT-4V can understand the meaning of each drawing and further guess the corresponding word successfully. We show two success examples in Figure 9 and Figure 10. GPT-4V can capture the visual concepts and guess the words successfully. But GPT-4V also makes mistakes. The only one failure case is presented in Figure 11. GPT-4V fails to find the

Math Picture Puzzle (1)

Prompt:

Solve this question

GPT-4V:

Let's solve the puzzle step by step:

From the first equation:

$$☹ + ☹ + ☹ + ☹ = 32$$

$$3☹ + ☹ = 32$$

$$\Rightarrow 3☹ = 32 - ☹ \quad \dots(i)$$

From the second equation:

$$☹ + 🐰 = 17 \Rightarrow ☹ = 17 - 🐰 \quad \dots(ii)$$

From the third equation:

$$☹ + 🍄 = 10$$

Using equation (ii), we get:

$$17 - 🐰 + 🍄 = 10 \Rightarrow 🍄 - 🐰 = -7 \Rightarrow 🐰 - 🍄 = 7 \quad \dots(iii)$$

Substituting equation (ii) in equation (i):

$$3(17 - 🐰) + 🍄 = 32$$

$$\Rightarrow 51 - 3🐰 + 🍄 = 32 \Rightarrow 🍄 - 3🐰 = -19 \Rightarrow 3🐰 - 🍄 = 19 \quad \dots(iv)$$

Now, subtracting equation (iii) from equation (iv):

$$3🐰 - 🍄 = 19 - 7 \Rightarrow 2🐰 = 12 \Rightarrow 🐰 = 6$$

Using equation (ii):

$$☹ = 17 - 6 \Rightarrow ☹ = 11$$

Using equation (iii):

$$☹ = 6 - 7$$

$☹ = -1$ (This doesn't seem logical since we can't have a negative value for a symbol in such puzzles.)

From the fourth equation:

$$🐰 + ☹ \times 🍄 = ?$$

$$6 + 11 \times (-1) = 6 - 11 = -5$$

However, the value for $☹$ as -1 doesn't make intuitive sense for this type of puzzle. There might be an error in the equations provided or the symbols might represent unconventional values.

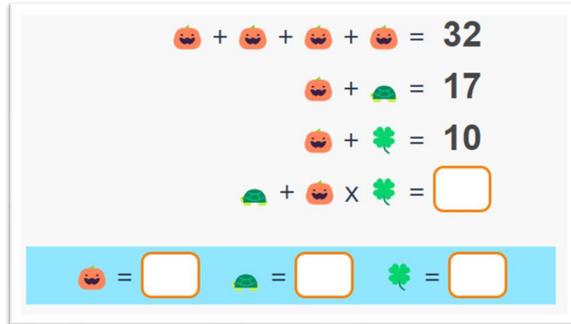


Figure 12: GPT-4V does not recognize the first math equation correctly and fails to solve the puzzle.

clues such as the claw in the image and recognize it as the washing machine. Draw & Guess requires the model to understand the visual concepts (recognizing the hat and the rabbit), recall the related commonsense knowledge (magic), and conduct reasoning to guess the words (magician). It could be possible to construct a visual reasoning benchmark by collecting more challenging instances.

Math Picture Puzzle. Previous work has shown that GPT-4 is good at solving math problems. Inspired by it, we are curious about GPT-4V's performance on the math picture puzzles. Surprisingly, the performance of GPT-4V is really low even though the math picture puzzles are much easier than the problems of the widely used datasets such as GSM-8K (Cobbe et al., 2021). We show an example in Figure 12. GPT-4V wrongly recognizes the first math equation leading to cascading errors and fails to solve the puzzle. We also ask Bard to solve the same puzzle (shown in Figure 13) and Bard fails to recognize all equations. We collect 10 math picture puzzles and test GPT-4V on them. GPT-4V only solves one problem correctly, which is out of our expectations consid-

ering the strong visual understanding ability and math problem-solving ability. We consider the reason could be that GPT-4V does not generalize well to this domain or GPT-4V utilizes some additional modules such as the OCR module to help itself but the additional modules do not work well in this domain.

5 Understanding of Other Modalities

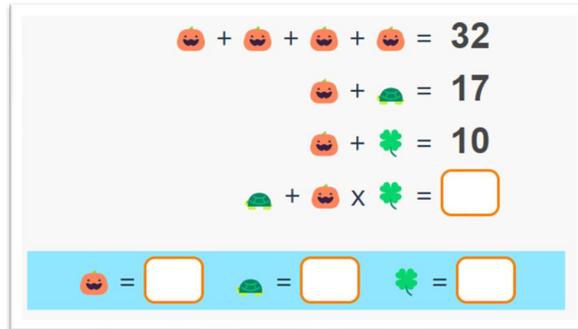
The research question that we investigate in this section is whether can GPT-4V perceive other modalities including depth, thermal, video, and audio. We manually construct test instances for each modality and evaluate GPT-4V's performance.

Depth. Based on the SUN-D dataset (Song et al., 2015), we create 20 test questions of asking for GPT-4V to select the scene label from the five choices for the given depth image. An example is shown in Figure 14. In this case, GPT-4V recognizes the chair through the depth image and correctly predicts the scene label as "living room". We prompt GPT-4V with 20 questions and 4 out of 20 questions are answered correctly, which reveals that GPT-4V struggles to understand the given

Math Picture Puzzle (2)

Prompt:

Solve this question


$$\begin{aligned} \text{Pumpkin} + \text{Pumpkin} + \text{Pumpkin} + \text{Pumpkin} &= 32 \\ \text{Pumpkin} + \text{Watermelon} &= 17 \\ \text{Pumpkin} + \text{Clover} &= 10 \\ \text{Watermelon} + \text{Pumpkin} \times \text{Clover} &= \square \\ \text{Pumpkin} &= \square \quad \text{Watermelon} = \square \quad \text{Clover} = \square \end{aligned}$$

Bard:

Let's solve the puzzle step by step:

To solve the puzzle, we need to find the missing numbers in the equations.

Equation 1: $32 + ? = 17$

Equation 2: $? + 3 = 10$

Equation 3: $? + ? = 15$

We can start by solving Equation 2. We know that $10 - 3 = 7$, so the missing number in Equation 2 is 7.

Now we can use this information to solve Equation 1. We know that $17 - 7 = 10$, so the missing number in Equation 1 is 10.

Finally, we can use the information from Equations 1 and 2 to solve Equation 3. We know that $15 - 10 - 7 = 2$, so the missing numbers in Equation 3 are 10 and 2.

Therefore, the completed puzzle looks like this:

$= + 32 + 10 = 17 + 7 + 3 = 10 + 10 + 2 = 15$

Figure 13: Bard can not recognize the math equations correctly and fails to solve the puzzle.

depth image and further training could be necessary before the application.

Thermal. We sample 20 test instances from LLVIP (Jia et al., 2021) and ask GPT-4V to find the number of pedestrians and their location in the thermal infrared image. We present an example in Figure 15. GPT-4V successfully detects all pedestrians in the image. We test GPT-4V on the sampled 20 instances and 9 out of 20 instances are solved. As for the remaining instances, GPT-4V also can detect correctly several pedestrians from the given images.

Video. It is hard to evaluate accurately and fairly GPT-4V's performance on video understanding considering the maximum number of the uploaded images to GPT-4V is 4 and too much information of the input video is lost. Nevertheless, we are still curious about the GPT-4V's performance on video understanding. We sample 20 test instances from the MSR-VTT dataset (Xu et al., 2016). We select four frames from each video and ask GPT-4V to generate the corresponding video description based on the four frames. We present an example in Figure 16. We find that GPT-4V tends to describe

the images separately and struggles to generate the caption for the whole video. We attribute it to the low sampling rate of the video and increasing the number of the uploaded images could be helpful. We conduct the experiment with 20 instances and 6 out of 20 instances are described correctly.

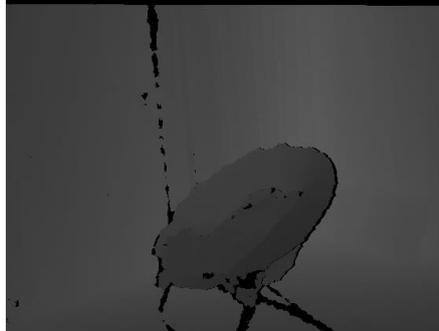
Audio. It is a common way to treat mel spectrograms as images and adopt pre-trained visual encoders to process the mel spectrograms (Wu et al., 2022). Hence, we attempt to evaluate GPT-4V's ability to perceive audio signals by converting the input audio waveforms into the mel spectrograms. Based on the ESC dataset (Piczak, 2015), we devise two tasks to assess GPT-4V's ability: (1) Given a mel spectrum diagram and four choices, the goal is to select an audio label from the four choices; (2) Given four mel spectrograms and four audio labels, the goal is to match the four spectrograms to the four audio labels. We show one example for each task in Figure 17 and Figure 18. As shown in Figure 17, GPT-4V knows the specific patterns of common animal sounds. It could be easy to eliminate "Pig" and "Cow" by the audio duration, but it is harder to decide which one of the other choices is the correct answer. The ground truth is "Dog"

Scene Understanding

Prompt:

This image is the depth image of scene. Please determine the category of this scene.

A. Office B. Kitchen C. Bedroom D. Bathroom E. Living room



GPT-4V:

Based on the depth image provided, it is challenging to definitively determine the category of the scene. However, the image appears to show a **chair-like structure** which could be associated with multiple room categories.

If I had to guess based on the provided options, I would suggest:

E. living room

However, this is a speculative answer as the image does not provide clear details about the entire scene.

Figure 14: GPT-4V can recognize the chair in the depth image and predict the scene label correctly.

Pedestrian Detection

Prompt:

This is an infrared image. Please describe the number of pedestrians and their location in the image for me.



GPT-4V:

In the infrared image, I see two pedestrians.

1. One pedestrian is near the left side of the image, walking on the street beside a tree.
2. Another pedestrian is on the right side of the image, closer to a fence or barrier and appears to be on a marked crosswalk or pavement.

The image is in grayscale and captures a night-time or low-light scenario using infrared technology.

Figure 15: GPT-4V can detect the pedestrians in the given infrared image correctly.

Video Captioning

Prompt:

Here are some pictures from a video. Could you explain what the video is describing?

GPT-4V:

These images depict fantastical cityscapes, likely from a movie or a series. The intricate architecture, the flowing waterfalls, and the towering structures give it a majestic and otherworldly vibe. The silhouette of a person in the second image adds to the sense of scale and wonder.

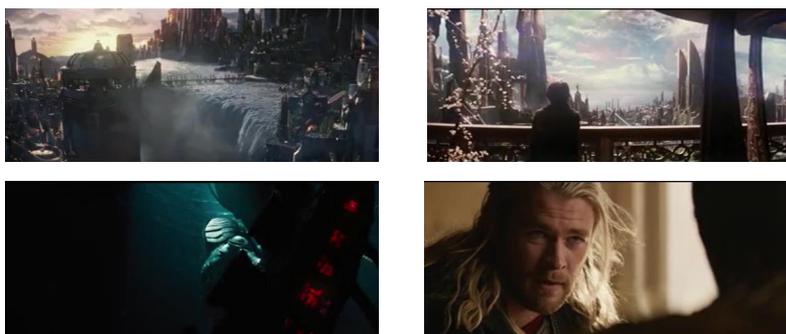


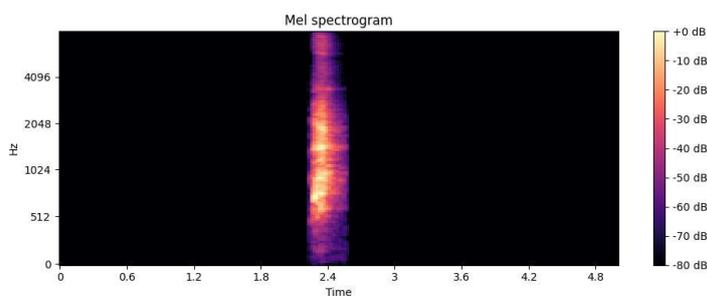
Figure 16: GPT-4V describes the image content well but struggles to generate the caption for the whole video. Increasing the number of the sampled frames could be helpful.

Audio Classification (1)

Prompt:

This image is the mel spectrum diagram of a piece of audio. Please determine the source of this audio.

A. Dog B. Rooster C. Pig D. Cow.



GPT-4V:

The image is a Mel spectrogram, which represents the distribution of frequencies in an audio signal over time. By looking at the image alone, I cannot definitively identify the source of the audio. However, based on common audio patterns: A rooster typically has a sharp and distinctive call that might show a pronounced and relatively short pattern in a spectrogram. Dogs can have a variety of barks, howls, and whines, with different frequency distributions and lengths. Pigs often produce grunts and squeals, which can be longer in duration and might show a more spread-out pattern in the spectrogram. Without listening to the actual audio, I can't provide a definitive answer. If you're looking for a guess based on the spectrogram's visual representation and the choices provided, I would lean towards B. Rooster due to the sharp and pronounced pattern in the spectrogram. However, this is just an educated guess.

Figure 17: The ground truth is “Dog” while GPT-4V wrongly recognize it as “Rooster”.

while GPT-4V wrongly recognizes it as “Rooster”. We construct 20 test instances covering 5 major

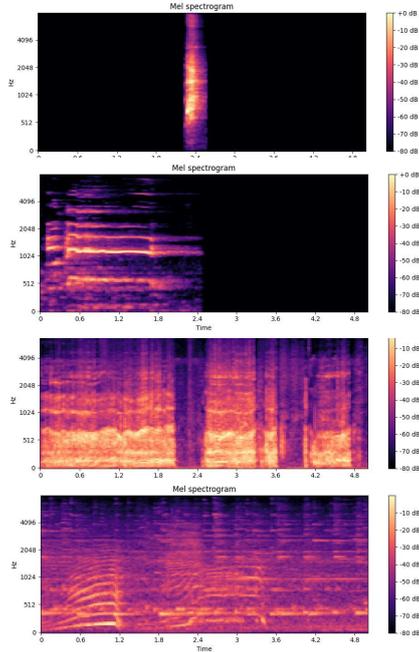
categories including Animals, Natural soundscapes & water sounds, Human(non-speech sounds), In-

Audio Classification (2)

Prompt:

These four images are mel spectrograms of four audio clips. Match the four sounds to the four options below

a: Pig b. Cow c. Dog d. Rooster



GPT-4V:

Matching audio clips to their visual representations as mel spectrograms can be challenging, especially without hearing the actual sound. However, based on typical patterns, I'll make an educated guess: 1. The first spectrogram shows multiple harmonic lines, which might indicate some sort of vocalization with a repetitive pattern. This could potentially represent the crowing of **d. Rooster**. 2. The second spectrogram is more consistent with fewer variations, which might indicate a constant sound or background noise. This might represent the grunting or ambient sounds of **a. Pig**. 3. The third spectrogram has a more irregular pattern, possibly indicating barking or other intermittent sounds. This might correspond to **c. Dog**. 4. The fourth spectrogram has patterns that suggest a mixture of consistent sounds with some harmonic overtones. This could be indicative of the mooing of **b. Cow**. Please note that these are speculative matches based on the visual representation and might not be accurate. Listening to the actual audio clips would provide a more accurate match.

Figure 18: GPT-4V fails to match the given four audios to their labels. The ground truth is “1. Dog 2. Rooster 3. Pig 4. Cow”

terior/domestic sounds, and Exterior/urban noises based on ESC. The result is that GPT-4V successfully recognizes 5 out of 20 instances, which is the same as the random selecting method. As for the second task, GPT-4V successfully matches 2 out of audios to their labels. We show an example for the second task in Figure 18. GPT-4V fails to match the given four audios to their labels, which indicates that although GPT-4V knows some common patterns of sounds, it is still challenging for GPT-4V to recognize the audio labels directly from the mel spectrograms.

6 Conclusion

In this paper, we quantitatively study GPT-4V’s performance on various tasks. According to the results, we find that although GPT-4V achieves high performance on standard English visual-centric benchmarks, it still can not perform Chinese text recognition. This observation suggests further in-depth evaluation on Chinese benchmarks is necessary for measure GPT-4V’s capability. We also observe that GPT-4V fails to solve easy math picture puzzles

even though it has strong visual understanding ability and math problem solving ability. The reason could be that GPT-4V does not generalize well to this domain. Another problem is that GPT-4V exhibits inconsistent refusal behavior when answering questions related to identity and sensitive traits such as gender, race, and age. This issue could lead to an obvious performance drop of GPT-4V and should be dealt with carefully in future studies.

As for the limitations, we acknowledge that GPT-4V’s performance could be different by adopting different prompting methods. For example, more specific instructions and better exemplars will improve its performance. We would like to explore utilizing other advanced prompts such as chain-of-thought prompting (Wei et al., 2022) in future work. We also acknowledge that more test instances for each task can make the estimated results more accurate, but we only sample a part of instances due to the high labor cost.

Nevertheless, it is the first attempt to quantitatively study GPT-4V’s performance on a wide range of tasks. In our study, we reveal the strengths

and limitations of GPT-4V. We hope our study can provide insights into future research and application.

References

- Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. 2019. No-caps: Novel object captioning at scale. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8948–8957.
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. Spice: Semantic propositional image caption evaluation. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part V 14*, pages 382–398. Springer.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Rongrong Ji. 2023. **Mme: A comprehensive evaluation benchmark for multimodal large language models**.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2016. **Making the v in vqa matter: Elevating the role of image understanding in visual question answering**. *International Journal of Computer Vision*, 127:398 – 414.
- Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3608–3617.
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709.
- Xinyu Jia, Chuang Zhu, Minzhen Li, Wenqi Tang, and Wenli Zhou. 2021. Llvip: A visible-infrared paired dataset for low-light vision. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3496–3504.
- Ruibo Liu, Jason Wei, Shixiang Shane Gu, Te-Yen Wu, Soroush Vosoughi, Claire Cui, Denny Zhou, and Andrew M. Dai. 2023a. **Mind’s eye: Grounded language model reasoning through simulation**. In *The Eleventh International Conference on Learning Representations*.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. 2023b. **Mmbench: Is your multi-modal model an all-around player?**
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204.
- OpenAI. 2023. **Gpt-4 technical report**.
- Karol J. Piczak. 2015. **Esc: Dataset for environmental sound classification**. In *Proceedings of the 23rd ACM International Conference on Multimedia, MM ’15*, page 1015–1018, New York, NY, USA. Association for Computing Machinery.
- Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. 2015. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Yang Wu, Zhenyu Zhang, Pai Peng, Yanyan Zhao, and Bing Qin. 2022. **Leveraging multi-modal interactions among the intermediate representations of deep transformers for emotion recognition**. In *Proceedings of the 3rd International on Multimodal Sentiment Analysis Workshop and Challenge, MuSe’ 22*, page 101–109, New York, NY, USA. Association for Computing Machinery.

- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The dawn of lmms: Preliminary explorations with gpt-4v(ision).
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. 2014. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78.
- Chenyu Zhang, Benjamin Van Durme, Zhuowan Li, and Elias Stengel-Eskin. 2022. Visual commonsense in pretrained unimodal and multimodal models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5321–5335.